

University of Southampton Research Repository ePrints Soton

Copyright © and Moral Rights for this thesis are retained by the author and/or other copyright owners. A copy can be downloaded for personal non-commercial research or study, without prior permission or charge. This thesis cannot be reproduced or quoted extensively from without first obtaining permission in writing from the copyright holder/s. The content must not be changed in any way or sold commercially in any format or medium without the formal permission of the copyright holders.

When referring to this work, full bibliographic details including the author, title, awarding institution and date of the thesis must be given e.g.

AUTHOR (year of submission) "Full thesis title", University of Southampton, name of the University School or Department, PhD Thesis, pagination

UNIVERSITY OF SOUTHAMPTON

FACULTY OF HEALTH SCIENCES

**SYMPTOMATIC DIAGNOSIS OF LUNG CANCER in a population referred to
lung-shadow clinic with high rates of chronic respiratory diseases: A
feasibility study**

by

Joanna Shim

Thesis for the degree of Doctor of Philosophy

May 2015

UNIVERSITY OF SOUTHAMPTON

ABSTRACT

FACULTY OF HEALTH SCIENCES

Thesis for the degree of Doctor of Philosophy

SYMPTOMATIC DIAGNOSIS OF LUNG CANCER

Joanna Shim

In the UK, 86% of lung cancer (LC) patients are diagnosed when curative treatment is not possible. Government guidelines recommend urgent chest X-ray referrals for patients presented with any 1 of 10 suggested LC symptoms. Little evidence currently supports these recommendations. Thus, the need for prospective studies to identify the predictive values of symptoms for LC diagnosis.

This study aimed to investigate the feasibility of a prospective study to identify symptoms that predict LC in a secondary care population with high rates of chronic respiratory disease by investigating 1) the content validity and acceptability to this population of a patient-completed questionnaire, and 2) identifying patient-elicited symptoms that predict LC.

The Identifying Symptoms that Predict Chest and Respiratory Disease (IPCARD) questionnaire was used to prospectively collect symptom, risk and comorbidity data in a chest clinic population investigated for LC (Patients aged ≥ 40). LC was identified six months following recruitment. Semi-structured and cognitive interviews explored the acceptability and content validity of the IPCARD questionnaire in this population. Multiple logistic regression was used to identify symptoms independently associated with LC in the clinic population, and in a COPD sub-group; at two levels of entry criteria ($p < 0.05$ and $p < 0.15$).

359 patients (70% of those eligible) completed the IPCARD questionnaire, and 77 (21.4%) were diagnosed with LC. Qualitative research indicated that participants found the IPCARD questionnaire acceptable, and its content validity was established in this secondary care population. Cough and breathing changes first indicated in the last three months, predicted LC in this referred population ($p < 0.05$). In the COPD sub-group, the symptom descriptor, unable to get enough air predicted LC ($p < 0.05$). At the relaxed criteria ($p < 0.15$), five symptoms predicted LC in the full clinic population; a hard/harsh cough without phlegm, increasing chest infections, and weight gain (last 12 months) were added to the previous model. The COPD sub-group at $p < 0.15$, breathing changes experienced (last three months), breathing changes first indicated within the last three months, unable to get enough air, and wheezing sensation, were predictors. Based on the more rigorous entry criteria ($p < 0.05$), the diagnostic criteria for the COPD sub-group (positive likelihood ratio (+LR)=1.91; Area under curve (AUC)=0.739) performed slightly better than the criteria for the full population (+LR=1.49; AUC=0.729) (at optimal cut-off).

The better performance of the COPD-specific model supports the need for an adequately powered study to investigate the predictive values of LC symptoms in homogeneous populations with specific respiratory diseases.

Contents

ABSTRACT	iii
Contents	v
List of tables	xiii
List of figures	xix
DECLARATION OF AUTHORSHIP	xxi
Acknowledgements.....	xxiii
Abbreviations	xxv
Chapter 1: Overview.....	1
1.1 Thesis layout and contents	1
1.2 Introduction.....	2
1.3 Aims	8
1.4 Research Questions	9
Chapter 2: Literature Review	11
LUNG CANCER	11
2.1 Pathophysiology	11
2.2 Epidemiology	12
2.3 Aetiology and Risk factors	15
2.3.1 Smoking.....	15
2.3.2 Occupational exposures	16
2.3.3 Socio-economic factors.....	17
2.3.4 Genetic and other factors	17
2.4 Signs and Symptoms.....	18
2.5 Early detection and screening for Lung Cancer.....	20
2.6 Evaluation of diagnostic tests	23
2.7 Lung Cancer staging in relations to treatment and prognosis....	24
CHRONIC OBSTRUCTIVE PULMONARY DISEASE (COPD)	26
2.8 Pathophysiology	26
2.8.1 COPD in relations to lung cancer	27
2.9 Epidemiology	30
2.10 Aetiology and risk factors	30

2.11	Signs and symptoms	31
2.12	Diagnosis.....	32
2.12.1	Use of symptoms-based questionnaire tools	32
2.13	Severity of COPD	33
2.14	Early detection and screening for Lung Cancer in COPD	34
	SYMPTOMATIC DIAGNOSIS OF LUNG CANCER.....	36
Chapter 3: Systematic Review		41
	Symptomatic diagnosis of Lung Cancer.....	41
3.1	Background.....	41
3.2	Aims	43
3.3	Methods.....	43
3.3.1	Search strategy.....	43
3.3.2	Inclusion and exclusion criteria	44
3.3.3	Study selection and quality assessment	45
3.3.4	Data extraction and analysis.....	46
3.4	Results.....	48
3.5	Results of Quantitative studies	49
3.5.1	Symptomatic prevalence	50
3.5.2	Symptomatic diagnosis.....	50
3.5.3	Key methodological strengths and limitations of the studies.....	61
3.6	Results of Qualitative studies	66
3.6.1	Time intervals between symptom onset and diagnosis.....	66
3.6.2	Pre-diagnostic bodily experiences (symptoms and health changes)	67
3.6.3	Diagnostic delays	68
3.7	Discussion	74
3.8	Conclusion.....	76
Chapter 4: Methodology		77
4.1	Research Design	77
4.2	Research Plan.....	78
4.2.1	Missing data.....	80
4.2.2	Quantitative analysis	80
4.3	Setting	82
4.3.1	Organisational structure of the clinic	82
4.4	Study Population	83
4.5	Inclusion Criteria.....	83
4.6	Exclusion Criteria.....	84

4.7	Recruitment methods	84
4.8	Development of the IPCARD Questionnaire	86
4.9	Identifying diagnoses (lung cancer and COPD)	87
4.10	Data entry and data cleaning	89
4.11	Ethics	89
4.11.1	Data protection and Confidentiality	90
Chapter 5: Study 1		91
5.1	Introduction.....	91
5.2	Methods	92
5.2.1	Sample size for qualitative research	92
5.2.2	Qualitative data collection (Interviews).....	92
5.2.3	Qualitative data analysis.....	93
5.3	Results	94
5.3.1	Recruitment to the Qualitative research	94
5.3.2	Characteristics of people who declined to be interviewed	96
5.3.3	Results of Qualitative Analysis	96
5.3.3.1	Acceptability of the questionnaire to a population of lung-shadow clinic attendees.....	97
5.3.3.2	Inconsistencies between symptom experiences reported in the interviews post-diagnosis and those recorded in the questionnaire pre-diagnosis.....	98
5.3.3.3	Content validity of the questionnaire in a population consisting of high rates of chronic respiratory diseases referred to lung-shadow clinic.....	108
5.4	Discussion and Conclusion of Study 1	113
Chapter 6: Study 2.....		117
6.1	Introduction.....	117
6.2	Methodology Section	118
	Missing data.....	118
6.2.1	Introduction	118
6.2.2	Rationale for missing data methodology.....	119
6.2.3	Multiple imputation (MI)	122
6.2.4	Identifying symptoms that predict lung cancer diagnosis	125
6.2.5	Method of variable selection.....	126
6.2.6	Variable selection with an imputed dataset.....	128
6.2.7	Assessing fit of the model	128
6.3	Method Section.....	130
	Missing data.....	130
6.3.1	Description of missing data.....	130

6.3.1.1	Missing data rates	130
6.3.1.2	Missing data patterns	131
6.3.2	Consistency check with MAR assumptions	132
6.3.3	Imputation model building	133
6.3.4	Imputation model and variable type.....	135
6.3.5	Number of imputations.....	136
6.3.6	Convergence	137
6.3.7	Diagnostics	137
6.3.8	Sample size calculation.....	138
6.3.9	Recruitment strategy	139
6.3.10	Quantitative data collection	139
6.3.10.1	Independent symptom variables.....	140
6.3.10.2	Independent socio-demographic variables	142
6.3.10.3	Independent epidemiological risk variables: clinical and behavioural risk factors.....	143
6.3.10.4	Comorbidities	144
6.3.11	Data entry and data cleaning	145
6.3.12	Quantitative data analysis.....	146
6.3.12.1	Univariate analysis	146
6.3.12.2	Bivariate analysis	147
6.3.12.3	Multivariate analysis	147
6.3.12.4	Variable selection procedure	147
6.3.12.5	Developing a set of diagnostic criteria for LC population	149
6.3.12.6	Weighted diagnostic criteria	151
6.3.12.7	Secondary sub-group analysis	151
6.4	Results Section.....	152
	Missing data	152
6.4.1	Frequencies of missing data	152
6.4.2	Distribution of missing data in symptom variables by socio- demographic variables	154
6.4.2.1	Gender	154
6.4.2.2	Age	156
6.4.3	Distribution of missing data in symptom variables by clinical outcome variables and clinical covariates	158
6.4.3.1	Cancer (Outcome variable)	158
6.4.3.2	COPD	160
6.4.3.3	Epidemiological Risk variable	162
6.4.4	Distribution of missing data in symptom variables by clinically relevant symptom covariates.....	162
6.4.5	Issues of collinearity.....	167

6.4.6	Result of convergence	167
6.4.7	Pattern of missing data	168
6.4.8	Diagnostics for imputations	168
6.4.9	Discussion and conclusion to missing data.....	173
	Main analysis.....	174
6.4.10	Recruitment strategy.....	174
6.4.11	Participant Descriptive Data.....	174
6.4.12	Dichotomising multiple-response variables.....	178
	Full population	180
6.4.13	Univariate analysis of the relationship between symptoms and lung cancer	180
6.4.14	Bivariate analyses (Comorbidities)	186
	Asthma.....	187
	COPD.....	187
	Arthritis.....	189
	Pneumonia	190
6.4.15	Multivariate analysis: Full population data	190
	Variable selection at $p < 0.05$	190
	6.4.15.1 Model 1; Symptoms, adjusted for age ($p < 0.05$)	190
	6.4.15.2 Model 2; Symptoms, adjusted for age, and risk variables ($p < 0.05$).....	191
6.4.16	Developing a set of diagnostic criteria ($p < 0.05$).....	192
6.4.17	Weighted set of diagnostic criteria	195
	Variable selection at $p < 0.15$	196
	6.4.17.1 Model 1; Symptoms, adjusted for age ($p < 0.15$)	196
	6.4.17.2 Model 2; Symptoms; adjusted for age and risk variables ($p < 0.15$).....	197
6.4.18	Developing a set of diagnostic criteria ($p < 0.15$).....	198
6.4.19	Weighted set of diagnostic criteria	201
	COPD.....	201
6.4.20	Sub-group analysis: COPD population.....	201
6.4.21	Univariate analyses: COPD sub-population.....	202
6.4.22	Bivariate analysis: COPD sub-population.....	207
	Asthma.....	207
	Angina.....	208
	Arthritis.....	209
	Allergy.....	210
6.4.23	Multivariate analysis: COPD sub-population.....	211
	Variable selection at $p < 0.05$	211

6.4.23.1	Model 1; Symptoms, adjusted for age ($p < 0.05$): COPD sub-population.....	211
6.4.23.2	Model 2; Symptoms, adjusted for age and risk variables ($p < 0.05$): COPD sub-population.....	212
6.4.24	A set of diagnostic criteria for a sub-population with COPD ($p < 0.05$).....	212
6.4.25	Weighted set of diagnostic criteria for a sub-population with COPD ($p < 0.05$).....	214
	Variable selection at $p < 0.15$	215
6.4.25.1	Model 1; Symptoms, adjusted for age ($p < 0.15$): COPD sub-population.....	215
6.4.25.2	Model 2; Symptoms, adjusted for age and risk variables ($p < 0.15$): COPD sub-population.....	216
6.4.26	A set of diagnostic criteria for a sub-population with COPD ($p < 0.15$).....	216
6.4.27	Weighted set of diagnostic criteria for a sub-population with COPD ($p < 0.15$).....	219
6.4.28	Results of diagnostic models: complete case vs. imputed	223
	Variable selection at $p < 0.05$	228
6.4.28.1	Model 1; Symptoms, adjusted for age ($p < 0.05$): Complete case analysis on full population data	228
6.4.28.2	Model 2; Symptoms, adjusted for age and risk variables ($p < 0.05$): Complete case analysis on full population data	228
6.4.28.3	Model 1; Symptoms, adjusted for age ($p < 0.15$): Complete case analysis on full population data	229
6.4.28.4	Model 2; Symptoms, adjusted for age and risk variables ($p < 0.15$): Complete case analysis on full population data	230
6.4.28.5	Comparison between complete case model and imputed model	231
6.5	Discussion	234
6.5.1	Missing data.....	234
6.5.2	Discussion of main findings.....	235
	Use of a relaxed criteria for variable selection ($p < 0.15$): full population	236
6.5.3	Sub-population COPD analysis	237
	Use of a relaxed criteria for variable selection ($p < 0.15$): COPD	238
6.5.4	Development of a set of diagnostic criteria.....	238
6.6	Conclusion.....	240
Chapter 7:	Discussion and Conclusion.....	241
7.1	Introduction	241
7.2	Study findings	241
7.2.1	Diagnostic symptoms in COPD population (sub-group)	244

7.3	Development of a set of diagnostic criteria	245
7.4	Methodological issues	245
7.5	Strengths and limitations of study	246
7.6	Implications for clinical practice and future research.....	249
7.7	Conclusion	249
Appendices		251
Appendix 1 Revised TNM staging system (2010) (American Joint Committee on Cancer (AJCC) 2002)		253
Appendix 2 Published paper.....		257
Appendix 3 Results of database search strategy (Supplementary Data)		269
Appendix 4 MRC Respiratory questionnaire		275
Appendix 5 IPCARD Questionnaire.....		279
Appendix 6 Introductory Letter.....		293
Appendix 7 Participant Information Sheet (A)		295
Appendix 8 Participant Information Sheet (B).....		297
Appendix 9 Schema for recruitment in Southampton clinic:		299
Appendix 10 Interview Schedule Outline.....		301
Appendix 11 Frequency and percentage of missing observations for all variables		305
Appendix 12 Imputation Model		309
Appendix 13 Collinearity diagnostics.....		311
Appendix 14 Missing data pattern		313
Appendix 15 Tetrachoric correlations to determine response cut-offs.....		343
List of References		351

List of tables

Table 2.1 Lung cancer statistics- number of deaths and crude rate per 100,000 of the UK population in 2012 (CRUK 2014).....	13
Table 2.2 Classification of COPD severity defined (GOLD 2001)	34
Table 3.1 List of electronic databases searched	44
Table 3.2 Summary of quantitative studies	53
Table 3.3 Study methodology- study design and characteristics of exposure data (symptom)	55
Table 3.4 Symptoms reported that were independently associated with lung cancer	57
Table 3.5 Symptoms measured in the quantitative studies using either predetermined list or questionnaire unless stated otherwise	58
Table 3.6 Diagnostic values (ORs, PPVs, and HRs) for symptoms reported in case-control studies and cohort studies.....	59
Table 3.7 Summary of qualitative studies	72
Table 5.1 Characteristics of interview participants	96
Table 6.1 Types of response categories for questionnaire items	140
Table 6.2 Risk variables.....	144
Table 6.3 Baseline comorbidities	145
Table 6.4 Investigation of non-differential classification of the proportion of missingness in each generic variable by gender	155
Table 6.5 Investigation of non-differential classification of the proportion of missingness in variables with > 10% missing data by gender	156
Table 6.6 Linear regression of generic symptom variables as explanatory/ predictor variables and age the dependent variable (continuous).....	157

Table 6.7 Linear regression of variables with 10% or more missing data as explanatory/ predictor variables and age the dependent variable (continuous)	158
Table 6.8 Investigation of non-differential classification of the proportion of missingness in each generic variable by lung cancer diagnosis	159
Table 6.9 Investigation of non-differential classification of the proportion of missingness in variables with missing data >10% by lung cancer diagnosis..	160
Table 6.10 Investigation of non-differential classification of the proportion of missingness in each generic variable by COPD diagnosis	161
Table 6.11 Investigation of non-differential classification of the proportion of missingness in variables with missing data >10% by COPD diagnosis.....	162
Table 6.12 Covariates that predicted the missingness in the variables with higher missing data > 10% at a statistically significant level ($p < 0.05$).....	164
Table 6.13 Covariates that predicted the missingness in generic symptom variables at a statistically significant level ($p < 0.05$).....	165
Table 6.14 Covariates with ORs > 2 but not statistically significant ($p > 0.05$) (variables with higher missing data > 10%).....	165
Table 6.15 Covariates with ORs > 2 but not statistically significant ($p > 0.05$) (generic symptom variables)	166
Table 6.16 Means, medians, and standard deviations of the complete case and imputed for severity variables	169
Table 6.17 Frequency distribution of ordinal variables (cut-off at 6).....	169
Table 6.18 Frequency distribution for binary variables (only 'yes' or '1' response presented)	170
Table 6.19 Frequency distribution for three-response categorical variables..	172
Table 6.20 Frequency distribution for four- response categorical variables ..	173
Table 6.21 Demographic details of participants (n=359)	176

Table 6.22 Clinical characteristics of participants (n=359); follow up of participant's diagnoses	177
Table 6.23 Univariate analysis of the relationship between symptoms and lung cancer ($p < 0.05$)	182
Table 6.24 Univariate analysis of the relationship between symptoms and lung cancer ($OR > 2.0$ or < 0.5 , or $p < 0.15$)	183
Table 6.25 Univariate associations between risk factors and lung cancer diagnosis	185
Table 6.26 Patient-reported comorbidities (full population data).....	186
Table 6.27 M-H χ^2 test for asthma.....	187
Table 6.28 M-H χ^2 test for COPD	189
Table 6.29 M-H χ^2 test for arthritis.....	190
Table 6.30 Model 1: Main effects model ($p < 0.05$), adjusted for age; using forward stepwise regression at $p(e) = 0.05$ and $p(r) = 0.10$	191
Table 6.31 Model 2: Main effects model with risk variables ($p < 0.05$); using forward stepwise regression; $p(e) = 0.05$ and $p(r) = 0.10$	191
Table 6.32 Analysis of the diagnostic performance of this set of criteria ($p < 0.05$)	193
Table 6.33 Model 1: Main effects model ($p < 0.15$), adjusted for age; using forward stepwise regression; $p(e) = 0.10$ and $p(r) = 0.15$	196
Table 6.34 Model 2: Main effects model ($p < 0.15$) with risk variables; using forward stepwise regression at $p(e) = 0.10$ and $p(r) = 0.15$	198
Table 6.35 Analysis of the diagnostic performance of this set of criteria ($p < 0.15$)	199
Table 6.36 Univariate analysis of the relationship between symptoms and lung cancer in a population with COPD ($OR > 2.0$, or < 0.5 , or $p < 0.05$).....	203

Table 6.37 Univariate analysis of the relationship between symptoms and lung cancer in a population with COPD (OR>2.0 or <0.5, or p<0.15).....	204
Table 6.38 Univariate analysis of the relationship between risk variables and lung cancer in a population with COPD	206
Table 6.39 Patient-reported comorbidities (sub-population data).	207
Table 6.40 M-H χ^2 test for asthma	208
Table 6.41 M-H χ^2 test for angina	209
Table 6.42 M-H χ^2 test for arthritis	210
Table 6.43 M-H χ^2 test for allergy	210
Table 6.44 Model 1: Main effects model (p<0.05), adjusted for age; using forward stepwise regression at p(e)=0.05 and p(r)=0.10	211
Table 6.45 Analysis of the diagnostic performance of this set of criteria.....	212
Table 6.46 Log odds ratios and derived weights from variables identified in the COPD model	214
Table 6.47 Model 1: Main effects model (p<0.15), adjusted for age; using forward stepwise regression at p(e)=0.10 and p(r)=0.15	215
Table 6.48 Analysis of the performance of the un-weighted set of diagnostic criteria for sub-group with COPD	217
Table 6.49 Log odds ratios and derived weights from variables identified in the model.....	219
Table 6.50 Analysis of the diagnostic performance of the weighted set of criteria in the sub-group analysis	220
Table 6.51 A table summary of the performance of the criteria for each version of the model.....	223
Table 6.52 Univariate analyses of symptoms and lung cancer as outcome for complete case and imputed data.	224

Table 6.53 Model 1 including variables at $p < 0.05$: Complete case	228
Table 6.54 Model 2 including variables at $p < 0.05$: Complete case	228
Table 6.55 Model 1 including variables $p < 0.15$: Complete case analysis	230
Table 6.56 Model 2 including variables $p < 0.15$: Complete case analysis	230
Table 6.57 Comparison of Model 1 ($p < 0.15$) between the complete case and imputed data.....	232
Table 6.58 Descriptive characteristics of the three miss_variables; Q21_BrChnges, Q46_ChInfectn, and Q52_Weight, by socio-demographic variable (age), outcome variable (lung cancer diagnosis), clinical covariates (COPD, smoking, family history)	233

List of figures

Figure 1.1 Schematic plan of the PhD in relations to the larger IPCARD study ..	7
Figure 2.1 Average number of new cases per year and age-specific incidence rates between 2009 and 2011 (CRUK 2014; ONS 2012).....	14
Figure 2.2 Trends in lung cancer age-standardised incidence rates, and smoking prevalence in Britain from 1948 to 2011 (CRUK 2014).....	16
Figure 2.3 Progression of cancer (classical model) (Zheng et al. 2011).....	21
Figure 2.4 Illustration of the physiology behind inflammation in COPD and lung cancer (de Torres et al. 2007).....	28
Figure 2.5 Mechanism of lung carcinogenesis following COPD (illustrations by Sekine et al. 2012)	29
Figure 2.6“Could it be COPD” questionnaire (GOLD 2004).....	33
Figure 3.1 Flow diagram of results of search strategy adopted from the QUOROM statement flow diagram (Moher et al. 1999)	49
Figure 4.1 Illustration of a confounder	81
Figure 4.2 Operational structure of lung clinic in SUHT.....	83
Figure 4.3 Participant recruitment process	85
Figure 6.1: Illustration of multiple imputation process (diagram modified from Humphries 2012)	124
Figure 6.2: Illustration of monotonic and arbitrary (non-monotonic) patterns of missing data for ‘k’ number of variables (adopted from Husmain 2008)	132
Figure 6.3 Example of questionnaire data with naturally occurring non-response	146
Figure 6.4 Example of generic symptom variable in questionnaire.....	153

Figure 6.5 Example of questionnaire structure for multiple response questions	178
Figure 6.6: Tetrachoric correlation for cough symptom, Q13a_Cgh.....	179
Figure 6.7: Tetrachoric correlation for breathing changes symptom, Q22a_BrChnges	180
Figure 6.8: Tetrachoric correlation for tiredness symptom, Q32_Tired	180
Figure 6.9 Receiver operating characteristic curve for the un-weighted set of criteria (p<0.05)	194
Figure 6.10: Receiver operating characteristic curve for the un-weighted set of criteria.....	200
Figure 6.11 Receiver operating curve for the un-weighted set of diagnostic criteria for sub-group with COPD (p<0.05)	213
Figure 6.12 Receiver operating curve for the unweighted set of diagnostic criteria for sub-group with COPD (p<0.15)	218
Figure 6.13: Receiver operating curve of the weighted set of criteria for the COPD sub-group.....	221
Figure 6.14 Receiver operating curves of the weighted and un-weighted set of criteria for the COPD sub-group.....	222

DECLARATION OF AUTHORSHIP

I, **Joanna Shim**, declare that the thesis entitled “**Symptomatic Diagnosis of Lung Cancer**” and the work presented in the thesis are both my own, and have been generated by me as the result of my own original research. I confirm that:

- this work was done wholly or mainly while in candidature for a research degree at this University;
- where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated;
- where I have consulted the published work of others, this is always clearly attributed;
- where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work;
- I have acknowledged all main sources of help;
- where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself;
- none of this work has been published before submission, or [delete as appropriate] parts of this work have been published as: [please list references]

Signed:

Date:.....

Acknowledgements

First and foremost, I would like to thank my participants, for their patience and generosity; taking the time to fill out my questionnaire in support of a better future through research.

Specific thanks go to the following for their support and professional guidance throughout my PhD:

- To the clinical staff members in Southampton General Hospital (medical out-patients unit) for your assistance in the recruitment process of my study. I couldn't ask for a more caring and supportive team. From the start, you've welcomed me into the team as your own.
- To Dr Michael Simon, for his invaluable guidance and knowledge in the systematic review process.
- To Dr Sean Ewings, for his unwavering patience in guiding me through my statistical journey.
- To my sponsors, the University of Southampton and the Brunei government for funding me.

A personal debt of gratitude goes to my friends and family for their love and support:

- To my PhD family, thank you for the daily encouragement and generosity you've shown me; be it time, resources, or cakes.
- To my amazing family, thank you for letting me pursue this, and supporting me unconditionally.
- To Jason, for believing in me.

Last, but by no means least, my sincerest gratitude goes to my supervisors:

- To Associate Professor Lucy Brindle, firstly for this opportunity- I only ever wanted to make you proud. I am very grateful for your patience, wisdom, and continuous academic support throughout this process. You have set an excellent example as a mentor, a researcher, and a role model.
- To Dr Steve George for his encouragement, and positivity in times of doubt. I have learnt so much from the both of you, and thank you for the countless opportunities you have given me. For that, I am eternally grateful.

Abbreviations

AIC:	Akaike's Information Criterion
AUC:	Area Under Curve
COPD:	Chronic Obstructive Pulmonary Disease
CT:	Computed Tomography
GOLD:	Global Initiative for Chronic Obstructive Lung Disease
HR:	Hazard Ratio
IPCARD:	Identifying Symptoms that Predict Chest and Respiratory Disease
LC:	Lung Cancer
LR:	Likelihood Ratio
MAR:	Missing at Random
MCAR:	Missing Completely at Random
MNAR:	Missing Not At Random
NICE:	National Institute of Clinical Excellence
NLR:	Negative Likelihood Ratio
NSCLC:	Non-Small Cell Lung Cancer
OR:	Odds Ratio
PLR:	Positive Likelihood Ratio
PPV:	Positive Predictive Value
ROC:	Receiver Operating Curve
SCLC:	Small Cell Lung Cancer
SUHT:	Southampton University Hospital Trust
SVC:	Superior Vena Cava

Chapter 1: Overview

1.1 Thesis layout and contents

Chapter One briefly introduces the rationale for the study; highlighting the implications of late lung cancer diagnosis, and the challenges associated with improving earlier symptomatic detection. Difficulties resulting from the lack of a clear “symptomology” in lung cancer, and implications of symptom overlap with comorbidities such as chronic obstructive pulmonary diseases (COPD). The epidemiology of primary lung cancer and COPD are presented in **Chapter Two** (narrative literature review) to provide the context for this study.

A systematic review of the literature available on symptomatic diagnosis of lung cancer was conducted and included in **Chapter Three**. Parts of this chapter have been published in the international journal, *Family Practice*.

Chapter Four of the thesis provides a high-level overview of the study design and methodology. The exploratory process of this feasibility study is structured into two studies; Study 1 (**Chapter Five**) and Study 2 (**Chapter Six**). A detailed account of the methodology and methods relating to each study (Study 1 and Study 2) is provided in the relevant study chapter (see Chapter Five and Six).

Chapter Five reports the methodology, methods and results from Study 1; this study used qualitative research methods to explore the acceptability and content validity of a patient-completed questionnaire in a secondary care population who had been referred to lung-shadow clinic for lung cancer investigation. Themes derived from participants’ accounts of their symptom experiences, and experiences of questionnaire completion, are presented and discussed.

The subsequent chapter, **Chapter Six**, details methodology, methods, and results of Study 2, which quantitatively identifies symptom(s) and combinations of symptoms that predict lung cancer in this referred population and in a sub-population with COPD (sub-group analysis). In this chapter, the development and evaluation of a set of diagnostic criteria is discussed. **Chapter Six** also addresses the handling of missing values in the questionnaire data.

Overview

The discussion chapter, **Chapter Seven**, interprets the results obtained from Study 1 and Study 2 with respect to the study aims, methodology, current literature, strengths, and limitations of the study.

Chapter Eight of the thesis draws conclusions on the basis of the discussion in Chapter Seven, explores the potential practical application of the results, and makes recommendations for further research.

1.2 Introduction

In the UK, there were 43,463 new cases of lung cancer diagnosed in 2011 (Cancer Research UK [CRUK] 2014). Although lung cancer is the most common type of cancer that a GP will encounter, the average GP will see about two cases annually (Hamilton and Sharp 2004; Scottish Cancer Registry, Information Services Division 2013; Parkin et al. 2005). Viewed in the context of the many thousands of patients that present with non-specific symptoms such as cough, breathlessness, and chest pain, the challenges associated with symptomatic diagnosis of lung cancer in primary care are clear, particularly in a population with existing comorbidities such as COPD (Neal et al. 2014; Mitchell et al. 2013).

The implications of these difficulties are delayed diagnosis, and poorer prognosis; 86% of people with lung cancer are detected at late stages when curative intervention is no longer viable with a UK 5-year survival rate of 8%, lower than its Western counterparts (Office for National Statistics [ONS] 2009; Health and Social Care Information Centre [HSCIC] 2011). Historically, lung cancer has been perceived as a 'silent' disease that remains asymptomatic until the disease has metastasised; forming the basis for most lung cancer screening trials (Peake 2015; Kumar and Clarke 2005).

Recent studies have suggested that lung cancer can be symptomatic even years prior to diagnosis independent of disease stage when diagnosed (Hamilton et al. 2005; Corner et al. 2005). These findings, along with evidence that later stage at diagnosis is partly responsible for poorer lung cancer survival in the UK have led to national attempts to address late diagnosis in the UK, which included the National Institute of Clinical Excellence (NICE) referral guidelines, and public awareness programs such as 'Be Clear on Cancer' (CRUK 2014).

However, the recommendations in the NICE guidelines; for urgent chest X-ray be offered to patients presenting with any one of the 10 unexplained symptoms, were based on weak evidence (Hamilton and Sharp 2004). There is a need for evidence about the predictive values of lung cancer symptoms in primary care to inform GPs referral decisions. The situation is further compounded by the absence of an identifiable lesion (e.g. palpable lump in breast cancer) and a vague symptom profile. Most of these symptoms are non-specific and could also suggest other differential diagnosis of lung and respiratory diseases such as COPD.

COPD precedes lung cancer diagnosis in 40% to 90% of the cases (Young et al. 2009), and those with COPD are generally diagnosed with lung cancer at an even later stage than individuals without COPD; resulting in poorer prognosis, and lower survival rate (Kiri et al. 2010; Powell et al. 2013). Explanations for late lung cancer diagnosis in COPD include symptom overlaps (presentations complicated by comorbidity), and the presentation of non-specific symptoms such as coughing and breathlessness, which made it difficult to distinguish between symptoms of lung cancer and COPD (Kiri et al. 2010). Furthermore, existing impairments to the lung as a result of the COPD, also reduces the suitability for curative, surgical intervention. Without controlling for smoking, patients with COPD are four times more likely to develop lung cancer (Kennedy et al. 1996; Diez-Herranz 2001; Wasswu-Kintu et al. 2005). Evidence regarding symptoms that predict lung cancer in patients with existing chest and respiratory comorbidities, such as COPD, would help to improve earlier diagnosis in this prevalent sub-group of lung cancer patients.

Much of the evidence that underpins current NICE referral guidelines, particularly in lung cancer, are weak and did not derive from large prospective studies in primary care, which reflects the large study sizes and costs involved in such longitudinal studies (Mulka 2005; NICE 2011; Scottish Executive 2002). Currently, the evidence on PPVs for use in primary care derives from retrospective studies that extract symptom data from routine primary care records. However, retrospective data from medical and GP records often tends to be diagnosis-driven rather than symptom-focused (Kroenke 2001). Symptom reporting is naturally closely related to the perceived diagnosis, which may influence the symptoms recorded in GP notes, resulting in recording bias. Data collected from secondary sources are also subject to missing data

Overview

(incompleteness in data) and recording errors. Other methodological issues associated with these studies have been highlighted in the Systematic Review (Chapter Three).

Questionnaire development and validation

The development of a valid and reliable questionnaire is essential to minimise measurement error; that is, where there is discrepancy between the respondent's characteristics and their responses in the questionnaire (Groves 2005). In general, the process of developing and testing questionnaires assumes the following sequential steps: 1) background, 2) conceptualisation of questionnaire, 3) format and data analysis, and 4) instrument testing (establishing the validity and reliability of the questionnaire in a pilot study) (Brancato et al. 2006; Radhakrishna 2007). The developmental process is iterative; making necessary changes to improve and testing the process before moving on to the next step, which is time-consuming.

Like any empirical study, the first step of the process explores the purpose, objectives, research questions, and hypothesis of the research. This also involves careful identification of the target audience, educational or literacy level, and accessibility to respondents (Radhakrishna 2007). A good understanding of the research through the literature, will form the basis for the next part of the development process; questionnaire conceptualisation.

Questionnaire conceptualisation is essentially the content construction part where items or questions are generated for the questionnaire (Brancato et al. 2006; Radhakrishna 2007). Independent variables would be identified and defined in this step. An aim of this step is to ensure the translation of the study objective into the content. Format and data analysis of the questionnaire concentrates on writing questionnaire items, selection of appropriate scales of measurement (e.g. nominal, ordinal, or ratio), questionnaire layout, question sequencing, formatting, and the proposed data analysis (Radhakrishna 2007).

Validity is the amount of systematic or built-in error in a measurement (Norland, 1990). The type of validity (content, construct, criterion, and face) depends on the study objectives. Content validity ascertains the appropriateness of the questionnaire and its relevance to the study purpose. It establishes that the content of the questionnaire reflects a complete range of the attributes under investigation (Pilot and Hunger 1999; DeVon et al. 2007).

Usually, once the conceptual framework of the research study is established, the questionnaire is reviewed by several experts relevant to the field to determine its content validity and ensure that the questionnaire is consistent with the conceptual framework (Pilot and Hunger 1999; DeVon et al. 2007). Closely related to content validity, face validity refers to the extent to which the measure **appears** to measure a certain criterion. Criterion validity is the extent to which a measure is related to an outcome, whereas construct validity refers to the degree of compatibility between the operationalisation a construct and theory (Finch et al 2002).

In addition to validity, a readability test can also be used to ensure readability. The questionnaire can then be tested on members not included in the target sample, and changes can be made accordingly before the questionnaire is piloted (Brancato et al. 2006; Radhakrishna 2007). Lastly, the reliability of the questionnaire is tested in the pilot study. Reliability refers to the accuracy and consistency of the instrument in measuring (Norland 1990). Computational statistical tests can be used to measure reliability.

Development and validation of the IPCARD questionnaire

The IPCARD self-completion questionnaire was designed for use in prospective studies to identify predictive values of symptoms for lung cancer diagnosis. The questionnaire was developed, initially, in a population recently diagnosed with operable lung cancer. Semi-structured interviews were used to inform the design of items that recorded the full range of symptoms experiences by patients with operable lung cancer in the two years before lung cancer diagnosis. The questionnaire also included items informed by the evidence base regarding symptoms that are independently associated with lung cancer reported in the literature, NICE Referral Guidelines, the International Primary Care Respiratory Group guidelines (Levy et al. 2006), risk information (Cassidy et al. 2008), and co-morbidities.

In depth analyses of the semi-structured interviews suggested that the use of medical symptom terminology could lead to patients normalising changes in health and denying the presence of symptoms (Brindle et al. 2012). Furthermore, findings of a previous qualitative study identified lay experiences of symptoms that distinguished between those with and without ovarian cancer (Bankhead et al. 2008). Generic symptom descriptors, used to categorise

Overview

symptoms potentially predictive of lung cancer (e.g. cough and breathlessness) refer to symptoms at a high prevalence in primary care consulting populations, and might not distinguish between those with lung cancer and those presenting with a range of other respiratory diseases. Therefore, in the development of the IPCARD questionnaire, non-medical and non-disease terminology, and lay descriptors of bodily sensations and symptoms, were included to more effectively elicit patients' experiences of health changes and increase the possibility of identifying symptoms that could distinguish between lung cancer and those without. Lay symptom descriptors were identified from the interviews with operable lung cancer patients (Brindle et al. 2012) and previous research with inoperable lung cancer patients (Corner et al. 2005). Lay descriptors of breathlessness were also identified studies with patients with asthma, COPD, interstitial lung disease, cardiac failure and lung cancer (Wilcock et al. 2002). An iterative process of semi-structured and cognitive interviews was used to redesign questionnaire items to improve content validity. Lay members of the consumer research panels in Southampton and Birmingham, chest physicians and radiographers at SUHT were also involved in questionnaire and study design.

The IPCARD questionnaire was validated in a population of GP referred chest X-ray attendees with a range of chest and respiratory diseases, and found to accurately record the symptom experiences of these individuals. Good content validity, data completion, and recruitment rates were established (Brindle et al. 2014; Brindle et al. 2015). Hence, the IPCARD questionnaire provides a valid method of capturing detailed information regarding a wide range of symptoms that could discriminate between those with lung cancer and those without lung cancer.

Research in symptoms

Presenting symptoms are those that prompt a person to seek health care but may not necessarily be the first symptom experienced. Many potential cancer symptoms presented in primary care are non-specific. Even though they could be caused by cancer, more often than not, their cause is benign; for example, they might be symptoms of COPD and other comorbidities (Neal et al. 2014; Mitchell et al. 2013). Thus, many coincidental associations of symptoms with cancer may be incorrectly interpreted as causal (Kroenke 2001; Mitchell et al.

2013). For GPs to be able to make informed judgements about the significance of symptoms, prospective symptom-based studies in primary care are needed to identify symptoms of lung cancer, and their positive predictive values (Summerton 2002).

Conducting research in early cancer diagnosis is challenging, and requires the careful consideration of the methodological approaches to sampling and measurement of symptomology, as recommended in the Aarhus statement (Weller et al. 2012). Cancer-related symptoms are very complex, and vary greatly depending on the method of documentation and elicitation of these symptoms (whether it is extracted by clinical checklist or clinical scales; patient-completed questionnaire, or spontaneous reporting). Current literature on symptomatic lung cancer diagnosis suggests that these factors have not yet been addressed in the study design of retrospective studies (Hamilton et al. 2005).

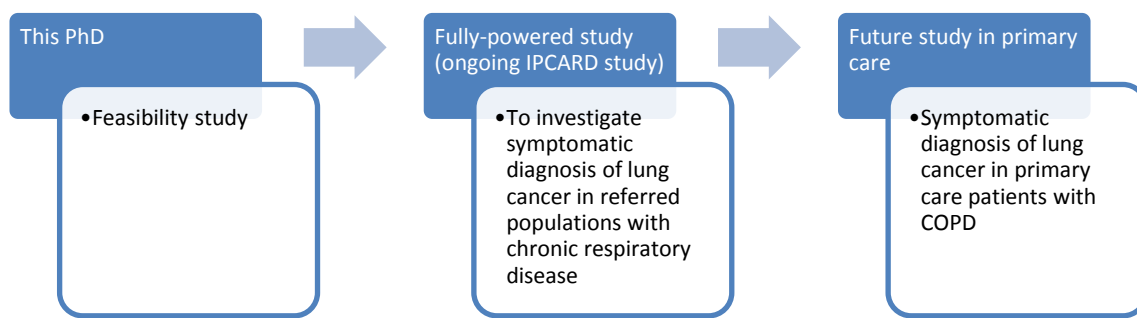
Hence, the present study proposes to use a previously designed participant-administered questionnaire (the Identifying Symptoms that Predict Chest and Respiratory Disease (IPCARD) questionnaire) to systematically elicit detailed information about patient-reported symptom experiences or changes in health status (Brindle et al. 2012).

Rationale of study

This feasibility study proposes to use the IPCARD questionnaire to systematically and prospectively record symptom experience of patients who had been referred to lung-shadow clinics. The study will estimate the discriminatory values (sensitivity, specificity, and likelihood ratios) of the symptoms recorded to explore the feasibility of symptomatic diagnosis of lung cancer in a population with high rates of chronic respiratory problems, and inform a definitive secondary care study (fully-powered); see Figure 1.1. Findings from this prospective, exploratory study would add to the current evidence on predictive value of symptoms for lung cancer diagnosis in secondary care populations (Kubik et al. 2001; Hoppe 1977), and contribute new evidence regarding the feasibility of the symptomatic diagnosis of lung cancer in a secondary population with COPD.

Figure 1.1 Schematic plan of the PhD in relations to the larger IPCARD study

Overview



There is currently no recommended lung cancer screening mechanism in the UK, even for high-risk groups. However, even with the existence of cancer screening programmes, 80-90% of cancers are diagnosed following symptomatic presentation to primary care (Hamilton and Peters 2007; Yoder 2006). Therefore, improving earlier symptomatic detection is of particular importance if lung cancer survival is to be improved in the UK.

1.3 Aims

- (1) To explore acceptability to patients and content validity of the IPCARD questionnaire in a population referred to secondary care on suspicion of lung cancer (a population referred to a lung-shadow clinic).
- (2) To explore the feasibility of using a patient-completed symptom questionnaire (IPCARD) to identify symptoms that predict lung cancer in a population referred to a lung-shadow clinic.
- (3) To identify patient-elicited symptoms that are independently associated with the diagnosis of lung cancer (a population with high rates of chronic respiratory disease).
- (4) To identify patient-elicited symptoms that are independently associated with the diagnosis of lung cancer in a population with COPD referred to secondary care on suspicion of lung cancer.

1.4 Research Questions

- (1) Is the IPCARD patient-completed questionnaire acceptable to a secondary care population that has been referred for lung cancer investigations?
- (2) Does the IPCARD questionnaire accurately capture the full range of symptoms and comorbidities experienced by those with COPD in a secondary care population that has been referred for lung cancer investigations?
- (3) Is it possible to identify lung cancer on the basis of symptoms in a population with high rates of respiratory disease that has been referred to secondary care on suspicion of lung cancer?
- (4) Is it possible to identify lung cancer on the basis of symptoms in a population with COPD that has been referred to secondary care on suspicion of lung cancer?

Chapter 2: Literature Review

LUNG CANCER

2.1 Pathophysiology

Lung cancer (LC) refers to a malignant growth or tumour in the respiratory tract and can be categorised into bronchial carcinoma and alveolar cell carcinoma. Lung cancer develops when epithelial cells in the airways of the lung become abnormal. A tumour results when the proliferated cells build up into a mass. The term bronchial or bronchogenic carcinoma was so named from the common occurrence of the tumour growths in the bronchial epithelium. Carcinoma refers to the malignant tumour that is formed from the epithelial cells. Bronchogenic carcinoma accounts for 95% of all primary lung cancer and hence is the most common lung malignancy in European countries (bronchial carcinoids comprise of the remaining 5% of lung cancers) (Kumar and Clark 2005).

Most lung cancer develops slowly and it is said that abnormal cell change can start years before it is detected (British Lung Foundation 2012). As a vital organ, healthy lungs possess a large reserve capacity to meet the body's demand for oxygen and do not have pain receptors. It is also this reserve capacity that allows the carcinoma to grow for years without compromising lung function or causing any pain symptoms; evading detection.

Histologically, based on the cellular differentiation, bronchial carcinoma has been expressed as squamous cell carcinoma, adenocarcinoma, large cell carcinoma and small cell carcinoma. However, for therapeutic purposes, these different types of bronchial carcinomas are broadly categorised into non-small cell lung cancer (NSCLC) and small cell lung cancer (SCLC) with estimated incidence percentage of 70% - 75% and 20% - 25% respectively. SCLC and NSCLC have distinct patterns of growth and spread (Feld et al. 1995).

The SCLC spreads aggressively and is highly malignant and very often non-amenable to surgical intervention when diagnosed. They are also treated differently. NSCLCs are said to have better prognosis than SCLC as they are slower growing (Feld et al. 1995). Hence, SCLC responds better to

chemotherapy and radiation, whereas NSCLC demonstrates better prognosis with curative surgery usually in the form of a lobectomy or pneumonectomy, which allows for full removal of the malignance. For SCLC, the first-line treatment is usually chemotherapy; sometimes in adjunct to radiotherapy (Feld et al. 1995). The median number of survival years with treatment for SCLC is one year (SIGN 2005) in comparison to 5 years for resectable NSCLC (Summerton 1999).

However, these are just general deductions of prognosis according to histological data, which are inconclusive and do not necessarily dictate the type of radical treatment given to the individual (Kumar and Clark 2005).

2.2 Epidemiology

With incidence rates of 1.35 million, lung cancer is to be the most common cause of cancer mortality resulting in an estimated 1.18 million number of deaths worldwide in a year (Parkin et al. 2005). More than 38,000 new cases of lung cancer are diagnosed each year in the UK. In 2011, there were 43,463 new cases of lung cancer in the UK (CRUK 2014).

Despite the improvement in diagnostics, and treatment modalities in lung cancer, patient prognosis remained poor with the highest 5-year survival period reported worldwide to be a bleak 18% (Landis et al. 1998; Jemal et al. 2003; Salomaa et al. 2005; American Cancer Society 2006). In the UK, lung cancer is still diagnosed at a later stage and has the lowest 5-year survival rate with figures of 5% to 8% in comparison to its European counterparts (Janssen-Heijnen et al. 1998; Coleman et al. 1999; 2003; CRUK 2004; Richards 2007; Coleman et al. 2011). In 2005, the National Lung Cancer Audit (LUCADA) programme was launched to collect epidemiology data on the incidence, nature, geographic distribution and treatment of lung cancer, aimed at improving patient care and outcomes (HSCIC 2006).

Lung cancer accounts for almost 7% of all deaths in the UK. Table 2.1 shows the incidences of mortality in the UK among men and women in 2011, which reported a total of 35,371 deaths; an increase from 34,859 lung cancer deaths in 2009. An average crude mortality rate of 69 suggests that there were about 69 lung cancer deaths for every 100,000 person of the population. With

considerations to the demographic changes, and population growth, deaths by cancer in general is still predicted to increase from 7.1 million in 2002 to 11.5 million in 2030 (Mathers and Loncar 2006).

Table 2.1 Lung cancer statistics- number of deaths and crude rate per 100,000 of the UK population in 2012 (CRUK 2014)

	Deaths	Crude rate
Male	19,304	76.5
Female	16,067	61.2
Total population	35,371	68.7 (average)

A higher incident rate for lung cancer is generally observed in men than in women with men accounting for almost 60% of the lung cancer cases in the UK (NICE, 2005). The male: female ratio for lung cancer was reported to be 39: 10 in 1975 but has since then declined considerably (CRUK 2009). New trends in lung cancer prevalence have emerged over recent years. One in particular is the increasing number of lung cancer diagnosis in women particularly in developed countries (Parkin 1998; Wingo et al. 1999), which could be related to the social evolution of cigarette smoking, and its consequent causal link to increased risk of lung cancer. In the Liverpool Lung Project, Field and Youngson (2002) observed a 30% rise in female smokers below the age of 75 years old between 1992 and 1995 within Liverpool.

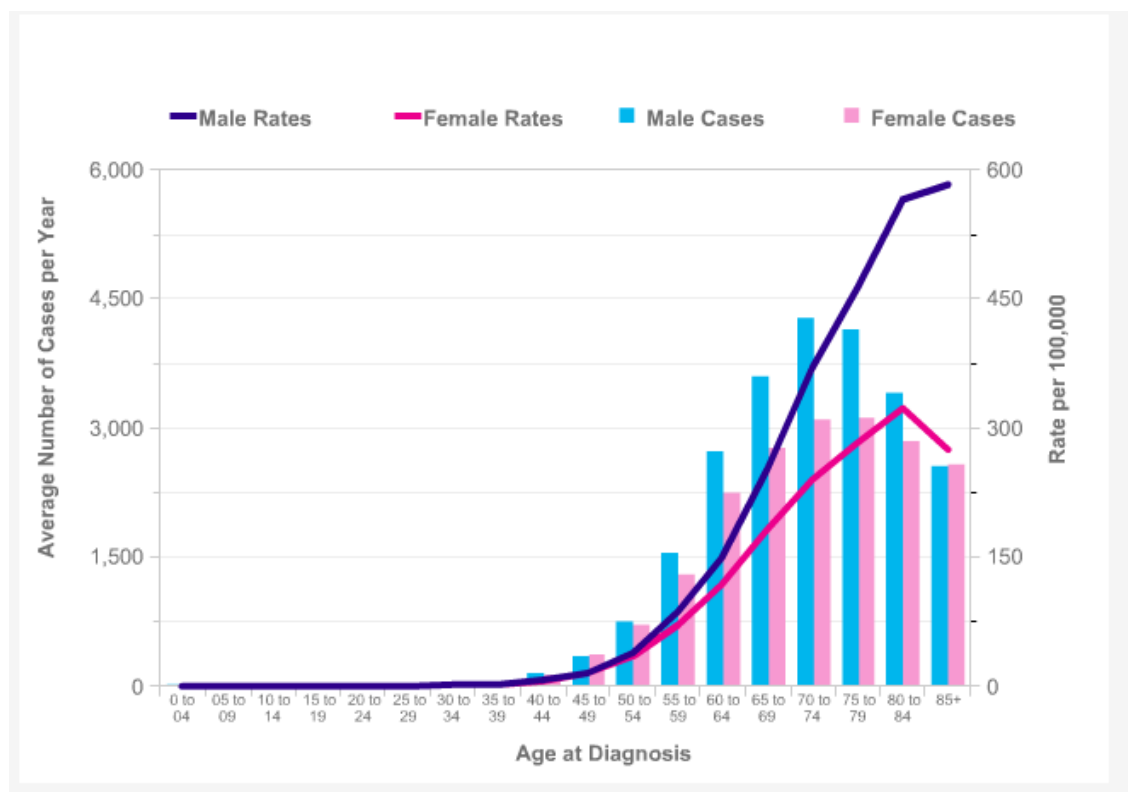


Figure 2.1 Average number of new cases per year and age-specific incidence rates between 2009 and 2011 (CRUK 2014; ONS 2012)

Mason (1941) described this as a disease of the middle ages, with the vast majority of lung cancer diagnosed in people over the age of 40. The graph shows the incidence rates of lung cancer for men and women in the UK and how lung cancer incidence is strongly related to age (see figure 2.1). Between 2008 and 2010, almost 75% of the new cases of lung cancer were diagnosed in men and women above the age of 65 (see figure 2.1). We can see the age-specific incidence rates rise steeply from age 40 for both genders, peaking at age 80-84. Male rates are generally higher than females, and this gap increases with age (steeper slope in males than females) (see figure 2.1).

Having said that, today's trend shows increasing incidences of lung cancer diagnosis observed in a younger population which could have much to do with the evolution of cigarette smoking in the younger population. The average start age of cigarette smoking in the UK has declined over the years as adolescents start smoking at a younger age (Field and Youngson 2002). Consequentially, this could translate in the age criteria of participants lowering in future lung cancer research and younger participants being included.

2.3 Aetiology and Risk factors

Substantial amount of research has been done on the aetiology of lung cancer and its associated risk factors such as smoking, age, occupational exposures, genetic predisposition, diet, and social class (Shields 2000; Steenland et al. 2001; Bofetta and Kogevinas 1999; ONS 1997). Having said this, the strength of interaction between cigarette smoking and lung cancer overshadows all other aetiological factors. 90% of lung cancer deaths in the UK have been attributed to smoking. Smokers are 15 times more at risk of lung cancer than non-smokers (Summerton 1999).

Conceptually, it should be noted that the total sum of the proportions of a disease attributable to various risk factors will not necessarily amount to 100% due to the numerous pathways in the carcinogenic process. This means that a 90% risk of lung cancer attributable to cigarette smoking for example, does not inevitably suggest that 10% is contributed to all other risk factors. This is based on the concept of interactions in the contribution to risk (Rothman 1986).

2.3.1 Smoking

Lung cancer research has already established a strong causal association between cigarette smoking and the disease prevalence which also considered different variables such as the starting age, duration, frequency, and level of smoking. Studies have also investigated the effects of passive smoking in correlation to lung cancer risk (Doll and Hill 1950; Parkin et al. 1995; Hackshaw et al. 1997). A meta-analysis by Hackshaw et al. (1997) reported an increased risk of incidence of lung cancer in non-smoking spouses of smokers by 24%, resulting in an increased relative risk factor of 1.24.

Figure 2.2 demonstrates how trends in lung cancer incidence rates reflect trends in previous cigarette smoking prevalence especially in men. Generally, a decline in smoking rates is followed by a decrease in lung cancer rates for some decades. It was proposed that smoking prevalence peaked in men before women, and therefore, the decrease in lung cancer rates observed in men is yet to happen in women (CRUK 2014).

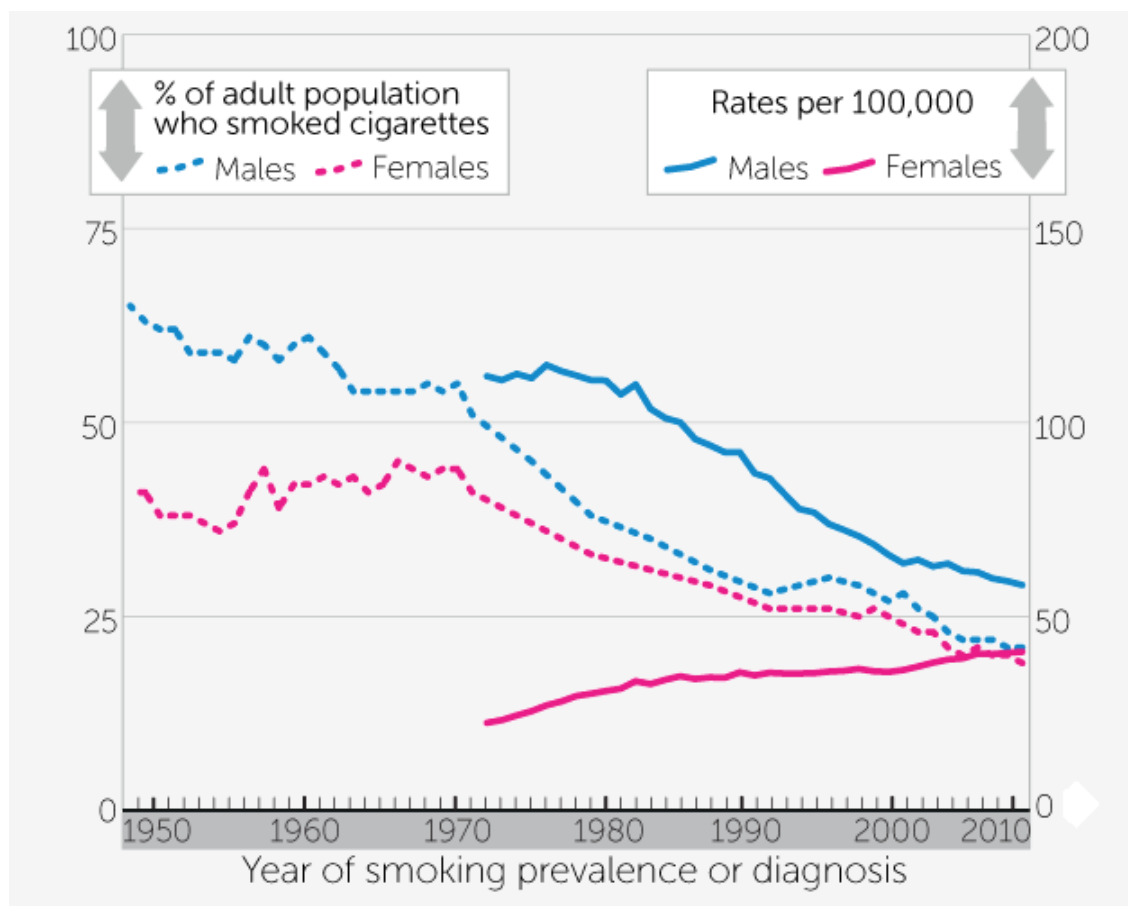


Figure 2.2 Trends in lung cancer age-standardised incidence rates, and smoking prevalence in Britain from 1948 to 2011 (CRUK 2014)

2.3.2 Occupational exposures

Occupational factors have also been related to exposure of carcinogenic chemicals in industrial and agricultural settings such as arsenic, petroleum based products, radiation, coal, and its products of coal combustion. Several studies looked into the effects of asbestos exposure attributing to lung cancer but results are inconclusive (Magnani et al. 2000; Hauptmann et al. 2002; Mastrangelo et al. 2008). However, it is noted that tumours relating to such occupational hazards are often adenocarcinomas, less commonly associated with cigarette smoking. Findings have shown a higher incidence in bronchial carcinoma among the urban population in comparison to the rural residents, even after adjusting to the higher smoking rate in cities. This had led to hypotheses associating environmental factors such as air pollution to lung cancer. In line with this, the British Lung Foundation has previously reported more than 40 carcinogenic compounds were found in air pollutants which

subsequently supported the premise that excessive exposure to air pollution could attribute to lung cancer (British Lung Foundation 1998; Tomatis 1990).

2.3.3 Socio-economic factors

Survival rates in lung cancer have been found to be inversely associated to socio-economic status (Pearce and Bethwaite 1997). Tomatis and colleagues (1997) attributed this observation to inequalities in healthcare between social classes which presented in the form of lack of access to healthcare and health education leading to a lack of understanding of the implications of the symptoms. This could be partly due to a juxtaposition of all the aforementioned factors (diet, environmental exposure, higher smoking incidences due to lack of information) in relations to different socio-geographical conditions.

2.3.4 Genetic and other factors

Other factors include family history and genetic predisposition where research showed increased risk in non-smoking family members of families with a presenting history of lung cancer (Tokuhata and Lilienfeld 1963). It was theorised that there exists tumour suppressor genes called oncogenes or a genetic predisposition to break down (metabolise) carcinogens, which could affect the development of lung cancer (Tokuhata and Lilienfeld 1963). Studies of familial predisposition to lung cancer have shown an association between rare autosomal genes (e.g. Mandel cancer susceptibility), and the development of lung cancer in young individuals (< 50 years). This gene was not found in older people who developed lung cancer, which suggests that the cause of cancer in these non-carriers was more likely to be long-term exposure to tobacco (Sellers et al. 1990). However, because genetic epidemiology in general looks at these smaller risk factors, one might not expect to find hugely significant findings to suggest definite relations in comparison to what is known about cigarette smoking. Also, these studies will not be able to rule out the collective effects of communal living within a household. To do so, careful consideration of the study design would be needed. Within familial aggregation of cancer due to inherited susceptibility, the variation of penetrance susceptibility alleles for lung cancer had been thoroughly researched (Eisen et

al. 2008). Highly penetrant mutations in known genes have been suggested to account for the increased risk of certain cancers.

2.4 Signs and Symptoms

It has been a common perception that the clinical presentation of bronchial carcinomas to be fatally silent; where lesions are usually inoperable at the time of symptom presentation (Kumar and Clarke 2005). People with lung cancer will eventually develop symptoms just that most present at the advanced stage of the cancer (Kumar and Clarke 2005). Only a minority (10%) of patients with lung cancer were reportedly asymptomatic at the time of diagnosis (Scagliotti 2001). 27% experienced symptoms related directly to the primary tumour (Summerton 1999; Scagliotti 2001). Majority of the patients (44%) experienced non-specific, systemic symptoms such as weight loss and fatigue (Yoder 2006; Beckles et al. 2003). Prognosis was agreeably better for those who were asymptomatic at the time of diagnosis as a result of spontaneous discovery of the malignancy (Beckles et al. 2003; Buccheri and Ferrigno 2004).

The presentation of lung cancer is highly variable which depends on factors such as tumour site and involvements of lymph nodes or other organs. This permits for a wide variety of possible symptoms which could involve distant organs in the case of secondary metastasis. For simplicity, only symptoms of localised lung cancer and systemic symptoms will be discussed. Chronic cough (cough that lasted more than 3 weeks) has been repeatedly noted as the first symptom of lung cancer and the most commonly reported symptom (Summerton 1999; Liedekerken et al. 1997; Buccheri and Ferrigno 2004). The frequency of cough in patients with diagnosed lung cancer varies between 21% and 87% depending on the stage of their cancer and the study design but chronic cough eventually occurs in up to 90% of lung cancer patients (Beckles et al. 2003; Liedekerken et al. 1997). However, most of the evidence on symptoms statistics provided in this non-systematic review by Beckles et al. (2003) are based on non-empirical studies, and/or relied on retrospective data of symptoms, not systematically collected for research purposes. This should be reflected on when drawing inferences from these studies.

With the increasing growth in tumour size, ulceration can occur and any discharge as a result of it will be presented as expectoration; or any bleeding

can manifest as haemoptysis. Haemoptysis or coughing up blood is another symptom that is highly suggestive of lung cancer with reported frequency ranging from 6% to 35%, and PPVs of 2.4 (Hamilton et al. 2005) to 6.4 (Hippesley-Cox and Coupland 2011). The simplified mechanism of the pathology involved the obstruction of the parent bronchus which could be incomplete or complete. In the event of incomplete obstruction with infection, purulent sputum can be coughed up. However, if the bronchus is fully obstructed, the affected lung collapses resulting in symptoms of dyspnoea and shortness of breath (Beckles et al. 2003; Kumar and Clark 2005).

Other common complaints include unintentional weight loss (reported frequency of up to 68%), wheezing and chest pain. Chest discomfort occurs in more than 50 % of patients at diagnosis (Beckles et al. 2003). Beckles et al. (2003) characterised the chest pain experienced in lung cancer as dull, persistent, poorly localised, and not associated to breathing or coughing. In the advanced disease, some of the symptoms are thought to be caused by the direct cancer invasion of the structures around the lungs or local lymph nodes. For instance, 8% to 15% of those with lung cancer experience pleuritic chest pain due to pleural involvement (Beckles et al. 2003). Other symptoms such as persistent hoarseness could result from malignant tumour affecting the vocal cords, and dysphagia or difficulty swallowing could be attributed by a tumour compressing on the oesophagus. Apical lung cancer in the apex of the lung or Pancoast tumour) could result in shoulder pain with or without referring arm and hand weaknesses. Another symptom, superior vena cava (SVC) syndromes; when the lung tumour compresses the SVC vein, will cause facial swelling and most of the symptoms mentioned. Chest pain or discomfort is another common symptom associated with lung cancer (Beckles et al. 2003).

Distant symptoms of metastasis can affect any parts of the body such as the brain, liver and bones causing associated symptoms. However, it is said that when systemic symptom (weakness, anorexia, fatigue, and weight loss) or metastatic symptoms emerge at the time of diagnosis, disease is usually in the advanced stage and prognosis is bleak (Buccheri and Ferrigno 2004). 3% to 10% of all lung cancer patients develop other paraneoplastic syndromes secondary to the carcinoma which include hypercalcemia, Cushing syndrome, neuromuscular pathologies e.g. peripheral neuropathy and polymyositis (Kumar and Clark 2005).

Literature Review

On physical examination, abnormal signs such as presence of any persistent localised chest sign (e.g. clinical features of lung collapse, consolidation, abnormal chest sounds by stethoscope or pleural effusion) will also help to identify the need for further investigation for lung cancer (SIGN 2005; NICE 2011).

At present, NICE guidelines recommend urgent referral for chest X-ray for those who presents with haemoptysis, or any persistent (defined as longer than three weeks) signs and symptoms of:

- cough
- chest/shoulder pain
- dyspnoea
- weight loss
- chest signs
- hoarseness
- finger clubbing
- features suggestive of metastasis from a lung cancer (e.g. in brain, bone, liver, or skin)
- Cervical/supraclavicular lymphadenopathy (NICE 2005).

The basis for this recommendation largely relied on Beckles's non-systematic review (2003) which included extrapolated results from case-control studies and cohort studies (Hyde and Hyde 1974; Andersen and Prakash 1982; Grippi 1990; Scagliotti 2001).

2.5 Early detection and screening for Lung Cancer

The classical biomedical model of cancer suggests that cancer is preceded by precursors that may evolve into early cancer, at which point it is still asymptomatic and not clinically detectable. Symptoms only appear when the cancer progressed into a bigger tumour, and into the advanced stage of the disease (Zheng et al. 2011). Figure illustrates the classical model of cancer development. Screening efforts are meant to target early detection of the cancer, when the disease is asymptomatic.

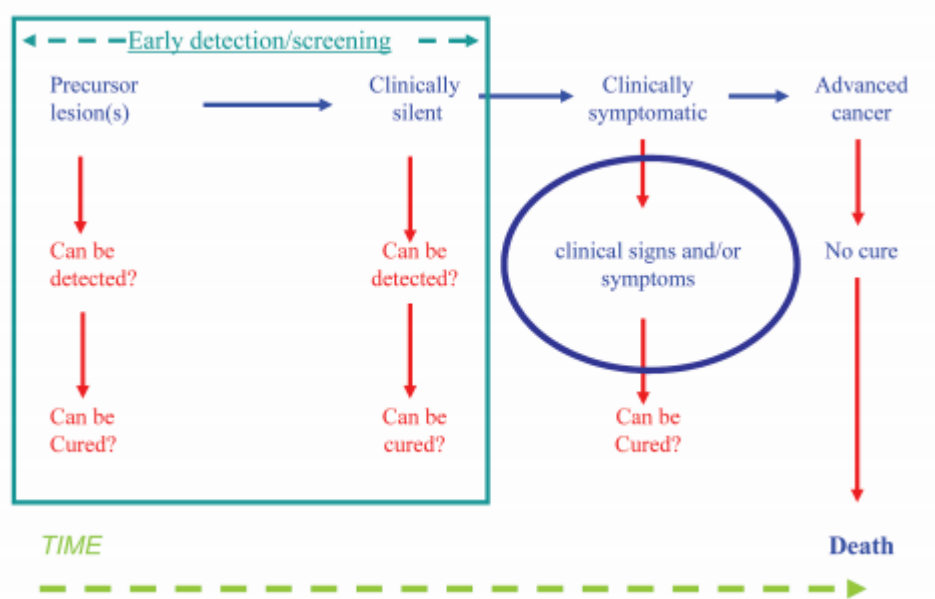


Figure 2.3 Progression of cancer (classical model) (Zheng et al. 2011)

Late diagnosis has been widely associated with the poor lung cancer outcomes and low survival rate of lung cancer in the UK, with 86% of patients being diagnosed with advanced stages of the disease (ONS 2009; HSCIC 2011). This has stimulated much interest in the possibility of screening asymptomatic individuals to improve early lung cancer detection (Porter and Spiro 2000). Screening trials have been driven by the belief that lung cancer has a long, silent period of latency before it is diagnosed; such that at the point of symptomatic presentation, the likely long-term survival benefits will be small (The National Lung Screening Trial Research 2011).

Although, routine screening programmes have been proven effective for certain cancers, and are currently recommended - faecal occult blood tests for colorectal cancer, mammograms, and self-breast examinations for breast cancer, and annual pap smears for cervical cancer screening (Hamilton et al. 2006; 2009; Jones et al. 2007; Ellis and Thompson 2005; Goff et al. 2004; Barton et al. 1999), the efficacy and feasibility of lung cancer screening is still inconclusive and debatable. The earliest randomised controlled trial in the 1960's compared the diagnosis, and resection rates of 6-monthly chest X-ray screening to chest X-ray in the beginning and the end of three years. No difference in mortality was found between the two groups (Brett 1968). This was followed by a Czechoslovakian Lung Cancer screening study (Kubik et al. 1990) and the Mayo Lung Project, a large randomised controlled trial (Fontana

Literature Review

et al. 1984; 1986). The Mayo Lung Project was one of the largest studies in lung cancer screening to compare the use of either chest X-ray and/or sputum cytology to standard care (Fontana et al. 1984; 1986). Findings from these early studies essentially triggered the nihilism around the effectiveness of lung cancer screening.

However, over the years, these studies have been heavily criticised for being underpowered and for having high contamination rates in the control groups. The Mayo Lung Project presented a study power of less than 20% to support a 10% advantage in improving lung cancer mortality. 73% of the participants in the control group had previous chest radiography in the final two years which posed as contamination of the control group. There were also issues of compliance between the screened groups (75% compliance) and the control group (50% were compliant to the advice given) (Fontana et al. 1984; 1986). Despite a significantly increased incidence in lung cancer diagnosis ($p = 0.016$) observed in the experimental (screened) group, no difference in the disease mortality can be found between the two groups (screened and control) in the Mayo Lung Project. Strauss (1997) suspected there to be a degree of over-diagnosis of clinically insignificant tumour in the screened group and also occurrences of under-reporting within the control group (where the control population died as a result of existing morbidities before diagnosed with lung cancer); adding doubt to the study design of the Mayo Lung Project. Also, the possibility of confounders could not be ruled out given that the study was largely heterogeneous, and did not provide information of exposures to other risk factors (such as occupational exposure and genetic predisposition) and existence of co-morbidities (e.g. chronic obstructive pulmonary disease). Hence, evaluation of the data provided by these studies suggested flaws in both the study design and interpretation.

Findings of the randomised prospective study in Czechoslovakia to screen for lung cancer in a population of heavy smokers (high-risk of lung cancer population) using chest X-rays and sputum cytology also revealed no significant changes to mortality between the screened (intervention) and non-screened (control) group; i.e. no measurable benefit (Kubik et al. 1990) replicating the findings of the Mayo Lung Project. None of these studies mentioned, had a 'no screening' control group. Another study, the Memorial Sloan-Kettering study included additional sputum cytology to annual chest X-

22

ray also showed a lack of improvement in lung cancer mortality (Melamed et al. 1984). Based on evidence from these studies, it was not possible to recommend a chest X-ray screening program for lung cancer.

The turn of the century welcomed new sparked interest in earlier lung cancer diagnosis looking at the use of helical computed tomography (CT) scanning and other biomedical markers following studies by Tockman and colleagues (1988; 1994) and the Matsumoto study by Sone et al. (1998). More recent studies have supported the effectiveness of using low-dose CT; Veronesi et al. (2012) reported a 20% reduction in mortality from lung cancer, and a 6% decrease in overall mortality (all causes). However, severe costs were incurred and over-diagnosis (identification of idle cancers that would never have become clinically apparent) was also a problem. Therefore, the view on the efficacy of low dose CT screening remains highly polarised.

Parallel to these studies, the Prostate, Lung, Colon and Ovarian (PLCO) study, which involved ten screening centres located across the United States to assess the effects of annual screening using modern chest X-rays was launched in view of the limitations of previous studies reported. The study recruited 77,445 subjects and 77,456 controls based on the National Lung Screening Trial (NLST) trial inclusion criteria. Their report showed no evidence of improved mortality over 13 years in the screened group or the high-risk subgroup, to suggest any true benefit.

To date, the benefits of screening is not fully justified, and the general consensus around the efficacy of different screening tools in detecting early lung cancer still remains inconclusive, which reflects the current absence of a lung cancer screening program nationally.

2.6 Evaluation of diagnostic tests

Sensitivity and specificity are used to evaluate a test to identify the presence or absence of a disease in diagnosis (Altman and Bland 1994b). Sensitivity in a test quantifies the ability of the test to correctly identify those patients with the disease and specificity refers to the ability of the test to correctly identify those patients without the disease. Sensitivity and specificity are not directly affected by the changes in the prevalence of the disease in the study population but can differ with population (Egger et al. 1997; 1998). A test with high specificity

would unlikely be beneficial in decreasing mortality due to referrals at a stage that is too advanced (investigation at too high risk) likewise for an overly sensitive test that picks up those at low risk level unnecessarily exposing individuals to harmful investigations.

For any given diagnostic tests, there is always a trade-off between sensitivity and specificity, where a cut-off value or threshold is chosen for a particular test that optimises sensitivity in relations to specificity. The difficulty comes in deciding these trade-off or cutting points which also depends on the intended population. For example, chest X-ray is often the first line of diagnostic imaging tool used in primary care despite issues of being less sensitive, and less specific, failing to detect 77% of all cancers (79% of the cancers $\leq 20\text{mm}$, and 50% of lesions $> 20\text{mm}$) particularly at a surgically curable stage (Sone et al. 2000). However, pragmatically, chest X-rays are cheaper and hold lower radiation exposure risk.

2.7 Lung Cancer staging in relations to treatment and prognosis

The following section explains how once lung cancer is diagnosed, the extent of the disease is determined. Staging is the process of classifying the severity of the tumour; according to its size and growth or spread from the original tumour (Summerton 1999). Accurate staging in lung cancer is important as it determines the resultant management pathway and ultimately, informs the expected survival rate (prognosis). Over the years, advancements in lung cancer outcomes (imaging technology, treatment, biopsy techniques) have improved in accurately defining the tumour characteristics, and preventing people with lung cancer from unnecessary test procedures.

Stages could range from stage I through to stage IV. The TNM classification system is internationally used to determine the stages of lung cancer and is applicable to both NSCLC and SCLC. The 'T' corresponds to tumour characteristics such as size, location and invasion; 'N' denotes the extent of regional lymph node involvement and 'M' is the metastasis status. A complete version of the TNM classification system has been included in the Appendix 1 (AJCC 2002). Having that said, clinicians sometimes classify SCLC as either limited or extensive stage disease mainly because relatively small differences

in the tumour size had been observed to have little impact on the treatment response in SCLC which could be due to the rapid spread of SCLC .

Aside from its obvious use to confirm the presence of malignant cells in a tissue, biopsies for histological assessment are more limited. A systematic review comparing fine-needle aspiration biopsy against core-needle biopsy reported lower sensitivity and specificity for both biopsies in identifying histologic subtypes. The ranges of sensitivity and specificity for diagnostic purposes were 56% to 89% and 7% to 57%, respectively, compared to 81% to 97.4% (sensitivity) and 75% to 100% (specificity) for general diagnostic purposes (benign or malignant) (Yao et al. 2012). Also, beyond the classification of 'small cell' and 'non-small cell' cancer, the classification of primary NSCLC in biopsy sample is known to be difficult to achieve due to the absence of distinguishable diagnostic features, and well-documented variations that exist within these tumours (Roggli et al. 1985; Thomas et al. 1993; Edwards et al. 2000). There is also a general lack of histology for about 8% to 16% of some of the lung cancer, often classified as 'others', 'undifferentiated', or 'unclassified' (Thomas et al. 1993).

Generally, the lower the stage, the better is the prognosis. It has been widely believed that earlier diagnosis of lung cancer correlates to better prognosis in the 5-year survival rate. This concept has been widely debated as some studies viewed delays to negatively affect prognosis (O'Rourke and Edwards 2000; Jensen et al. 2002) while others formed more neutral standpoints (Salomaa et al. 2005; Billings and Wells 1996). One study found that 6 out of the 29 patients diagnosed with lung cancer (21%) deteriorated from potentially curable to incurable (inoperable) while waiting for treatment on the waiting list (O'Rourke and Edwards 2000). However, this was a relatively small study sample and also the severity of the disease in the early stage and the aggressiveness of the tumour should be considered. In another retrospective study based on the re-evaluation of 132 patient records, Salomaa et al. (2005) found no association between the time delay and the lung cancer stage. Delays were found to have less effect on prognosis in advanced cases of lung cancer when compared to early stage lung cancer (Salomaa et al. 2005). The choice of treatment of lung cancer (radiation, chemotherapy, surgery and/or combination of modality) is highly dependent on an individual's clinical stage.

Most studies agree that earlier detections increases operability rates which impacts survival rates for the disease (Salomaa et al. 1998; Strauss et al. 1997; Moody et al. 2004; Rogers 2006). There is still a debate regarding the centralisation of surgery in lung cancer. Early stage disease treated with radical surgery can increase the 6-month survival rate up to 85%, compared to the 45% observed in the non-surgical groups at early cancer stage (CRUK 2009). However, a recent retrospective cohort study described a vast amount of research exploring the different predictor variables (which included gender, age, pre-operative stage, type of lung surgery, type of cancer, co-morbidities, and lung function) of survival rate or lower complication rate following lung cancer surgery (Roth et al. 2008). Therefore, radical surgical intervention does not necessarily warrant longer lung cancer survival rate.

CHRONIC OBSTRUCTIVE PULMONARY DISEASE (COPD)

2.8 Pathophysiology

There had always been a wide variability in the definition of COPD and criteria for its diagnosis which makes it difficult to compare findings of COPD studies (Weiss et al. 2003).

Most recently, the World Health Organisation (WHO) has defined COPD as a lung disease characterised by chronic obstruction of lung airflow that interrupts normal breathing and can only be partially reversible. Previously independent defining terms, chronic bronchitis and emphysema, are now included in the diagnosis of COPD instead of being regarded as separate disorders (WHO 2011). The definition of COPD has since evolved into the co-existent of these two disease which leads to the heterogeneity observed in COPD symptom manifestations. The course of COPD will take one through phases of sudden deterioration known as acute exacerbation which is often triggered by an infection or noxious stimuli. Hence, generally, COPD is regarded as airflow limitation triggered by an abnormal inflammatory response and is thought to worsen with time (Edelman et al. 1992). The evaluation of spirometric parameters should be used to confirm diagnosis in the presence of suggestive symptomology (Weiss et al. 2003).

Symptomatically, COPD is characterised by both chronic bronchitis and emphysema which have different mechanisms of pathology (Edelman et al.

1992). Physiologically, the main characteristic anomalies that can be found in chronic bronchitis include:

- Hypertrophy of goblet cells and sub-mucosal glands of the airway
- Hyperplasia of goblet cells and sub-mucosal glands
- Infiltration of the mucosa with inflammatory mediators
- Hyperplasia of smooth muscle and
- Frequent inflammation of the respiratory bronchial

With inflammation, scarring and remodelling occurs which encourages thickening of the bronchial epithelium which inevitably constricts the airway and limits airflow (Edelman et al. 1992).

Emphysema is commonly defined anatomically by the permanent dilation of the distal airways to the terminal bronchioles as a result of the destruction of the walls of the alveolar (Viegi et al. 2007; Schwartz et al. 2009). This reduces the surface area of the walls which impedes effective gaseous exchange during respiration. This ultimately impacts the elasticity of the walls, and reduces the structural support of the airways which increases the likelihood of airway collapse which further restricts airflow (Edelman et al. 1992).

2.8.1 COPD in relations to lung cancer

Little is still understood about the relationship underlying COPD and lung cancer. Their disease mechanisms are characteristically poles apart where lung cancer is described as anti-apoptotic that involves uncontrolled cell proliferation and maintaining angiogenesis (growth of blood vessels), COPD attributes to apoptosis, destruction of the cellular matrix, impeding angiogenesis, and subsequent cell death. Existing theories speculate an inflammatory process, and remodelling phase of the lung architecture underlying COPD which could pre-dispose the development of lung cancer (Potton et al. 2009). The increase release of growth factors, and metalloproteinases from the matrix remodelling was said to encourage carcinogenesis (by cell proliferation) and malignant transformation of the bronchiole epithelium in a process called epithelial-mesenchymal transition (EMT) (Potton et al. 2009). In this process, an epithelial phenotype changes to a

mesenchymal phenotype (Dasari et al. 2006; Boyer et al. 2000). This was later supported by CT evidence showing a strong correlation between lung cancer and emphysema in particular even after adjusting for confounder such as smoking and other respiratory co-morbidities (de Torres et al. 2007).

The study hypothesised that normal lung homeostasis maintains a small turnover rate of cells with a few macrophage around the alveolar to remove any particles (de Torres et al. 2007). On exposure to cigarette smoking, inflammatory markers (neutrophils and macrophages) were thought to be further recruited and activated in the attempt to remove the foreign carcinogens. Emphysema occurs as a result of this inflammation which results in the degradation of the extracellular matrix and cell death. Bronchoalveolar stem cell (BASC) would attempt to repair these alveolar cells, and restore the enlarged airspace. However, over a chronic period of time, repeated stimulation of BASCs might have resulted in the uncontrolled proliferation of cells within a carcinogenic environment rendering the cells to be malignant. Thus, the result is bronchial carcinoma. This physiological effect is illustrated in Figure 2.4.

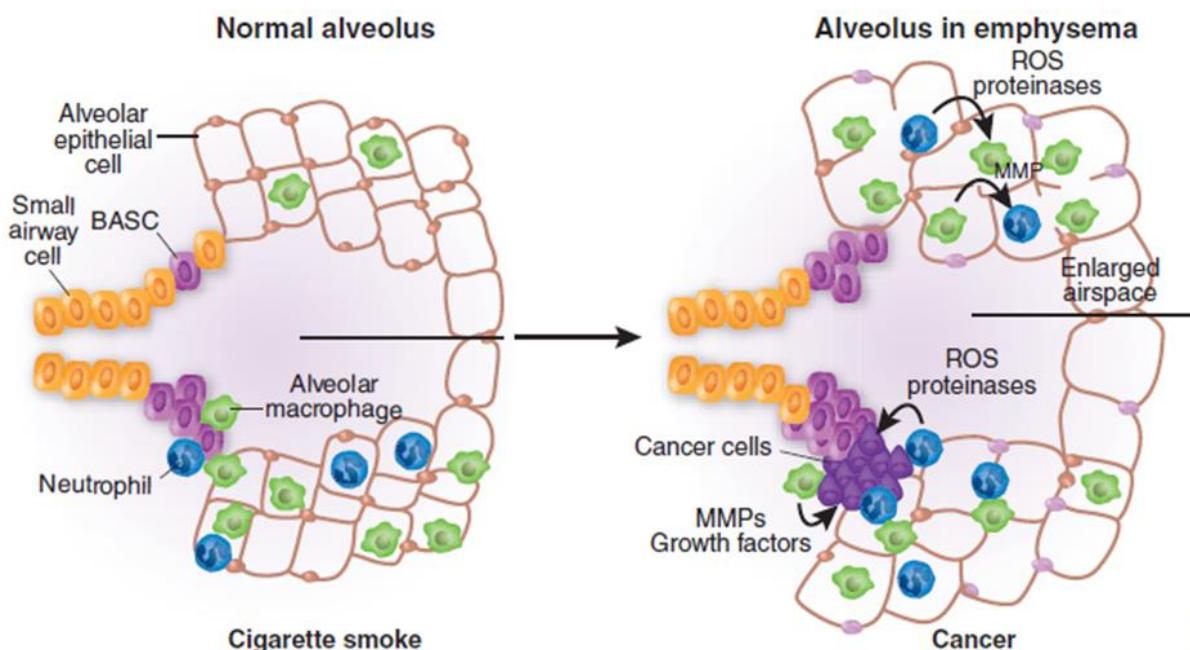


Figure 2.4 Illustration of the physiology behind inflammation in COPD and lung cancer (de Torres et al. 2007)

Similar mechanism behind lung carcinogenesis in association with COPD due to a stimulus (cigarette smoking) was proposed in a later review by Sekine et al. (2012) (shown in Figure 2.5).

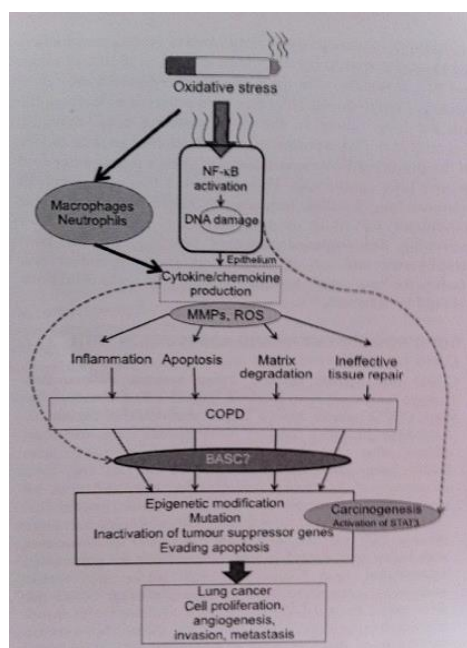


Figure 2.5 Mechanism of lung carcinogenesis following COPD (illustrations by Sekine et al. 2012)

Data from case control studies and cohort study were suggestive but inconclusive on the role of COPD as a definite risk factor in the development of lung cancer (Brownson et al. 1998; Petty et al. 1998; Mannino et al. 2003). However, most of the studies agreed to a causal relationship between the COPD and lung cancer to a certain extent (Brownson et al. 1998; Wu et al. 1995; Mannino et al. 2003). In a population based case-control study among women in Missouri (United States), Brownson et al. (1998) found consistently elevated lung cancer risk associated with history of COPD (odds ratios 2.1 to 2.7) to suggest an association. Biological evidence associating lung disease and lung cancer explained that damage to the lung from chronic lung diseases might have impaired the airway clearance ability and compromised immunity, resulting in an increased susceptibility to exposure of lung carcinogens such as cigarette smoking or other occupational exposure (Wu et al. 1995; Papi et al. 2004). This increased exposure is said to promote pathological changes leading to squamous cell neoplasia (Papi et al. 2004). In essence, the pathogenic role of COPD as a risk factor in the development of lung cancer is

not well defined as no one is sure of a model that fully explains the underlying mechanism relating the two pulmonary diseases other than the known phenomena of airway obstruction caused by muco-hypersecretion in response to a non-specific pulmonary inflammation. Hence, more research will still be needed to look at the issue at a molecular and practical level.

2.9 Epidemiology

Approximately 900,000 people in the UK are diagnosed with COPD (British Lung Foundation 2003; NICE 2004; NHS Choices 2010) but this prevalence of COPD was estimated to be much higher ranging from 3.4 to 3.7 million with 75% of the population either misdiagnosed or remained undiagnosed (Stang et al. 2000; Shahab 2006; NHS Choices 2010).

2.10 Aetiology and risk factors

Abundant research has suggested a strong causal association that exist between smoking and COPD (Edelman et al. 1992; Tockman et al. 1987; Mannino et al. 2002; Petty et al. 1997; 1998). There is sufficient epidemiology evidence to suggest a causal relationship that links COPD morbidity to active smoking (Tockman et al. 1987; Mannino et al. 2002; Viegi et al. 2007). The frequencies of respiratory symptoms increases with smoking history and status (pack years). A definite model for dose-response relationship between smoking frequency and disease severity has not yet been identified but generally, ex-smokers experience higher cumulative remissions of respiratory symptoms than persistent smokers (Eagan et al. 2004). Furthermore, there are studies to suggest that COPD due to environmental risk factors in developed Western countries was hardly ever severe enough to cause considerable obstruction to the lung in non-smokers (Buist 1988).

Other risk factors for the development of COPD include age, occupational exposures, air pollution, familial aggregation and genetic susceptibility to COPD.

2.11 Signs and symptoms

Clinically, symptoms of COPD are characterised as chronic productive cough, excessive sputum production, wheezing, and dyspnoea on exertion (laboured breathing) (Ferguson et al. 2000; Anthonisen et al. 2001; Rabe et al. 2007). Clinical features vary with each individual and are often under-reported which can make COPD diagnosis difficult. The presentation of symptoms also depends on the combination of risk factors an individual is exposed to (Weiss et al. 2003; Rabe et al. 2007).

The condition for chronic bronchitis in COPD is characterised as the presence of cough and secretions on most days for ≥ 3 months over 2 consecutive years (Siafakas 2006). Prevalence rates of chronic cough and/or secretion production in a population of Italy reportedly ranged from 14-44% and 6-17% in males and females, respectively (Viegi et al. 2007). Coughing is strongly related to smoking and changes in its incidence or remission are associated with changes in the smoking status. Haemoptysis can occur in severe stages of COPD especially during chest infections (Pauwels and Rabe 2004). Clinical literature has also suggested that chronic persistent cough often precede dyspnoea in COPD (Rabe et al. 2007). A brief trajectory of the symptoms of COPD from early onset COPD progressing into advanced COPD identified persistent cough in most of the participant experiences.

Dyspnoea is highly significant in COPD and is often the symptom that concerns people the most. It is usually persistent and progressive in COPD where a declined FEV1 of less than 30% predicted would result in breathlessness on minimal exertion (Viegi et al. 2007). These are generally symptom characteristics of COPD. Most studies would agree to the use of pulmonary spirometric testing to confirm the basis of COPD diagnosis.

Respiratory symptoms may exist for years before the development of airflow obstruction which is usually irreversible in COPD (Weiss et al. 2003). Airflow obstruction is diagnosed as a forced expiratory volume in 1 second $< 70\%$ or 0.70 of the predicted and forced expiratory volume in 1 second/ forced vital capacity ration $< 70\%$. Physical signs are usually more likely to present in severe COPD. They could reveal tachypnoea (rapid breathing), pursed lip breathing and increased work of breathing (on exertion). On auscultation, clinicians may

find decreased breath sounds which may relate to airflow limitation and possible inspiratory crackles. Wheezing during unforced tidal breathing is more specific to airflow limitation.

Some of the literature suggests that in the advanced stages of COPD, more systemic symptoms such as weight loss, anorexia, and/or psychiatric health issues (e.g. depression and anxiety) are common symptoms (Schols et al. 1993; Pauwels and Rabe 2004).

2.12 Diagnosis

The Global Initiative for Chronic Obstructive Lung Disease (GOLD) has produced guidelines to highlight the use of spirometry (gold standard) when possible as the first protocol to confirm the diagnosis of COPD where indicated (Kornmann et al. 2003). In the absence of spirometry, the GOLD guideline advises the assessment of symptom (symptom diagnosis) using symptom-based questionnaires to identify people at high risk of COPD, followed by routine monitoring for disease progression. Symptom indicators are not proven to be diagnostic on their own, but the presence of multiple symptom indicators, and a history of exposure to risk factors (e.g. environmental pollution, cigarette smoke, occupational dusts) is said to increase the probability of COPD. Another diagnostic tool to detect COPD is the use of sputum cytology. Samples are analysed and interpreted by a sputum cytopathologists (Petty et al. 1996).

2.12.1 Use of symptoms-based questionnaire tools

Functional spirometry has been proposed to be fairly effective in detection early COPD. However, like all screening programs, there were concerns about its practicality as it is not necessarily available in all primary care clinics, and inevitably the cost has to be considered. Following this, attempts were made to identify specific individuals at higher risk of having abnormal spirometric findings such as the Global Initiative for Chronic Obstructive Lung Disease (GOLD) classification (GOLD 2004) (see Table 2.2). However, reviews have questioned the validity of this classification with regards to the usefulness of the new addition of a stage 0 (Vestbo and Lange 2002). GOLD stage 0 represents those symptomatic with no evidence of airflow limitation, and is

intended to include those 'at risk' of developing airways obstruction in the future. Vestbo and Lange (2002) found that stage 0 was inappropriate for its intended purpose in a randomly selected general population particularly among smokers. Figure 2.6 presents a questionnaire designed by GOLD to detect people with COPD (GOLD 2004).

Could It Be COPD

Chronic obstructive pulmonary disease (COPD) is a major cause of illness, yet many people have it and don't know it.

If you answer these questions, it will help you find out if you could have COPD.

This short interactive questionnaire was developed by GOLD and has been scientifically evaluated and shown to identify people who are more likely to have COPD.

Questionnaire

Question 1 Do you cough several times most days?	Yes <input type="radio"/> No <input type="radio"/>
Question 2 Do you bring up phlegm or mucus most days?	Yes <input type="radio"/> No <input type="radio"/>
Question 3 Do you get out of breath more easily than others your age?	Yes <input type="radio"/> No <input type="radio"/>
Question 4 Are you older than 40 years?	Yes <input type="radio"/> No <input type="radio"/>
Question 5 Are you a current smoker or an ex-smoker?	Yes <input type="radio"/> No <input type="radio"/>

Figure 2.6 "Could it be COPD" questionnaire (GOLD 2004)

2.13 Severity of COPD

Peak flow alone is said to have low specificity for COPD diagnosis as decreased peak expiratory flow could also suggest other pulmonary diseases. The severity of COPD is classified as at risk, mild, moderate, and severe according to the GOLD (2001). Spirometric testing can determine the diagnosis of COPD and is also used in the staging of disease severity. Table 2.2 explains the

Literature Review

GOLD's classification of COPD severity based on findings of the FEV₁ and FEV₁/ forced vital capacity (FVC) ratios.

Table 2.2 Classification of COPD severity defined (GOLD 2001)

Stage	Classification	Characteristics
0	At risk	Presence of chronic respiratory symptoms, smoking history and normal spirometry
1	Mild	FEV ₁ ≥ 80% predicted FEV ₁ /FVC < 70%
2	Moderate	FEV ₁ :30-80% predicted FEV ₁ /FVC < 70%
3	Severe	FEV ₁ < 30% predicted FEV ₁ /FVC < 70%

2.14 Early detection and screening for Lung Cancer in COPD

With studies on chronic respiratory diseases as risks for lung cancer in both smokers and non-smokers showing a general increased risk of 30% to 60%, an association between COPD and lung cancer is plausible (Koshiol et al. 2009; Kennedy et al. 1996; Diez-Herranz 2001; Wasswu-Kintu et al. 2005). The Environment And Genetics in Lung cancer Etiology (EAGLE) population-based case control study found the risk of lung cancer to be more than two times higher in the COPD population (odds ratio 2.5) even after adjusting for confounders (smoking, demographic, and socioeconomic variables) (Koshiol et al. 2009). Some studies attributed this association to smoking, an explicit confounder linking COPD and lung cancer (Alavanja et al. 1992; Alberg 2003) whilst others refuted this theory in view of findings of strong associations between COPD and lung cancer in never smokers (Mayne et al. 1999; Brownson and Alavanja 2000; Turner et al. 2007). Having stratified for smoking status, Brownson and Alavanja (2000) found negligible differences in lung cancer risk due to emphysema (OR 2.7) and chronic bronchitis (OR 1.7). These findings were consistent with Mayne et al. (1999); emphysema and chronic bronchitis

were significantly associated with increased lung cancer risk in men and women non-smokers. Sun et al. (2007) found that adenocarcinoma, less commonly found in smokers, was still strongly associated with COPD which also supported the hypothesis. More recently, evidence of strong associations found chronic bronchitis, and lung cancer to be strongest among lighter smokers in the EAGLE study which suggested a common molecular feature that is independent of smoking (Koshiol et al. 2009). Koshiol et al. (2009) suggested a possible inflammatory response as a result of an infection which led to COPD exacerbation and rapid decline of lung function but the never-smoked population of the study was too small to offer conclusive evidence. Nevertheless, there are studies that have found exposures to risks such as smoking and occupational asbestos; which often associates with airway inflammation and COPD, can be linked to lung cancer (Barnes et al. 2003; Brody and Spira 2006).

Given what is known about lung cancer, the concept of 'effective screening' might be employed in the effort to diagnose earlier lung cancer. This differs from conventional systematic testing of asymptomatic individuals to identify those at potential risk of a specific disease. Rather, it identifies a selection of individuals who present with an increased index of suspicion that justifies their screening. For instance, fewer cases of lung cancer were found in the screening of a large cohort of young smokers and non-smokers (Sone et al. 1998) when compared to the screening of older committed smokers as seen in the study of Henschke et al. (1999). Albeit both studies had different sample sizes, and are not directly comparable, the underlying concept assumes that the appropriate choice of population screened will markedly ensure cost-effectiveness where economic factor is a major influence.

One study evaluated the use of a questionnaire tool to screen for patients at high risk of lung cancer in groups with (COPD) and without airflow obstruction over a 5-year period (Bechtel et al. 2009). 430 patients who attended the primary care clinics were identified as high-risk. High-risk was characterised as being above the age of 50 years; either smoked or had previously smoked, exposed to asbestos, and/or had a family history of otolaryngology-related cancer. Eight out of the 126 high-risk patients with airflow obstruction were eventually diagnosed with lung cancer; 50% of which had early staged lung cancer and survived 5 years. The high-risk group without airflow obstruction

Literature Review

(n= 304) were not screened further had 10 lung cancer cases during follow-up; mostly at the advanced stage of the disease (60%). This suggested that lung cancer was found in a higher percentage of patients with airflow obstruction (6.8%) compared to those with normal spirometry (3.1%) which concurred with findings of previous studies (Wilson et al. 2008; Skillrud et al. 1986). Wilson et al. (2008) found independent statistically significant association between emphysema and lung cancer with odds ratio of 3.14 (confidence intervals 1.91-5.15); after adjusting for gender, smoking exposure, and additional adjustments for GOLD class (mild, moderate and severe COPD).

Although on face value the results suggested that lung cancer has a higher prevalence in people with COPD, the findings of Bechtel et al. (2009) could be skewed as not every participant was formally tested. Like most cohort study, missing data resulting from participants lost to follow-up and non-cancer related deaths during the 5-year period might introduce selection bias, and limit the evidence for recommendations for targeted lung cancer screening following risk assessment using a questionnaire and spirometry.

More recently, there had been more evidence to suggest an association between lung function and the increased risk of developing lung cancer (Van den Eeden and Friedman 1992; Mannino et al. 2003; Young et al. 2007; Calabro et al. 2010). Calabro et al. (2010) compared people with decreased FEV1 (<90%) with people with healthy lung functions (FEV1 \geq 90%) and found that people with even a slight reduction FEV1 resulted in an increased likelihood of lung cancer (odds ratios 2.4) having adjusted for age, sex, and smoking variables. Therefore, lung function was a significant predictor of increased lung cancer risk. This study was much larger (n=3,869) with higher numbers of lung cancer cases in comparison to the other studies hence, well-powered (Calabro et al. 2010). It is known that lung function deterioration is an eventual physical manifestation of deteriorating COPD. This added risk of lung cancer within a population that is known to receive late diagnosis in lung cancer warrants the need for studies to identify ways to improve the survival rate of lung cancer in those with chronic respiratory diseases.

SYMPTOMATIC DIAGNOSIS OF LUNG CANCER

Despite lung cancer being one of the largest cause of cancer death, it is very under-researched compared to other cancers. There is still not a lot is known

about symptomatic diagnosis of lung cancer particularly within sub-population with specific comorbidities such as COPD.

The annual incidence rates of lung cancer were reportedly four times higher in those with COPD compared to the general population (Kiri et al. 2010). With COPD preceding lung cancer diagnosis in 40% to 90% of the cases (Young et al. 2009), many individuals with COPD in primary care experienced delays in the diagnosis of lung cancer, which contributed to advanced stage of the disease at diagnosis, and poorer survival rates (Kiri et al. 2010). Existing impairments to the lung as a result of the COPD further limits their chances of curative surgery, which might explain the poorer 3-year survival for lung cancer in those with prior COPD. Patients with COPD and lung cancer had almost 50% lower survival rates than those without COPD (15% vs 26%; $p < 0.01$). These primary care findings were collected using the UK General Practice Research Database (GPRD) (Kiri et al. 2010).

Studies have not specifically explored reasons for the delays in lung cancer diagnosis in the COPD population but one of the reasons proposed by the British Lung Foundation was that some of the lung cancer symptoms are also common to COPD. People who have underlying respiratory diseases such as COPD frequently experience cough as a symptom but it is often the change in their cough from their “usual” cough that indicates lung cancer diagnosis (Hamilton and Sharp 2004). Thus, changes in pre-existing coughs should be highlighted. One study published in 1977; recruiting 6137 patients, found prolonged cough as a symptom had low PPV (0.03) to predict lung cancer (Boucot et al. 1977). Another prospective study in Germany followed 329 consecutive individuals with chronic persistent cough over a period of two years did not find any lung cancer cases (Kardos and Gebhardt 1996). Bjerager et al. (2006) looked at the possible reasons for the delay in medical referral of symptomatic patients, identified that symptoms were often attributed to the co-morbidity e.g. COPD rather than potential lung cancer in patients with the pre-existing disease. However, findings in this study represent those of a relatively small sample population ($n=84$), which had also limited the type of statistical analyses allowed (as a result stratified analyses were not possible). Furthermore, their data from interviews with GPs were retrospective, which could be subjected to recording and recall bias. Then again, the population-

Literature Review

based design of the study could have improved the generalisability of the results (Bjerager et al. 2006).

A study on lung cancer presentation found that while haemoptysis generally associates with prompt referral to a specialist, more common symptoms such as cough often obscures the process of getting appropriate medical referral as they tended to be neglected for a while, to consider other differential diagnosis (Buccheri and Ferrigno 2004). If this is true for those without chronic respiratory problems, we can reasonably infer that it would only be more difficult to suspect lung cancer in patients with COPD (a disease that commonly presents with cough) and refer them to the appropriate services. Hence, the American College of Chest Physicians (ACCP) practice guidelines recommends that patients with normal chest imaging with a history of COPD and smoking, presenting with haemoptysis should still be closely monitored for suspected lung cancer to enhance timely diagnosis and care (Spiro et al. 2007). However, there is still insufficient evidence that examines the reason for a delay in reporting lung cancer symptoms.

Recommendations can be found in the British National Institute for Clinical Excellence (NICE) guidelines regarding symptomatic chest X-ray referral (NICE 2005). A recent NICE guideline (2011) acknowledges that the evidence on symptoms has not been updated and reviewed since 2005. Furthermore, most of these recommendations are based on grade D studies; evidence based on non-empirical studies and formal expert opinions, and some extrapolations from case control studies, and/or cohort studies. Evidence on lung cancer symptoms in the NICE guideline was largely based on a non-empirical study by Beckles et al. (2003). The current knowledge-base is mostly on symptoms at diagnosis and post-diagnosis (measuring treatment effect).

This might have been again due to the conventional theory that there is a lack of clinical presentation (asymptomatic) in the early stages of neoplastic pulmonary growths (Mason 1041; Spiro and Silvestri 2005). However, recent studies have disputed this view and indicated that many of the lung cancer patients can in fact experience symptoms for many months prior to diagnosis, with reported delays up to 2 years regardless of their disease stage at diagnosis (Corner et al. 2005; Buccheri and Ferrigno 2004; Tod et al. 2007). Furthermore, majority of the patients (80-90%) in the UK have been diagnosed

following symptomatic presentation in primary care (Hamilton and Peters 2007; Yoder 2006). These findings have informed the current national guidelines, and mark the beginning of major public awareness programs (NICE 2014). This included the government's National Awareness and Early Diagnosis Initiative (NAEDI), a key component of the 2007 cancer reform strategy, aimed at early symptomatic diagnosis, and promoting symptom recognition (Department of Health [DoH] 2007; CRUK 2009).

In order to determine the diagnostic values of lung cancer symptoms, predictive values, and likelihood ratios of symptoms for lung cancer diagnosis are needed (Altman and Bland 1994a). Predictive value of the symptom is the proportion that is clinically useful as it informs clinician the likelihood of a patient having lung cancer given that the individual has a particular symptom(s) but predictive values are population-dependent (cannot be generalised to other populations). Sensitivity represents how well the clinical test correctly identifies all patients with the disease, and specificity refers to the ability of the test to correctly identify those without the disease. Although sensitivity and specificity are not dependent on the population, they can however, be affected by the spectrum of disease (Lalkhen and McCluskey 2008).

A population based case control study by Hamilton et al. (2005) was the first to advance the evidence base for diagnostic values of symptom in early detection of lung cancer. The study calculated the likelihood ratios for symptoms from data on 247 eligible lung cancer patients against 1235 controls matched for age, sex and general practice. Positive predictive values of the symptom variables were identified.

With the study being adequately power (85%), findings showed that symptoms of haemoptysis, dyspnoea and abnormal spirometry were calculated to pose high risk for lung cancer. Haemoptysis had the highest PPV value of 2.4%, which supports current guidelines recommending urgent chest X-rays for patients with haemoptysis. However, haemoptysis as a symptom was relatively uncommon; only reported in 20% of the cases, compared to the other low risk symptoms (PPV <1%). Therefore, this means that primary care clinicians must still make their referral decisions in patients with symptoms weakly associated with lung cancer. A combination of more than one symptom presentation

Literature Review

generally increases the PPV for lung cancer and therefore, the risk of lung cancer. Similarly, re-attendance or persistent cough symptom was found to increase the PPV value. This suggests that a set of variables will predict better than individual signs and symptoms. The final model included a combination of symptoms gathered from the multivariate analysis which forms an initial guide to clinicians, and guidelines developers on the selection of high risk patients for rapid referral.

As the first study to examine all the pre-diagnostic features of lung cancer in a primary care setting, Hamilton et al.'s study (2005) was undoubtedly imperative and helped to strengthen the weak evidence base that previously supported guidelines. However, the study also had its limitations. Their findings were based on retrospective records of GPs which may be subjected to recording bias due to the potentially wide variance across the different GP practices. There is also a risk of over-inflated PPVs, where symptoms were recorded more carefully if lung cancer was already suspected, associated with retrospective data collection. Detailed report of the PPVs of symptoms indicative of lung cancer and the methodological issues of previous studies can be found in the systematic review included in this thesis (Chapter Three).

There has been much research into better understanding lung cancer in hope of a breakthrough in this prevailing endemic. Evidence already suggests that better survival outcomes are associated with earlier detection of lung cancer which can be symptomatic in the early stages (Richards 2007). As the EURO CARE report has shown, the reported differences between the European countries are significant as studies have indicated that the lower survival rates in the UK were attributable to late staged lung cancer at time of diagnosis in comparison to its European counterparts (with consideration of factors such as effects of age, and national expenditure). Improving the evidence base for earlier symptomatic diagnosis is a priority if the proportion of patients diagnosed with operable disease is to increase, and lung cancer outcomes are to be improved in the UK.

Chapter 3: Systematic Review

Symptomatic diagnosis of Lung Cancer

3.1 Background

This systematic review was carried out to gather sources of the available evidence in lung cancer symptoms, and provide an update of the literature on symptomatic diagnosis of lung cancer. Parts of this work have since been published in *Family Practice* and can be found in Appendix 2 (reprinted with permission from Oxford Journals). In the UK, lung cancer (lung cancer) is still the second most common cancer, accounting for 13% of all new cases of cancer (CRUK cancer statistics report, 2009). Over the last four decades, survival rates have only improved slightly with most lung cancer being diagnosed at late stages when curative intervention is no longer viable in both developing and developed regions of the world (GLOBOCAN 2008; Janssen-Heijnen and Coebergh 2003). Evidence from large population-based studies has since associated relatively lower survival in some countries with delays in diagnosis (Richards 2007; Walters et al. 2013).

The current UK NICE guidelines recommend urgent chest X-ray referrals for patients experiencing any persistent symptoms that might indicate lung cancer. Ideally, where prospective diagnostic evidence in primary care is available, it is used to support guidelines in cancer referrals. However, such evidence base is not available for lung cancer. Ideally, recommendations and guidelines would prefer to apply the highest level of evidence, level-A evidence comprising of systematic review and/or meta-analysis of randomised controlled trials (RCTs) where possible.

Most of the symptom-based research studies found in lung cancer were either retrospective cohort studies or case-control studies that had relied on clinical records, and databases. These studies are limited by methodological issues regarding symptom data collection; including the possibility of recording bias by clinicians, which might have implications for the predictive values obtained. Cases of newly diagnosed lung cancer had been found to be symptomatic at the time of their diagnosis when presented in primary care (Fergusson et al. 1996).

Systematic Review

In the absence of experimental RCTs, a meta-analysis of well-executed diagnostic studies (to include studies that have been recruited from a comparable patient population and has considered important potential biases) offer stronger evidence for further inferences to be made compared to a meta-analysis of observational studies (Deeks 2001). Studies of diagnostic accuracy describes the relationship between test results, and disease using probabilistic measures such as PPVs, diagnostic odds ratios, and likelihood ratios which provide useful evidence of the value of a diagnostic tool (Deeks 2001). Essentially, the predictive value of a test may be of highest significance in deciding whether or not to refer a patient in primary care as it represents the probability of a serious malignancy in a person with a symptom or test result (Hamilton 2009).

A thorough search of the available literature described in the narrative review in Chapter Two, has shown very little evidence could be found on diagnostic studies relating to symptoms in lung cancer. Two systematic reviews have addressed the diagnostic value of symptoms, in lung cancer; Hamilton and Sharp (2004) and Shapley et al. (2010). Hamilton and Sharp (2004) reviewed features of symptomatic lung cancer across studies, and estimated the likelihood ratios of some of the symptoms in lung cancer diagnosis. The estimated likelihood ratios reported were based on referred populations in secondary care settings with hardly any research to be found in primary care populations. The study of Shapley et al. (2010) identified symptoms, signs, and non-diagnostic test results that were highly predictive of specific cancers (where positive predictive values (PPVs) $\geq 5\%$ were reported). The review analysed all higher quality evidence of symptoms that predicted lung cancer in an unselected primary care population and reported two studies for lung cancer. Only haemoptysis was identified as having high PPVs ($\geq 5\%$) in lung cancer diagnosis.

Due to the lack of high quality research in primary care populations in the most recent systematic review identified, this systematic review will investigate the diagnostic value of symptoms for lung cancer regardless of national health care system or spectrum of disease, and identify any new primary care evidence since 2010. It is acknowledged that these estimated diagnostic values of the symptoms reported in this review are limited in its generalisability to any

other populations. The review also included qualitative studies to explore the symptom experience of people diagnosed with lung cancer, and identify factors associated with patient reporting of symptoms that might have implications for the design of future diagnostic studies. This qualitative component could also reveal any non-classical symptoms or characteristics of symptoms experienced before lung cancer diagnosis not investigated in diagnostic studies.

3.2 Aims

- (1) To critically evaluate and summarise existing evidence providing the diagnostic values of symptoms for lung cancer.
- (2) To provide an overview of the characteristics of pre-diagnostic symptom (symptoms reported before diagnosis) characteristics reported in qualitative lung cancer studies.

3.3 Methods

3.3.1 Search strategy

Electronic databases were searched from their commencement to July 2012 using search terms 'lung cancer*', and 'sympto*' for title abstracts. The key terms were exploded to include alternative MeSH descriptors such as 'Lung Neoplasms', 'Signs and Symptoms', and 'Differential Diagnosis'. Search hits were filtered using qualifying restrictions: diagnosis (DI), epidemiology (EP), and etiology (ET) for 'Lung Cancer', 'Symptoms' and 'Differential, Diagnosis'. Details of the complete search strategies of all the databases performed are included in Appendix 3 (supplementary data).

Similar search methods were applied for all the electronic databases listed in Table 3.1.

Table 3.1 List of electronic databases searched

Electronic databases
<ul style="list-style-type: none"> • MEDLINE (from 1946 to July Week 1 2012) • Embase (1946 to Week 28 2012) • Cumulative Index to Nursing and Allied Health Literature (CINAHL) (1981 to 16th July 2012) • Multi- database: Embase (1980 to 2012 Week 2); Ovid MEDLINE (1946 to July Week 1 2012); Ovid MEDLINE Daily Update (July 13 2012); Ovid MEDLINE In-Process & Other Non-Indexed Citations (July 13 2012) • Cochrane Library

In addition, the contents pages of four journals (two for quantitative and two for qualitative studies) between 1st January 2009 and 31st December 2011 were hand-searched: Thorax and the British Journal of General Practice (BJGenPrac) for quantitative studies, and the Psychooncology and the European Journal of Cancer Care (EJCC) for qualitative studies. 3041 duplicates were removed using EndNote reference manager. Based on the search strategy, journals with the highest number of relevant papers for quantitative and qualitative studies were selected. This generated a total of 3017 papers (1830 quantitative and 1187 qualitative). The final update was performed in July 2012 with a total of 9054 articles to be assessed for inclusion. All records were retrieved and screened for relevance.

3.3.2 Inclusion and exclusion criteria

The following criteria were used to determine the eligibility of studies for this review:

- Quantitative study design- Studies that reported diagnostic values (positive predictive values (PPVs), hazard ratios, odds ratios, and/or likelihood ratios) for the symptom, sign, or test, or provide the necessary information needed to calculate these values (2x2 contingency tables could be reconstructed).
- Qualitative study design- Studies that explored the trajectory of symptoms from when symptoms were first experienced before diagnosis, or studies that described the onset of symptoms or first symptoms at

presentation to clinician (primary care) were of interest for the purposes of the review.

- Participants- Only adult populations recruited from hospitals, outpatient clinic, specialist clinic, specific community or the general population.
- Outcomes- The group with the positive outcome (lung cancer) must have had a confirmed diagnosis of lung cancer that met diagnostic standards set by the health service provider.
- Other- Studies written in a language other than English, German, Spanish, Malay and Chinese were excluded. Studies on multi-site cancers were included provided that lung cancer was distinguished from other cancers in reporting of results.

Exclusion criteria for the review were:

- Study design- Studies that reported symptoms post-treatment were excluded. These included studies on the management of symptoms in advanced lung cancer, studies measuring the effect of toxicity, and quality of life studies on symptom burden where baseline reported only post-treatment symptoms which will not provide diagnostic values. Single case studies, case reports, editorials, symposiums, reviews (literature), practical guidelines were excluded.
- Participants- Studies that reported symptoms of metastatic cancer, where lung cancer is the secondary cancer were also excluded.

3.3.3 Study selection and quality assessment

The initial screening of the titles and abstracts was carried out independently by the first reviewer (JS). A second and third reviewer (LB and MS) each checked a random sample of 75 (2.5%) of the abstracts. All papers shortlisted were retrieved in full. A second reviewer (LB) checked 100%, and a third reviewer (MS) 25%, of the full papers that were shortlisted to ensure that they met the eligibility criteria. Methodological quality was assessed using the Scottish intercollegiate Guidelines Network (SIGN) checklist for cohort studies and case-control studies, and the Consolidated Criteria for Reporting Qualitative research (COREQ) checklist was used to assess qualitative studies. In addition

to the usual criteria for reporting standards, the SIGN checklist for the quantitative diagnostic studies included items on disease prevalence, characteristic of the reference standard, and values of sensitivity, specificity, predictive values and likelihood ratios to evaluate the diagnostic accuracy of studies, which was relevant to the purpose of this systematic review. Other quality assessment tools such as the STARD (Standards for Reporting of Diagnostic Accuracy) statement and Newcastle-Ottawa Quality Assessment (NOQA) Scale for case control studies were considered. However, these tools were found to have less emphasis on diagnostic accuracy in comparison to the SIGN checklist. For example, issues around confounding, usually applicable to diagnostic test studies were not assessed in the STARD checklist and NOQA scale.

The COREQ checklist for qualitative studies was developed to promote explicit and transparent reporting of interviews; the most common method of data collection used in the qualitative studies included in this review. The 32-item checklist in COREQ had been established to report important aspects of the research team, study methods, context of the study, findings, analysis and interpretations in qualitative studies (Tong et al 2007).

The study did not use the rating or scoring system to the items on both checklists due to the wide variability across the studies; in concept and methodology. Internal validity and quality of study was based on the item responses, and considered accordingly to the individual study to determine the level of evidence offered. Effectively, the quality standard of the study was decided by the reviewers collectively based on methodological variations, and any potential biases that may affect the findings. Disagreements or uncertainties about satisfaction of quality criteria were discussed with the second and third reviewers (LB and MS) and consensus achieved.

3.3.4 Data extraction and analysis

The reviewers requested raw data of potentially relevant data from main authors to ensure a comprehensive inclusion of existing literature. Data on the type of study, characteristics of the study population, duration of follow-up and the effect sizes were extracted systematically and tabulated for each study that

met inclusion criteria. If diagnostic values were not reported, the positive and negative likelihood ratios (PLR and NLR), sensitivity and specificity were calculated using the 2x2 contingency tables (Altman and Bland 1994b). PPVs were also calculated where possible.

There are many sources of heterogeneity that are important in a systematic review, such as clinical heterogeneity; where patient populations are not the same, and/or methodologic heterogeneity; where studies are conducted differently (Garg et al. 2008). Across the five quantitative studies included in this review, the research methodologies varied in methods for measuring symptoms, population sample, and spectrum of disease to address different research aims. Although the diagnostic value presented (ORs and PPVs) (see Table 3.6) may appear similar, these values are dependent upon symptom measurement methods and, particularly the PPVs, are dependent on the population, and therefore, diagnostically it would not be meaningful to pool the data from primary and secondary care studies. Furthermore, a standard measure of effect sizes (e.g. odds ratio, risk ratio, or hazard ratio) across the studies is required in a meta-analysis, which was not available in the current review; for example, there were no odds ratio of one symptom in all five studies (Garg et al. 2008; Thompson 1994).

Moreover, it has been suggested that where there are reasonable clinical and methodological variability across a small number of studies (<5 studies), it might not be appropriate for the outcome data to be pooled (Garg et al. 2008). Limited number of studies can be a problem in a meta-analysis as the overall pooled N (sample size) might be small, which could lead to wide confidence intervals or imprecision (Guolo and Varin 2015). Such instances can cause the pooled estimate to cross the null hypothesis, making it difficult to draw a conclusion in either direction. For these reasons, the quantitative data were not suitable to be mathematically combined or pooled to present in a meta-analysis.

The study used narrative summaries (narrative review approach) to synthesise the qualitative data. As the field of mixed-methods systematic reviews is still in its infancy, there is currently no consensus on how such reviews should be formally analysed or how synthesis of data should be conducted (The Joanna Briggs Institute 2014). There are only few published reviews that can be

Systematic Review

considered mixed-methods in that included data are combined in a single synthesis or united in a secondary “final” synthesis (The Joanna Briggs Institute 2014; Harden 2010).

A primary criterion for the development of an integrated mixed-method systematic review is that both quantitative and qualitative data have to be similar enough to be assimilated into a single synthesis (The Joanna Briggs Institute 2014; Harden 2010). Only then both types of data can be assimilated into a single synthesis, whereby 1) the quantitative data is converted into themes, codified, and then presented along with qualitative data in a meta-aggregation, or 2) qualitative data is converted into numerical format, and included with quantitative data in a statistical analysis to collate both data (The Joanna Briggs Institute 2014).

Where the current review was not able to formally use any of the integrative methods of analysis due to the diversity within the included studies, qualitative and quantitative data were analysed separately, and then the “total” results discussed in a narrative discussion (Bruinsma et al. 2012).

3.4 Results

6,037 papers were retrieved using the search strategy. Result of the search strategy and selection process is shown in the flow diagram below (Figure 3.1). Duplicates were removed using EndNote reference manager. The final update was performed in July 2012. In total, 9054 articles (including the 3017 hand-searched journals) were assessed for relevance. Out of which, 11 studies (5 quantitative and 6 qualitative) were eligible for inclusion in the final review.

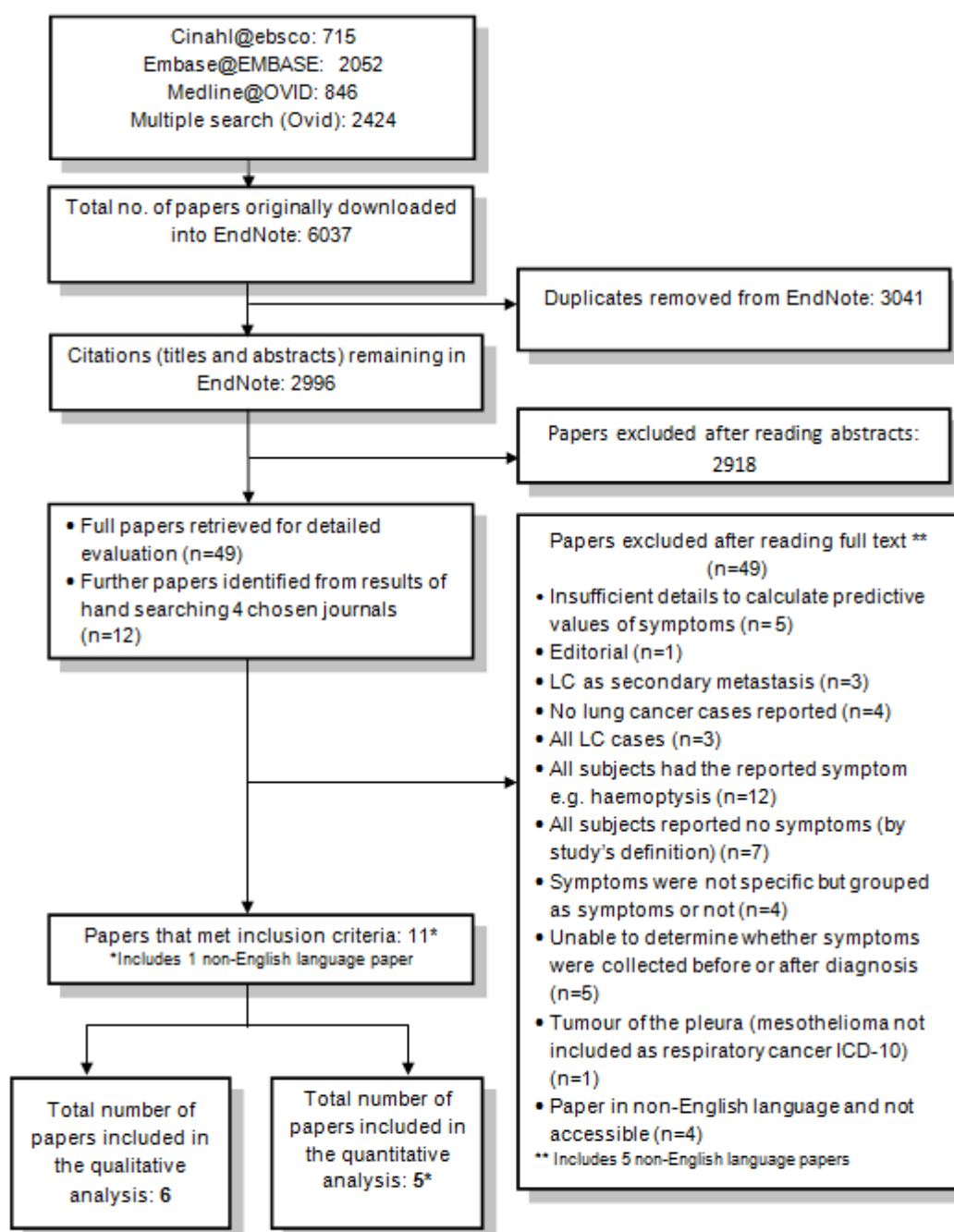


Figure 3.1 Flow diagram of results of search strategy adopted from the QUOROM statement flow diagram (Moher et al. 1999)

3.5 Results of Quantitative studies

Table 3.2 and Table 3.3 summarises the main characteristics of the quantitative studies to provide an overview of the studies. The duration of symptom onset to diagnosis generally ranged from 6 months to no more than 2 years before

diagnosis. Two cohort studies followed up the lung cancer diagnosis at one year period intervals (Jones et al. 2007; Hippenley-Cox and Coupland 2011). Four studies obtained their diagnostic values of symptoms from GP records and/or general practice databases (QResearch database) (Hoppe et al. 1977; Hamilton et al. 2005; Jones et al. 2007; Hippenley-Cox and Coupland 2011) whilst the remaining study obtained its symptoms data directly from the participants using a form of standardised questionnaire e.g. Medical Research Council Respiratory Questionnaire (Kubik et al. 2001).

3.5.1 Symptomatic prevalence

In general, cough, dyspnoea and haemoptysis were the most commonly measured symptoms, based on patient-reported symptoms, and those recorded by GP (shown in Table 3.4). In the four studies that reported symptom frequencies, systemic symptoms such as appetite loss, weight loss, fatigue, and fever/flu were less frequently reported (Hoppe 1977; Hamilton et al. 2005; Hippenley-Cox and Coupland 2011). However, symptom frequencies will vary depending on the method of symptom reporting (directly by patient or recorded by GP), and the characteristics of the lung cancer population (early or late-staged lung cancer), and the control groups (healthy user effect).

Table 3.4 and Table 3.5 show symptoms that were recorded in the study and which of these predicted lung cancer ($p \leq 0.05$). Hamilton et al. (2005) excluded some of the symptoms associated with lung cancer in the analysis of their study which included 'hoarseness' and 'shoulder pain from Pancoast tumour' (refer Table 3.5). The study reasoned that they had found these selected presentations to be rare and possibly only a few of these presentations were observed in the lung cancer cohort. They also added that the symptom hoarseness was associated with features of advanced late-stage lung cancer (Hamilton et al. 2005).

3.5.2 Symptomatic diagnosis

The diagnostic values of pre-diagnostic symptoms for lung cancer were calculated from all five studies where possible (using 2x2 contingency tables).

Table 3.6 present the PPVs and odds ratios of the individual symptoms reported in the case-control studies and cohort studies, respectively. In two of the studies, p-values or confidence intervals were not provided (Hoppe 1977; Kubik et al. 2001). Most symptoms reported in the study of Hamilton et al. (2005) had reasonably high odds ratios (OR 4.40 to 16.24). However, even with the high odds ratios, the likelihood of a patient presenting with that symptom having lung cancer (as indicated by the PPV) was low (<1.0 for most symptoms).

PPVs of 2.0 and higher ($2.4 \leq \text{PPV} \leq 7.5$) were consistently found for the symptom haemoptysis across the studies. Jones et al. (2007) evaluated the diagnostic value of haemoptysis in lung cancer, reported an increase in PPV from 5.8; at six months after symptom onset, to 7.5; three years after the first symptom occurrence was recorded in the male group. Similar increase was observed in the female cohort: 3.3 to 4.3. A gender difference in PPV was also identified (Jones et al. 2007). However, no statistically significant differences in gender-specific PPVs were identified in Hippenley-Cox and Coupland's study (2011); the remaining three studies did not carry out separate analysis for men and women (Hoppe 1977; Kubik et al. 2001; Hamilton et al. 2005).

Hippenley-Cox and Coupland (2011) calculated the hazard ratios (HRs) of individual symptoms in the derivative cohort to develop a model strategy that predicts lung cancer risk in a population. Only variables of HR <0.80 or >1.20 were included into the final model. The study also observed the highest haemoptysis had a higher hazard ratio in comparison to the other symptoms for the final model for lung cancer in both genders (Male HR: 21.5 and Female HR: 23.9), followed by weight loss, appetite loss, and cough (symptom variables included in the final model). In their final model, the top 0.5% of their risk score produced a PPV of 9.5 (8.8 to 10.3).

The remaining two studies did not identify any statistically significant associations of symptoms with lung cancer (Hoppe 1977; Kubik et al. 2001).

Odds ratio for each symptom was calculated from the data obtained from Hoppe's study (1977). Symptoms; flu and general wellbeing, had the lowest odds ratio in this study. This could be explained by the low specificity of these symptoms; could also be symptoms of a wide range of other ailments, and less likely to be related to lung cancer. However, the review will not be able to know

Systematic Review

for certain, as we do not know the population frequencies for all the symptoms. Chronic bronchitis had the highest odds ratio of 1.51 which can be said to concur with findings of the many studies establishing a positive association between COPD and increased lung cancer risk (Koshiol et al. 2009; Kennedy et al. 1996; Diez-Herranz 2001; Wasswu-Kintu et al. 2005).

Table 3.2 Summary of quantitative studies

Study (year)	AIM	Country (city)	Study Design	Length of Follow up /years	Symptom onset to diagnosis	No. recruited /No. eligible (%)		Symptoms collected as 1° or 2° part of study	Reference standard - outcome	No. participants with LC outcome (%)
Hamilton et al. (2005)	To identify the prediagnostic features of LC and to calculate the PPV of symptoms and signs for LC in an unselected population	UK	Case control study	N/A	≤ 2 years	Case: 247/260 (94.6%)	Control: 1235	1°	Histological records to confirm cancer	247 (N/A)
Hoppe (1977)	To study the epidemiology of LC and to determine signs and symptoms that might contribute towards better diagnostic accuracy	Germany	Cohort	N/A	< 6 mnths	20,829 (N/A)		1°	Mediastinoscopy, bronchography and cytological examination of bronchial excretions.	12360 (62.2)
Hippesley-Cox and Coupland (2011)	To develop and validate an algorithm to estimate absolute risk of LC, incorporating symptoms and baseline risk factors	UK	Cohort study	2	2 years	2,406,127 (N/A)		1°	Study outcome (LC) was obtained from patient's GP records, coded using the UK diagnostic codes or from the ONS COD record, coded using the ICD-9/ ICD-10).	3785 (0.16)
Jones et al. (2007)	To evaluate the association between alarm symptoms and the subsequent diagnosis of LC in a large primary care population through the estimation of PPV of alarm symptoms	UK	Cohort study	5	≤ 3 years	762,325 (N/A)		1°	GP diagnoses- no further clarification	2930 (men) 1882 (women)
Kubik et al. (2001)	Role of active smoking in LC risk of women, and to record info on other variables including symptoms	Czech Republic	Case-control study	N/A	< 2 years	Case: 140/161 (87%)	Control: 280	2°	Microscopically confirmed incident LC	140 (N/A)

1°	Primary study	ONS	Office for National Statistics
2°	Secondary study	COD	Cause of death
		BMRC	British Medical Research Council
		PPV	Positive predictive values

Table 3.3 Study methodology- study design and characteristics of exposure data (symptom)

Study (year)	Country (city)	Study Design	Data source of LC	Sample characteristic	Period symptom data was collected before diagnosis	Method of recording symptom
Hamilton et al. (2005)	UK	Case control study	GP records	Population of Exeter aged above 40 years with cases cohort diagnosed with primary LC	≤ 2 years	Coded using the ICPC-2 coding system
Hoppe (1977)	Germany	Cohort study	Medical records	People attending a chest clinic in the state of Northrhine-Westphalia on suspicion of LC	<6 mnths	-
Hippesley-Cox and Coupland (2011)	UK	Cohort study	GP record (using ICD-9 or ICD-10 codes), patient's electronic record (EMIS), QResearch database (including 564 practices in England and Wales)	Population of primary care patients registered to a GP practice in England and Wales; 30-84 years	2 years	Established predictors were recorded from patient's electronic records using a symptom checklist
Jones et al. (2007)	UK	Cohort study	GP records, patient's record	Primary care population registered to a GP practice aged ≤100 and had reported the occurrence of haemoptysis before	≤ 3 years	Occurrences of alarm symptom recorded from patient's record
Kubik et al. (2001)	Czech Republic	Case control study	Participant-reported	Cases were female Czech patients with confirmed LC receiving treatment at the hospital in Prague. Controls were women spouses, relatives or friends of other patients of the same hospital (aged 25-84)	< 2 years	Medical Research Council Questionnaire on Respiratory Symptoms (interviewer-administered).

GP	General practice
ICPC	International Classification of Primary Care
ICD	International Classification of Diseases
EMIS	Egton Medical Information Systems

Table 3.4 Symptoms reported that were independently associated with lung cancer

Study (year)	Criteria for inclusion in the final model	Symptoms													
		Cough	Haemoptysis	Weight loss	COPD	Dyspnoea	Chest pain	Appetite loss	Wheezing	Cough and phlegm	Fatigue	Hoarseness	Heaviness in chest	Fever/flu	General unwell
Hamilton et al. (2005)	p<0.01	✓	✓	✓	✓	✓	✓	✓			✓				
Hippesley-Cox and Coupland (2011)	HR <0.9 or >1.2	✓	✓	✓	✓	✓		✓			✓	✓			
Kubik et al. (2001)		✓				✓			✓	✓					
Jones et al. (2007)			✓												
Hoppe (1977)		None of the symptom diagnostic values were found to be statistically significant													

✓ symptom recorded in the study

Table 3.5 Symptoms measured in the quantitative studies using either predetermined list or questionnaire unless stated otherwise

Study (year)	Symptoms											
	Cough	Haemoptysis	Weight loss	Dyspnoea	COPD	Chest pain	Appetite loss	Wheezing	Cough + phlegm	Asthma	Fatigue	Hoarseness
Hippesley-Cox and Coupland (2011)	✓ ^a	✓	✓	✓ ^a	✓		✓			✓	✓ ^a	✓ ^a
Kubik et al. (2001)	✓			✓				✓	✓			
Jones et al. (2007)		✓										
Hamilton et al. (2005)	All symptoms recorded by GP notes except hoarseness, and shoulder pain due to Pancoast tumour											
Hoppe (1977)	All symptoms recorded in medical notes											

✓ Symptom recorded in the study

^a In the past 12 months

Table 3.6 Diagnostic values (ORs, PPVs, and HRs) for symptoms reported in case-control studies and cohort studies

	CASE-CONTROL STUDIES				COHORT STUDIES							
	Kubik et al. (2001)		Hamilton et al. (2005)		Hoppe (1977)	Jones et al. (2007)				Hippesley-Cox and Coupland (2011)		
	Cases	Controls	Cases	Controls		6 months after symptom recorded		3 years after symptom recorded		PPV	Only HR <0.8 or >1.2 included HR (95% CI)	
	140	280	247	1235								
Time period before diagnosis	2 years		2 years		6 months	6 months after symptom recorded		3 years after symptom recorded			2 years	
Symptoms	PPV* (95% CI)	OR (95% CI)	PPV* (95% CI)	OR (95% CI)	OR (95% CI)	PPV		PPV		PPV		
	-	-	0.40 (0.3 - 0.5)	4.40	-	M	F	M	F	-	M	F
Cough	0.44 (0.28-0.6)	1.92	-	-	-	-	-	-	-	-	1.47 (1.23-1.75)	1.9 (1.56-2.32)
Chronic cough with phlegm	-	-	-	-	1.51 (0.85-2.67)	-	-	-	-	-	1.51 (1.34-1.69)	1.82 (1.57-2.11)
Chronic bronchitis	-	-	0.66 (0.5 - 0.8)	6.99	-	-	-	-	-	-	-	-
Dyspnoea	-	-	2.40* (1.4 - 4.1)	16.24	0.96 (0.32-2.82)	5.8 (5.0-6.7)	3.3 (2.3-4.3)	7.5 (6.6-8.5)	4.3 (3.4-5.3)	6.4 (5.9 - 7.0) ^a	21.5 (19.3-23.9)	23.9 (20.6-27.6)
Haemoptysis	-	-	0.82 (0.6 - 1.1)	4.92	0.79 (0.34-1.82)	-	-	-	-	-	-	-
Chest pain	-	-	1.1 (0.8 - 1.6)	8.14	-	-	-	-	-	-	6.09 (5.33-6.95)	4.52 (3.8-5.38)
Weight loss	-	-	0.87 (0.6 - 1.3)	5.69	-	-	-	-	-	-	4.71 (3.69-6.1)	4.14 (3.15-5.45)
Appetite loss	-	-	0.43 (0.3 - 0.6)	3.07	-	-	-	-	-	-	-	-
Fatigue	-	-	1.6 (0.9 - 2.9)	9.39	-	-	-	-	-	-	-	-
Abnormal spirometry	-	-	-	5.45 (3.81 - 7.79)	-	-	-	-	-	-	-	-
Worsening cough	-	-	-	-	0.58 (0.21-1.6)	-	-	-	-	-	-	-
Flu/Fever	-	-	-	-	0.76 (0.31-1.85)	-	-	-	-	-	-	-
General unwell												

* Positive predictive values (PPVs) in %

^a for a separate validation cohort of the same study

3.5.3 Key methodological strengths and limitations of the studies

The assessment of symptoms; which in itself are highly subjective, requires careful consideration of the methodology and study design. As mentioned earlier, the studies included in this review displayed several key methodological weaknesses which included:

- Lack of standardised data collection
- Retrospective study design
- Recording bias
- Selection bias
- Likely potential confounders
- Limited generalisability

Data collection issues

A variety of data collection methods were applied across the studies.

Symptoms data were either gathered directly from the participants or indirectly from medical notes/databases (refer to Table 3.3); with both methods of data source presenting with their own strengths and limitations.

Four of the five studies extracted their symptom data retrospectively using medical, and GP records or medical databases (Hoppe 1977; Hamilton et al. 2005; Jones et al. 2007; Hippenley-Cox and Coupland 2011). This tends to be the method chosen for most cohort studies of large sample population as methodologically, it could generate the large sample sizes. Kubik et al. (2001) used a standardised questionnaire, the MRC respiratory questionnaire (see Appendix 4, to prospectively record symptoms. A strength of most studies was that they included symptoms reported and recorded at the time of presentation to the clinician, as far as indicated, which reduced the likelihood of recall bias or retrospective reinterpretation of symptoms following diagnosis.

However, reports of symptoms based on these records can be restricted as they only reflect the occurrence of symptoms at the time of reporting and often tend to be diagnosis-focused (Kroenke 2001) rather than symptom-focused; potentially only including symptoms thought to be relevant to a differential diagnosis, resulting in partial recording of symptoms. Therefore, medical notes or records can be subject to recording bias, as clinicians may have recorded symptoms more thoroughly if lung cancer was suspected (Hamilton et al. 2005). The use of prospectively completed checklists or questionnaires might avoid recording bias if administered systematically to patients. However, this stands the risk of only recording common chest symptoms (see Appendix 5 for MRC questionnaire), and excluding systemic symptoms (Kubik et al. 2001). Retrospective symptom reporting is naturally closely related to the perceived diagnosis, which means that patients may not have reported symptoms they considered unrelated to the lung cancer or may not have recognised them as symptoms (further discussed in the qualitative analysis of this review).

Retrospective study design

In most of the studies (Hoppe 1977; Hamilton et al. 2005; Jones et al. 2007; Hippenley-Cox and Coupland 2011), the symptoms were recorded/measured before diagnosis but were obtained from previously recorded data (medical records) not systematically recorded for research purposes. The limitations of retrospective study design for these studies largely relates to recording bias (where prospective patient-reported symptoms are not being interpreted or recalled in light of a diagnosis).

Although a prospective study design is preferred to a retrospective study, it is often difficult to achieve in lung cancer cohorts due to costs involved in large prospective study as mentioned earlier. The concern with retrospective data largely relates to the recall bias and recording error described earlier.

Incomplete reporting of symptoms and recording bias

There will always be some degree of bias in the different methods of data collection available. As already highlighted, the method of extracting symptoms from GP notes could be subject to recording bias. In general, recording bias could over-inflate the likelihood ratios of symptoms; where GPs/clinicians are more likely to record symptoms if lung cancer is suspected. In the same way, the ratios could also be under-estimated if patients with lung cancer under-reported their symptoms, which would then go unrecorded (Hippenley-Cox and Coupland 2011). Likewise, the under-reporting of family history of lung cancer or misreporting past smoking habits could result in an under-estimation of ratios (Hippenley-Cox and Coupland 2011); for example, some ex-smokers may regard themselves as never smokers after many years which could skew the hazard ratio for ex-smokers towards the 1.0.

Four of the five quantitative studies evaluated within this review used clinical records to ascertain symptoms. Therefore, the odds ratios and likelihood ratios obtained might lead to the under or over-estimation of the true predictive values of symptoms.

In the attempt to reduce recording bias in their study, Hamilton et al. (2005) argued that adequate matching would reduce this variation as similar GP recording characteristic would apply to both cases and controls. The study considered age, sex, and more importantly, GP practices when matching their

Systematic Review

cases to controls (Hamilton et al. 2005). However, this will not be adequate to completely eliminate bias as the study would not be able to know whether or not the GP had suspected lung cancer, which would have influence the way symptoms were recorded (Hamilton et al. 2005).

Selection bias

Studies using clinical records that capture data on every patient in a region are less likely to be subject to selection bias than studies that rely on the recruitment or inclusion of individual patients (Hamilton et al. 2005; Jones et al. 2007; Hippesley-Cox and Coupland 2011). However, the eligibility criteria for lung cancer diagnoses in some of the clinical records-based studies included in this review could have some selection effect. For example, Kubik et al. (2001), using a prospective study design (recruitment rate 87%) with the potential for self-selection bias, only included histologically-confirmed lung cancer cases. Therefore, it is possible that those with advanced lung cancer were more likely to be excluded (advanced lung cancer patients will less likely be subjected to further invasive investigation to confirm diagnosis).

Selection bias could be an issue in case-control groups when it comes to matching of the groups. If the cases are significantly different to the chosen control, this would create a falsely magnified difference in effect sizes of the findings. Kubik et al. (2001) recruited their controls from friends and relatives of hospitalised patients with smoking related lung diseases with the assumption that most of the controls could be exposed to long-term passive smoking. This could result in an under-estimation of the odds ratios. However, it is almost impossible to fully control for smoking confounders (as it relies on accuracy and specification of smoking information) in lung cancer epidemiology studies (Rothman et al. 2008).

Likely Confounders

Confounders are variables that could be associated with both the exposure (a symptom in this case) and the outcome, but is not a causal factor in itself (Rothman 2012). It is important to consider and control for potential confounders or risk getting inaccurate results; as the confounder could partially attribute to the observed effect of an effect size on a disease status. Two of the studies did not adjust or controlled for the most likely confounders

(COPD and smoking status) (Hoppe et al. 1977; Jones et al. 2007). Smoking status, and COPD have independent associations with lung cancer and might cause symptoms; 40% to 90% of lung cancer is preceded by COPD (Young et al. 2009). Therefore, COPD and smoking might account for the observed effect of a symptom on the disease outcome (lung cancer) in these studies.

Some of the likely confounders considered in the studies included smoking status, smoking intensity, co-morbidities, gender, and age (Kubik et al. 2001; Hamilton et al. 2005; Hippesley-Cox and Coupland 2011).

Limited Generalisability

All studies recruited adults of both genders above the age of 18 (see Table 3.3) except Kubik et al. (2001); this study only represented women ≥ 18 , in the Czech Republic presenting at a secondary care centre at the time of recruitment. This inevitably reduces the generalizability of their findings to the general population.

Most of the studies were of unselected primary care populations (Hoppe 1977; Hamilton et al. 2005; Jones et al. 2007; Hippesley-Cox and Coupland 2011). Therefore, the results might be transferable to consulting patients in comparable primary care populations. However, the spectrum of disease will differ between referred secondary care and primary care populations and so ratios (LRs, ORs, HRs) and PPVs obtained in secondary care populations cannot be compared to primary care ratios. Similarly, already referred primary care populations (e.g. Hoppe 1977) are unlikely to be generalisable to patients presenting with symptoms in primary care; Hoppe (1977) retrospectively extracted data from medical records in primary care but specifically selected only those who were referred to chest clinics on the basis of possible lung cancer. These latter two studies involving populations referred to secondary care only identified one symptom with an independent association with lung cancer (cough with phlegm) (Hoppe 1977; Kubik et al. 2001).

3.6 Results of Qualitative studies

In this current review, six qualitative papers (O'Driscoll et al. 1999; Corner et al. 2005; 2006; Levealahti et al. 2007; Tod et al. 2007; Molassiotis et al. 2010) were reviewed. Table 3.7 summarises the qualitative studies according to the characteristics of the sample population, data collection method and data analysis performed. The purpose of the review of the qualitative studies was to explore patient's interpretation of symptoms before diagnosis, and whether patients recalled any change in symptom with disease duration to provide a more complete overview of symptoms experienced before lung cancer diagnosis. All studies were retrospective but interviews were conducted close to the time of diagnosis, and therefore, likely to have represented the patient's symptom experience before diagnosis (O'Driscoll et al. 1999; Corner et al. 2005; 2006; Levealahti et al. 2007; Tod et al. 2007; Molassiotis et al. 2010).

Collectively, these studies also provided narratives of the patient experience in their health pathway and explored factors that may have influenced the process of diagnosis with the aim to improve earlier referral for people with possible lung cancer.

Two types of references to time intervals can be found in the literature: the 'time period' (Corner 2005) and the 'total delay' (Tod et al. 2007; Molassiotis et al. 2010). Both terms refer to the time interval between symptom onset and the diagnosis, and use the same information to derive (date of symptom onset and date of diagnosis). Whilst the term 'period' suggest a quantified duration of the length of the event, the term 'delay' is more implicative of an avoidable, unnecessary time lag that deterred diagnosis, which the reviewer may not necessarily completely agree with.

3.6.1 Time intervals between symptom onset and diagnosis

Using an interval event chart to demonstrate the time intervals between key events in the pathways to diagnosis, operable and inoperable lung cancer patients recalled health changes that occurred 12 months (median score) before the time of diagnosis (Corner et al. 2005). Using non-parametric testing, the study found no significant differences between the two groups ($p>0.05$) across the time intervals between key events, and time of diagnosis. However,

the study was under-powered; with only a small sample size of individuals with operable lung cancer (n=7).

3.6.2 Pre-diagnostic bodily experiences (symptoms and health changes)

Participants reported experiencing a broad spectrum of bodily experiences prior to diagnosis (Corner et al. 2005; 2006; Levealahti et al. 2007 Tod et al. 2007; Molassiotis et al. 2010). Both systemic (lethargy, weakness, fatigue, weight loss, and appetite change), and chest and respiratory symptoms were reported.

Levealahti et al. (2007) evaluated the biographical narratives of 37 individual patients with inoperable primary lung cancer, (for those with staged cancer: 19 were late staged IIIb-IV and five were diagnosed at earlier stages). Participants in this study were interviewed a year after they were first diagnosed and had reportedly survived between two months to over three years after their interviews. This suggests a varied survival prognosis possibly relating to some earlier stage of the lung cancer, albeit inoperable, when diagnosed.

Often, systemic symptoms have been related to advanced stages of the disease in previous literature but several studies have reported systemic symptoms in patients with operable lung cancer or patients in the earlier stage of the disease (Corner et al. 2005; Levealahti et al. 2007; Molassiotis et al. 2010). Also, not all the participants in Corner's study (2005) experienced chest-related symptoms (e.g. symptoms suggested in lung cancer referral guidelines).

O'Driscoll et al. (1999) focused on exploring specifically the experience of breathlessness of 52 patients who had been diagnosed with lung cancer. The study reported breathlessness to be continuous in eight patients (15%) whilst the remaining 44 (85%) experienced intermittent breathlessness, often brought on by a trigger. Most patients were also said to have experienced breathlessness at initial presentation and time of diagnosis. However, not much can be commented regarding any changes in breathlessness with the disease progression as it was not systematically recorded by the nurse researchers in the study. According to their patient assessment notes, breathlessness was described as shortness of breath (73%), inability to take a deep breath (21%), and tight band around chest (21%). There will be some

degree of variance to the patient-reported symptom descriptors observed, which are open to the clinician's interpretation in clinical practice and potential recording bias (relating back to the issue of retrospectively collected data).

Patients were unable to discern between 'normal' and 'symptomatic', particularly systemic symptoms (Corner et al. 2005; Levealahti et al. 2007); where symptoms were attributed to occurrences in daily living or 'everyday' bodily changes rather than a health problem (Corner et al. 2005; Levealahti et al. 2007; Molassiotis et al. 2010). In light of this 'normalisation' of symptoms, symptoms then had to become severe before they were presented to the GP/clinician (Corner et al. 2006).

Difficult-to-interpret changes in pre-existing conditions such as chronic respiratory diseases further complicated the diagnostic process, and the appraisal of health change (Molassiotis et al. 2010). Some of these issues overlap with the themes derived from the studies when investigating reasons for delay in lung cancer diagnosis which will be further discussed in the following sections.

3.6.3 Diagnostic delays

Diagnostic delays essentially refer to the total time taken from the onset of the first symptom presentation to the time of diagnosis. 'Delay' is the same as symptom duration but just varies in way the data is interpreted for the intended study. Total delay relates more to understanding the reason for the delay and the analysis of the different components of delays in the process to diagnosis. Therefore, to avoid any misunderstanding and to be consistent with the literature, the review will be using the term 'delay' in the same context as proposed by Burgess et al. (1998); that is participant had waited for more than 3 months before seeking medical advice after being aware of the symptom(s). Five of the six studies reported some form of delay in lung cancer diagnosis, and the reasons for it (Corner et al. 2005; 2006; Levealahti et al. 2007; Tod et al. 2007; Molassiotis et al. 2010).

In the literature, this total delay or diagnostic delay (Corner et al. 2005) can be generally divided into patient delay (appraisal delay, scheduling delay, self-management, stigma), and provider or treatment delay (misdiagnosis, waiting for specialist appointment, professional miscommunication, administrative

errors). Burgess et al. (1998) also identified almost similar types of delay in the diagnosis of cancer: firstly, patient's failure to act on suspicious symptoms, and there is clinician generated delay and finally, there is system generated delay from long waiting times for appointments or test results. Previous studies have suggested GP and system-resulted failure in delaying the diagnosis of lung cancer (Bowen et al. 2002; Koyi et al. 2002) but more recently, Corner et al. (2005, 2006) highlighted that patient delay potentially has larger implications in influencing the timing of diagnosis.

Patient Delays

Patient delay usually refers to the time from the time of the symptom onset when the individual is aware of it until the time he/she seeks healthcare. This delay is described as a deterrence generated by the patient's failure to act on potentially suspicious symptoms for reasons that will be discussed later.

Following unexpected findings of the possibility of patient delays in Corner's study (2005), they carried out a further analysis of their previous interviews to better understand how the individual's action in response to their health change may have influenced the timing of the lung cancer diagnosis. Their findings revealed that all the 22 participants (operable and inoperable) had reportedly delayed seeking medical attention following suspicious symptoms. A median time of 9 months elapsed from the first symptoms to the time of the first consultation (Corner et al. 2005). This delay formed a large proportion of the participant's interview accounts of the months preceding diagnosis.

Reasons for this seemingly lack of awareness about their health were explored and themes revealed that patient did not perceive the symptoms to warrant medical attention. Levealahti et al. (2007) and Molassiotis et al. (2010) also found concurrent themes in their studies. Molassiotis et al. (2010) revealed that most patients did not appraise their initial symptom as serious or had attributed the symptom to other causes, an interval also referred to as appraisal delay in Andersen's model of total patient delay (1995). This was a common observation which also resonated in other studies (Corner et al. 2006). Patients were expecting more severe and extreme symptoms for such a fatal disease such as lung cancer; creating a lack of symptom awareness even amongst those at high risk (Tod et al. 2007; Molassiotis et al. 2010).

Participants were not sure what was 'normal' or they felt the symptoms were minor. Some of the symptoms were attributed to occurrences in daily living or 'everyday' bodily changes rather than a health problem; a process of 'normalisation'. For some of the participants, there appeared to be a lack of association between the bodily experiences as signs of ill-health; viewing them as unrelated occurrences (Corner et al. 2005). Levealahti et al. (2007) reported that some of the participants were not able to recognise the bodily changes as 'symptoms' especially systemic symptoms. Symptomology was almost normalised in relations to daily occurrences (Levealahti et al. 2007; Molassiotis et al. 2010).

Corner et al. (2006) did not find that participants were intentionally delaying seeking help for the symptoms they experienced as a result of fear nor did they think participants were in denial. The study also found that haemoptysis was the only symptom perceived as alarming enough to prompt the participant to seek immediate healthcare. Levealahti et al. (2007) also noted a similar theme in their findings. Tod et al. (2007) and Molassiotis et al. (2010) described how patients waited for their symptoms to resolve through self-medication, and/or self-monitoring until symptoms were regarded severe enough that they sought help which is explained as a threshold effect in Corner's study (2005).

Generally, most of the literature agreed that the existence of co-morbidity (especially those respiratory related) often complicated the process of becoming aware of a new and different disorder (masks new symptom) (Corner et al. 2006; Tod et al. 2007; Molassiotis et al. 2010). There was a tendency to attribute symptoms to other acute or chronic conditions (Tod et al. 2007); in part because of the absence of a tangible symptom or lack of specificity in the symptoms experienced which led to the misinterpretations of the implications of the health changes (Tod et al. 2007; Molassiotis et al. 2010).

Treatment Delays (Provider delay)

Molassiotis et al. (2010) defined this delay as the time patient made first contact with healthcare until their diagnosis to more specifically address the delays experienced by patients due to failure of the healthcare system. This could be through inefficiency or long waiting lists for tests and appointments,

miscommunication between clinicians, or a contradictory discourse between the patient and the clinician which all lead to a delayed correct diagnosis.

Corner et al. (2005) reported a median delay of two months for the time between GP visit or other service and the time of diagnosis. Nevertheless, their findings indicated that once medical action was initiated, the events leading up to a diagnosis was fairly quick for most patients except for three (two were operable and one had inoperable lung cancer). However, this data only applies to the one study with a fairly small sample size.

It is generally difficult to evaluate provider delay as this period would include the waiting times for appointments, and the natural time taken for investigative test results. Taking into account for the urgency of these diagnostic procedures, the fairly short period of provider/treatment delay would be even smaller in effect. However, this is not to disregard the group of symptomatic people who were diagnosed more than 3 months after their first consultation with a doctor. Molassiotis et al. (2010) noted general negative attitudes and assumptions on GP's part in detecting changes in symptoms reported, and misdiagnoses.

Table 3.7 Summary of qualitative studies

Study (year)	AIM	Study design	Method of data collection	Period of data collection	No. recruited/ eligible (%)	Sample characteristics	Method of data analysis
O'Driscoll et al. (1999)	Recorded detailed notes of how patients described their breathlessness and its impact have been analysed and presented in order to offer descriptive material regarding the experience of breathlessness in LC. These data were from a previous study looking to develop and evaluate a breathlessness intervention.	Retrospective	Assessment notes recorded by nurse research-practitioner regarding the reported experience of the patient, or how both nurse and patient agree to describe the symptom and the patient's coping strategies (a collaboration).	Notes were recorded at each subsequent visit to the nursing clinic.	52	30 men (58%), and 22 women (42%) with a LC diagnosis who had completed chemotherapy or radiotherapy and experienced breathlessness aged ranging from 33-76 years (mean age: 60 years). Patients must have been healthy enough to be able to provide adequate material for data analysis.	Content analysis- with frequency counts; descriptive analysis
Corner et al. (2005)	To develop a detailed picture of the pathway to diagnosis by mapping the pre-diagnosis symptom history and the events leading up to diagnosis of a group of patients diagnosed with LC.	Retrospective	Directed interviews- semi structured (entitled "what happened to me?") and structured approach	20 patients- interviewed between 3 days and 4 weeks post-diagnosis, and 2 patients were interviewed 2 to 3 months after diagnosis.	22/30 (73)	12 males and 10 females; aged 42-82 recently diagnosed with LC to map. A third (n=7) of the patients had operable LC and the remaining (n=15) had inoperable LC.	Thematic analysis
Corner et al. (2006)	To further analyse the data of previous study (Corner et al 2005) to re-address unexpected findings/themes that emerged so as to better understand how individuals through the way they responded to their health changes, might have influenced the timing of their LC diagnosis	Retrospective	Results were from the same study conducted (Corner et al. 2005) therefore the study design and sample characteristic are the same (refer to the study above).				

Study (year)	AIM	Study design	Method of data collection	Period of data collection	No. recruited/ eligible (%)	Sample characteristics	Method of data analysis
Levealahti et al. (2007)	To explore how people with inoperable LC frame and conceptualize the onset of their sickness; testing a theoretical construct of viewing illness as a biographical continuity over existing theory of disruption	Retrospective narratives	Semi-structured explorative interviews with open questions	Within 1 year post diagnosis	37	37 patients (21 women and 6 men) aged 48-86 who survived the first year post LC diagnosis. 24 had their LC staged (19 were diagnosed at stage IIIb- IV, and 5 people at earlier stages) but all inoperable LC. Participants survived between 2 months to over 3 years after the interview.	Narrative analysis- deductive approach
Tod et al. (2007)	To identify factors influencing delay in reporting symptoms of LC	Retrospective	Semi-structured interviews. Questions were informed by the study of Corner et al. (2003).	Within 1 year after LC diagnosis of 6-18 months	20	Purposive sample of people diagnosed with LC in the previous 6 months or longer. 18 participants were diagnosed 6 months ago and 2 were 18-month survivors. 8 females and 12 males with age ranging from 47 - 81 years.	Framework analysis (Ritchie & Spencer 1994)
Molassiotis et al. (2010)	To map the pathway from initial persistent change in health to diagnosis of cancer and explore the aptient and system factors mediating this process	Retrospective accounts	In-depth interviews opened with broad question, asking patients' to recall when they first became aware of a change in their health	2-3 weeks after initial cancer diagnosis and prior to or at initiation of treatment	75 cancer diagnosis; 14 (18.7%) were LC diagnosis	In general, the sample characteristic represent patients from seven diagnostic groups, including LC (n=14). Patients aged from 18 to 93 years (mean 58.5 years). Consecutive sample consisted of attendees of outpatient cancer clinics at a large centre hospital.	Content analysis- with frequency counts; descriptive analysis

3.7 Discussion

This review has updated evidence obtained from previous studies, and unlike previous reviews, was not limited to primary care studies. As there is a lack of evidence about early symptom epidemiology in lung cancer, the inclusion of studies in secondary care and non-UK health systems could potentially enable the identification of diagnostics values of symptoms across a broader spectrum of the disease. However, diagnostic values of symptoms obtained in non-UK or non-primary care health services research are not generalisable to symptoms presented in UK primary care (Summerton 2002).

This is also the first review to have included qualitative studies of the available evidence in symptomatic detection of lung cancer. The qualitative studies report lung cancer patients' recollections of their symptom experience before lung cancer diagnosis, with the potential for re-interpretation of symptoms in light of a diagnosis. These studies indicated that patients experienced bodily changes including systemic and non-systemic symptoms, months before diagnosis. Although they reported the occurrence and interpretation of symptoms, the studies did not report the characteristics of symptoms experienced before lung cancer diagnosis in any detail. These retrospective reports indicated that some symptoms were not interpreted as being serious or alarming enough to prompt help seeking, and therefore, were not presented to primary care clinicians until late. Furthermore, less severe and normalised symptoms might never be presented to primary care clinicians by consulting patients (Brindle et al. 2012). This could possibly have implications for the symptoms recorded in the GP clinical notes. Therefore, primary care records-based studies might not provide predictive values of early symptoms not fully elicited by clinicians, or not thought serious enough by the GP/clinician to record. The literature also suggests that symptoms that might precede lung cancer further complicated the process of recognising new symptoms for both patients and clinicians, when there is a pre-existing co-morbidity such as COPD.

The evidence for the predictive values of some of the symptoms for the diagnosis of lung cancer is still weak. Based on the PPVs (PPV>5%), there is little evidence to suggest that symptoms other than haemoptysis consistently predicted lung cancer. This is in keeping with previous studies and reviews

(Hamilton and Sharp 2004). However, individual studies within this review have identified other symptoms that were independently associated with a lung cancer diagnosis such as appetite loss, weight loss, fatigue, and fever/flu presentations, some of which; for example, appetite loss despite increasing lung cancer risk four to five-fold (Hippesley-Cox and Coupland 2011), are currently not included in the UK NICE guidelines as grounds for referral, and may be worthy of further investigation. That being said, stronger claims beyond this cannot be made due to the methodological difficulty of non-comparability of findings across the different studies that extend beyond the spectrum of disease. For example, methods of recording symptoms were highly divergent between studies.

Four of the five quantitative studies reviewed used routine data sets; i.e. symptoms recorded in GP notes or electronic medical databases, with the potential for recording bias, and incomplete presentation of symptoms (Hoppe 1977; Hamilton et al. 2005; Jones et al. 2007; Hippesley-Cox and Coupland 2011). Kubik et al. (2001) did prospectively and systematically record symptoms using a pre-defined checklist, which then raised the possibility of omitting symptoms with diagnostic value for lung cancer. Studies have not yet to date, systematically recorded systemic symptoms as well as non-systemic respiratory symptoms prospectively for research purposes. The possibility of recording biases by GPs/clinicians, and the incomplete recording of symptoms in published prospective studies, limits the value of current evidence regarding the diagnostic value of symptoms for lung cancer.

Gaps in the literature

At present, there is not enough evidence to suggest a signs and symptom profiles from the literature. Undoubtedly, a symptom profile specific to early stage of the disease would be useful as reduction in mortality is related to early stage lung cancer diagnosis. However, there is no definitive evidence to establish a clear temporal relationship between symptoms across the stages of the disease (e.g. hoarseness might be associated with more advanced stage of lung cancer (Hamilton et al. 2005)).

Beyond reporting the frequency of symptoms, very few studies explored and even fewer reported the severity or change in frequency of the symptoms which might occur over time. Hamilton et al. (2005) reported episodes of

cough with each subsequent clinic consultation suggesting persistence of the symptom but not severity. It is very plausible for severity and frequency of symptoms to change over time especially with disease progression. The level of reporting of the disease staging and histology also varied across the studies. Some studies provided more details of the types of lung cancer histology than others but none of the studies included in this review reported the stages of the disease (e.g. I-II, Ia-Ib or TNM). One study reported the percentage of lung cancer cases with operable lung cancer (Hoppe 1977) but none of the quantitative studies distinguished between symptoms of operable and inoperable lung cancer.

Most of the symptoms reported in the literature used generalised terminologies and/or medical jargons e.g. haemoptysis and dyspnoea. So far, there had not been any published empirical study that explored symptomatic diagnosis for lung cancer using patient-reported lay symptom descriptors.

3.8 Conclusion

This review highlights the need for prospective diagnostic studies that systematically elicit, and record symptoms. Patients easily under-appraise the significance of their symptoms suggests the need for questionnaires that elicit a range of 'normalised' symptoms. This applies to all future prospective studies.

Chapter 4: Methodology

4.1 Research Design

This exploratory study was designed to investigate the feasibility of using a patient-completed questionnaire (IPCARD) to identify the predictive value of symptoms for lung cancer in a secondary care population with high rates of respiratory and chest symptoms. Results from the feasibility study were to inform the design of an adequately powered secondary care study to identify the predictive value of symptoms for lung cancer diagnosis in a population with COPD.

This feasibility study investigated:

- (1) The acceptability and content validity of a patient-completed symptom, risk and comorbidity questionnaire in a population referred to secondary care following abnormal CXR (or) the occurrence of potential lung cancer symptoms (patients referred to a lung-shadow clinic). (Study 1)
- (2) The feasibility of identifying patient-elicited symptoms that predict lung cancer in the secondary care population that had been referred to the lung shadow clinic population. (Study 2)
- (3) The feasibility of identifying patient-elicited symptoms that predict lung cancer in a sub-group with COPD that had been referred to the lung shadow clinic. (Study 2)

Study 2 was a prospective study design. The quantitative symptom, risk and comorbidity data were collected before diagnosis using the IPCARD questionnaire (refer to Appendix 5). The same questionnaire had previously been used to identify symptoms that predict chest X-rays suspicious for lung cancer, in a primary care population. This was a lower risk population with fewer patients with COPD in its cohort compared to the current population; a group with higher rates of chest and respiratory diseases that had been referred to secondary care for lung cancer investigations on the basis of abnormal radiological findings and/or presentation of symptoms that might suggest lung cancer according the NICE guidelines for referral.

Methodology

Although the IPCARD questionnaire was already validated, the questionnaire was not used extensively in patients with COPD, and had not been validated for use with a population with COPD and lung cancer pre-diagnosis. Earlier IPCARD feasibility work had not carried out in-depth comparative analysis of those with and without lung cancer in the COPD population. The current study was focussed on COPD, and therefore, a detailed qualitative analysis between those with and without lung cancer in a COPD population was required to evaluate the content validity in the COPD subgroup. This ensures that the questionnaire would be able to distinguish symptoms between the two groups; those with COPD only, and those with COPD and lung cancer (if there is a difference in symptoms).

Although evidence of symptoms that predict lung cancer in a referred secondary care population with high rates of chest and respiratory problems have little clinical use in primary care, predictive values of symptoms that predict lung cancer in the COPD sub-group with a similar spectrum of disease in primary care might be clinically more relevant. However, it should be noted that symptoms that identify patients with lung cancer in one population might not do so in another population with a different spectrum of disease (Hamilton et al. 2005; Ransohoff and Feinstein 1978; Whiting et al. 2004).

4.2 Research Plan

There were two parts to this research project, Study 1 and Study 2.

Study 1 was a qualitative study using semi-structure and cognitive interviews. Semi-structured interviews were used to evaluate the acceptability of the IPCARD questionnaire to a pre-appointment population in a secondary care clinic that was being investigated for potential lung cancer. Interviews were semi-structured to take on a more conversational, and yet, focused approach to explore participants' health experience leading up to their clinic appointment. The inclusion of broad, open-ended questions where the researcher could follow relevant leads produced richer data that enabled the understanding of symptom experiences of individuals with COPD; how they interpreted their health changes and how they described these events.

Cognitive interviews were also used to investigate the content validity of the use of the questionnaire in sub-groups with COPD, and COPD with lung cancer.

An interview schedule was designed with key questions, and a symptom checklist, to be used as a reference, and as prompts if necessary. The intention was to use the schedule more spontaneously during the interview rather than the need to refer to the schedule explicitly. The use of cognitive interview was appropriate to find out content validity (Willis 2005).

The interviews were carried out retrospectively, at least a week after diagnosis, and could be subjected to recording bias; particularly where the effect of the outcome (the diagnosis) might influence the individual's perception of the symptoms. However, the decision to conduct the interview retrospectively was agreed to minimise any anxiety that the interview might cause whilst a participant was awaiting imaging results following their attendance at the lung-shadow clinic. Furthermore, interviews carried out within a week of diagnosis can minimise the potential for recall bias from the time the questionnaire was completed.

Study 2 systematically recorded patient-elicited symptoms, risks and comorbidity data using the IPCARD questionnaire to investigate the feasibility of identifying symptoms that predict lung cancer in a secondary care population with a high prevalence of chest and respiratory symptoms, as well as in a sub-population with COPD. This feasibility study was carried out prospectively (before diagnosis). Study 2 was designed to minimise the problem of recording bias, making it a methodologically robust and rigorous study.

The use of IPCARD to systematically record patient-reported symptoms addresses the potential recording bias observed in current evidence using GP notes or CPRD-based studies; see Systematic Review (Chapter Three). People with lung cancer might place more emphasis on symptoms, or interpret symptoms in light of a diagnosis than those not diagnosed with cancer. The lack of a systematic method of data collection in studies that rely on GP records and medical notes could also result in similar recording bias. Primary care clinicians or GPs are more likely to report symptoms when lung cancer is suspected.

Methodology

Patient-elicited symptoms are more likely to be complete and accurate than data drawn from medical notes, or clinician-reported symptoms that used a questionnaire or a checklist. The threshold for symptom reporting is higher in the latter study design, and is likely to under-represent the experiences of patients.

The quantitative work (Study 2) ran alongside Study 1 to investigate the questionnaire use in a secondary care population with COPD that had been referred to a lung-shadow clinic, to inform future COPD studies.

4.2.1 Missing data

Despite the strengths of the study methodology (prospective systematic data collection), there was a proportion of missing data because it is information from a patient-completed questionnaire. Missing data was less of a concern in previous studies that relied on GP notes and patients' records because the absence of a symptom recorded just means no symptom; however, in the case of the questionnaire study, missing data is more of a problem. Analytical methods are required to address the missing data using imputations, which will be discussed in greater length in Chapter Six.

4.2.2 Quantitative analysis

The quantitative data analysis consisted of multiple regression methods that were used to identify symptoms independently associated with lung cancer, adjusting for potential confounders such as age.

Symptom variables for selection were identified using univariate analyses at two statistical significance levels ($p < 0.05$, and $p < 0.15$) to be included into the multivariate analyses. As this study is exploratory rather than explanatory, the criteria for variable selection was relaxed to identify potential symptoms that could be further investigated in the fully-powered study.

Confounders are defined as factors that are associated with both exposure and outcome but do not have any causal effects (dos Santos Silva, 1999) (see Figure 4.1).

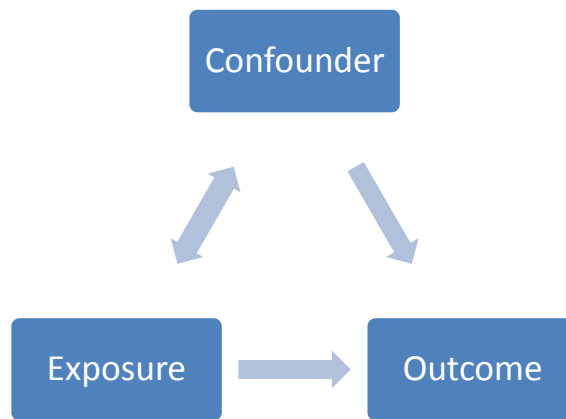


Figure 4.1 Illustration of a confounder

Smoking, is a cause of lung cancer, and might directly cause symptoms that are also associated with, and appear to be caused by, lung cancer. Adjusting for smoking would exclude the possibility that smoking rather than lung cancer was the cause of the symptom. However, as smoking is a key cause of lung cancer (a risk factor), adjusting for smoking in a multivariate model might also obscure a relationship between lung cancer and a symptom, where lung cancer is the cause of the symptom. In other words, the relationship between lung cancer and the symptom might not be independent of smoking. It is possible to explore the relationship between the symptom and lung cancer within the strata of the third variable (smoking behaviour). As this study is underpowered to carry out separate analyses in those who, for example, do and do not smoke, initially models will be developed that do not include key risk factors (causal variables) for lung cancer (Model 1). Risks factors such as smoking behaviour and COPD, that are also causes of lung cancer, will be added into a final model (Model 2). Potential confounders, which also include comorbidities associated with lung cancer, might be found not to have any confounding effect in the final model, and will therefore be removed.

Overlooking confounders can overestimate or underestimate the effects of the true association between an exposure and an outcome variable. However, determining whether a variable is a confounder can be difficult in practice. Once confounders are identified, they can be adjusted in the regression modelling; adding a confounding variable as a predictor in the final model will control for the effect of the confounding.

Methodology

Risk variables were recorded in the questionnaire to also look at interactions, but this study was not adequately powered to do this. However, interactions will be investigated in the larger, well-powered secondary care IPCARD study.

Multivariate analyses were used to identify independent relationships between symptoms and lung cancer, but individual symptoms alone are neither particularly predictive, nor clinically useful. Therefore, a simple score (weighted and un-weighted score) was developed for each model to distinguish between people with lung cancer and those without lung cancer. The set of diagnostic criteria is easy to apply and straightforward. Weighting the criteria takes into consideration the relative effect of each variable.

The performance of the scores were interpreted using the positive likelihood ratio at the optimal cut-off decided by Youden's index (1950), and the area under the curve (AUC) of receiver operating characteristics (ROC) curves.

4.3 Setting

This feasibility study was the first stage of a larger IPCARD symptomatic lung cancer detection project (The IPCARD Chest Clinic Study) led by Dr Lucy Brindle based within the Faculty of Health Sciences in the University of Southampton. This PhD project recruited subjects from a regional chest clinic (the lung shadow clinic) in Southampton University Hospital Trust (SUHT).

4.3.1 Organisational structure of the clinic

The clinic was set up for patients suspicious for lung cancer; which included patients with suspicious radiology imaging (abnormal CT or chest X-ray imaging), and those referred with symptoms potentially indicative of lung cancer. Patients are seen by a respiratory physician who schedules appropriate investigations which usually take place within the next 7 days (biopsy, CT or PET). Following the appropriate investigation(s), patients will then be seen by the respiratory consultant again in clinic within the next 7 days to receive a diagnosis. This pathway is demonstrated in Figure 4.2.

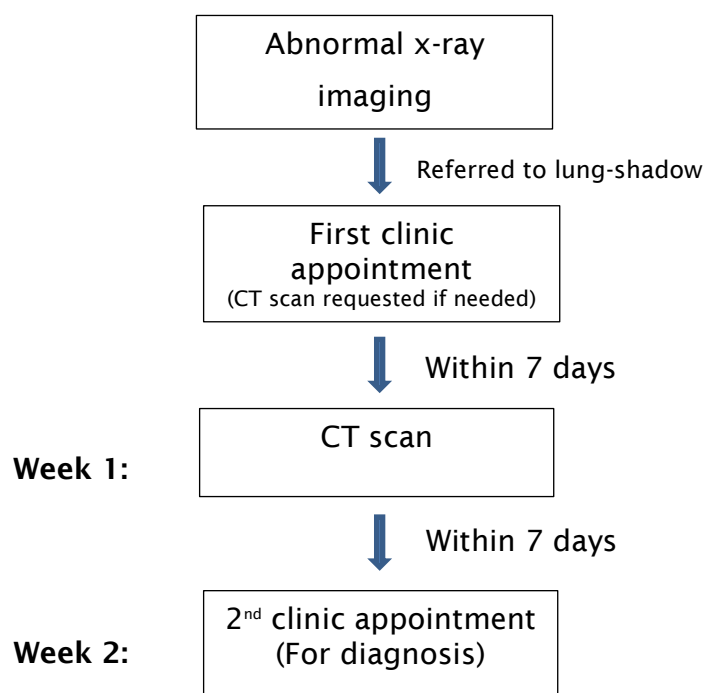


Figure 4.2 Operational structure of lung clinic in SUHT

Clinical collaborators at each site have granted permission for recruitment within their lung-shadow clinics.

4.4 Study Population

All participants would have been referred to their regional lung-shadow clinic following abnormal lung imaging results (chest X-ray) or on suspicion of lung cancer. It was anticipated that about 60% to 70% of the participants would have been referred by their GPs in primary care, including those patients admitted via the two week wait (2ww) rule.

The population of **Study 1** is a sub-group of **Study 2**, whereby those recruited to Study 1 had first completed the IPCARD questionnaire during Study 2.

4.5 Inclusion Criteria

- Participants must be aged 40 years and above.
- New attendees of the lung-shadow clinic (considered new patient cases of the clinic).

Methodology

The incidence of lung cancer is very low in people below the age of 40.

Therefore, subjects eligible for recruitment must be aged 40 years and above and have been referred to a lung-shadow clinic on the basis of abnormal lung imaging results (for investigation of possible lung cancer).

4.6 Exclusion Criteria

The study's exclusion criteria were:

- Subjects below the age of 40.
- Subjects who have received a diagnosis of 'probable lung cancer' (>90%) before their attendance at the chest clinic.

Previously diagnosed patients were excluded to minimise recall bias and recall error. The study aimed to collect prospective data.

4.7 Recruitment methods

Following NHS ethics and R&D approval, the researcher attended the clinic on the designated clinic days (half a day a week at each health care site). At the start of the clinic, the clinical nurse specialist approached the eligible attendees of the lung-shadow clinics and asked if they would be happy to speak to a researcher. The researcher would then introduce the study to potential participants whilst they were waiting to be seen by the respiratory clinician. Participants who expressed an interest in participation were given a study pack. The information packs consisted of an introductory letter (appendix 6), participant information sheets (appendix 7 and 8), consent form incorporated into the questionnaire, and a prepaid return envelope.

Throughout the recruitment process, the researcher ensured that the participants were fully informed about the study, and provided ample opportunities for participants to ask questions before giving out the questionnaire. If participants did not wish to fill in the questionnaire at the clinic, they were given the option to take the questionnaire home and return it by post within 24 to 48 hours if he/she would like to take part in the study. Therefore, the consent form had been designed for self-completion, which meant that the consent form would not require the researcher's signature (appendix 5). This option allowed potential participants to take the form away

from the clinic, and also provided participants more time to carefully consider and understand the participant information sheet (appendices iii and iv) before deciding whether or not to participate, thus further facilitating informed consent. Previous IPCARD feasibility studies have demonstrated that the questionnaire items can be self-completed.

Flow diagrams of the steps of the recruitment process according to the organisational structure of each individual lung-shadow clinic (Southampton lung-shadow clinic) can be found in appendix 9. A summary of the recruitment process is illustrated in Figure 4.3.

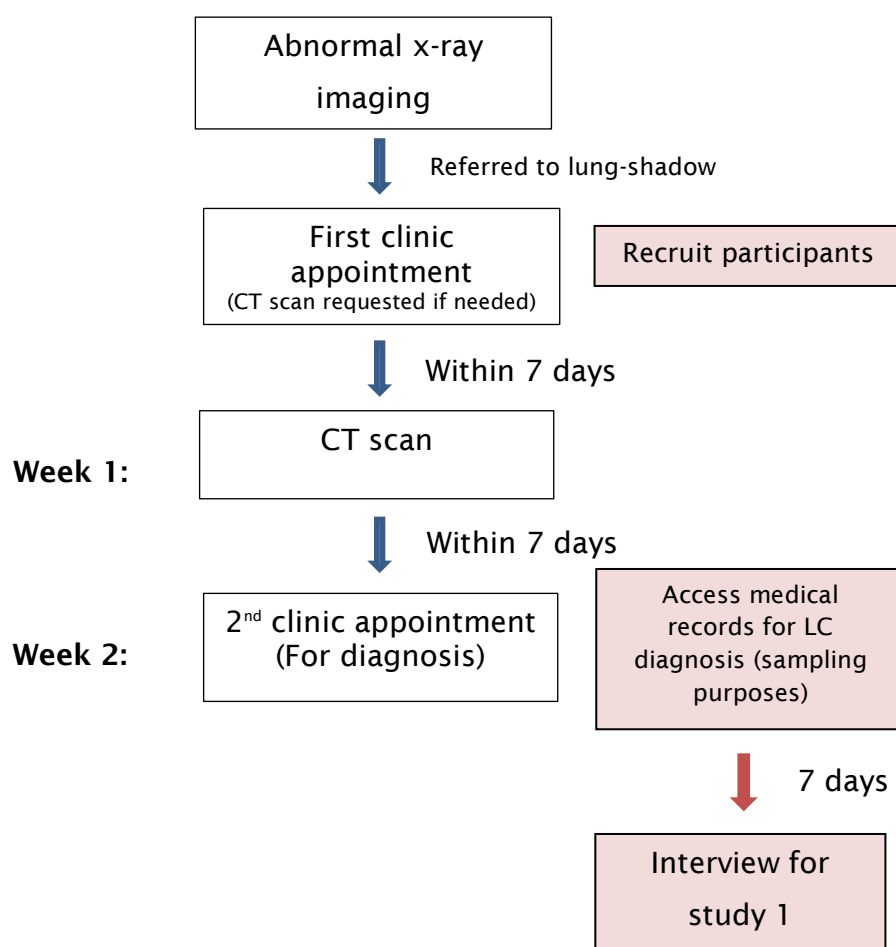


Figure 4.3 Participant recruitment process

Study 2 aimed to consecutively recruit lung-shadow clinic attendees to gather data relating to prospective symptoms, risks and co-morbidities data. Pragmatic consecutive recruitment was used to obtain a sample that was as representative as possible of the lung-shadow clinic population.

Methodology

Giving participants the opportunity to complete the questionnaire at the lung-shadow clinic was likely to increase the quantity of prospective data gathered and the participants' response rate. The IPCARD Feasibility Study has indicated that response rates decline when participants are required to post the questionnaires back to the researcher. Participants who completed the questionnaire in the clinic would still be able to withdraw from the study if they wished to do so. As some clinic attendees might not have time to complete the questionnaire prior to their appointment, or may not wish to remain and complete the questionnaire following their consultation, an option to post back the questionnaire was necessary, and therefore provided.

The study recorded the date questionnaires were returned, and any personally identifying information of the questionnaires was removed in preparation for automated data entry. For data to be prospective, questionnaires received that were postmarked after the date of diagnosis were excluded from the study.

4.8 Development of the IPCARD Questionnaire

The study used a patient-completed questionnaire called the Identifying Symptoms that Predict Chest and Respiratory Disease (IPCARD) questionnaire (Brindle et al. 2015). This symptoms, risk and co-morbidities questionnaire was designed to elicit information about symptom experiences or changes in health status that might have the potential to distinguish between those with lung cancer and non-malignant disease that commonly occurs in those with a smoking history. Conceptually, the IPCARD questionnaire was developed with newly diagnosed early stage lung cancer patients and later refined with chest X-ray clinic attendees (IPCARD Feasibility Study).

The questionnaire contained items identified in the International Primary Care Respiratory Group guidelines (designed to identify COPD, and to distinguish between COPD and asthma (Levy et al. 2006)), and included lay descriptors of breathlessness identified in studies with patients with asthma, COPD, interstitial lung disease, cardiac failure, and lung cancer (Wilcock et al. 2002). Lay descriptors of symptoms experienced by those in the late-staged lung cancer and early-staged lung cancer were also identified through qualitative research with predominantly late stage lung cancer patients (Corner et al. 2005) and operable lung cancer patients (Brindle et al. 2012). These symptoms

descriptors were used to record symptoms and health changes in 11 areas: cough, chest/shoulder pain, breathing, joint/bone aches, weight loss/gain, tiredness, haemoptysis, hot/cold sweats, voice changes, eating and taste change, skin condition, and any other changes in health in the last two years.

Relating to format and data analysis, most of the items in the IPCARD questionnaire were fixed response items. This allowed items to be converted into categorical variables, which would fit into multivariate analyses. The questionnaire also enabled automated data entry using an optical mark reader format. There were two open-ended questions that ask about other health changes in the last two years and family history of certain ill health. The IPCARD questionnaire also recorded established epidemiological risk factors (smoking history, family history of lung cancer, known occupational exposures and COPD), and more recently identified risk factors (pneumonia in the previous five years, and previous malignancy) (Cassidy et al. 2008). The questionnaire entailed all aspects of health from everyday health to changes in the health and symptom experience prior to diagnosis.

The IPCARD Questionnaire was developed in an operable lung cancer population, and further validated in a heterogeneous GP referred chest X-ray population with high rates of chronic respiratory disease and chest and respiratory symptoms. Therefore, content validity was investigated in the operable lung cancer population and in the GP referred chest X-ray population. The questionnaire was found to record a comprehensive range of symptoms experienced by both populations. The response and data completion rates were high (>70%), and investigation of test retest reliability indicated substantial to outstanding agreement for most generic symptom items; agreement was moderate for 'breathing changes current or in the last 3 months' in the chest X-ray population (Brindle et al. 2015).

4.9 Identifying diagnoses (lung cancer and COPD)

All diagnoses were extracted from medical notes; secondary care electronic records of those who had attended the lung-shadow clinic, after a six-month follow-up period. Most diagnoses would have occurred within 2 months of investigation (in line with Department of Health's (DoH) recommendation for a maximum 62-day wait from referral to first treatment for all cancers). The

Methodology

decision to follow up at six months was based on anecdotal clinical experience (clinical collaborator- Dr Anindo Banerjee) and research experience (IPCARD Feasibility Study) to ensure no lung cancer cases were missed. It was important that the study successfully identified all the lung cancer cases; also at six months, the confirmation of lung cancer type and stage was more likely to be available. Follow-up also determined the operable and inoperable lung cancer cases, which were categorised into; ever-surgically referred, surgically-removed lung cancer, and non-surgically treated.

A sub-group of patients with a COPD diagnosis were defined by:

(1) Abnormal spirometry post-bronchodilator: defined as a forced expiratory volume in 1 second (FEV1) <70% of predicted and forced expiratory volume in 1 second/ forced vital capacity (FEV1: FVC) ratio <0.70. This is based on the American Thoracic Society's standards agreed by most respiratory guidelines.

AND/OR

(2) Clinical diagnosis by the respiratory physician in the lung-shadow clinic, based on the presence of symptoms such as recurrent cough (productive), wheezing, dyspnoea in conjunction with pre-disposing risk factors such as smoking, age and family history.

AND/OR

(3) Evidence of emphysema including structural damage presented on lung imaging results.

Two definitions of COPD will be used to define the COPD cohort, for which two separate analyses will be performed. The group defined by the first definition included all COPD participants who meet any one of the three criteria for confirmed diagnosis. The other group (definition two) included those who meet either the FEV criteria or symptom criteria, and excluded those defined on the basis of CT evidence of emphysema alone. The rationale for this divided classification was to minimise differences in disease spectrum between primary and secondary care identified COPD populations, thereby it was unlikely that patients diagnosed with COPD in primary care would be diagnosed on the basis

of a CT scan alone. Heterogeneity across patient sub-groups limits the comparability and transferability of the findings.

4.10 Data entry and data cleaning

The type of data provided in the IPCARD questionnaire included binary, categorical, and ordinal responses. The diagnosis of lung cancer (outcome) was coded a dichotomous variable (yes=1; no=0). Data were cleaned in Microsoft Excel, and checked for consistency.

4.11 Ethics

Ethical approval for the study was granted from [Southampton and South West Hampshire (B) Ethics Committee REC no. 12/SC/0490]. The study protocol had also gone through review from the Cancer Care Directorate Protocol Review Committee in SUHT. The study had been given permission to recruit from October 2012 to February 2014.

The proposed study used a very similar study design and patient information to that used in the NIHR funded IPCARD Feasibility Study (ethics REC no. 10/H6504/72). The study design and patient information sheet was designed with service user-representatives in Southampton, and the Pan Birmingham CLRN over a period of 4 years. The study and questionnaire design involved lay members of the consumer research panels in Southampton and Birmingham, chest physicians and radiographers at SUHT, and participants in the pilot study which developed the questionnaire. The IPCARD Feasibility Study uses the IPCARD questionnaire to record the symptom experiences of those referred by their GP for a chest X-ray. Many of the ethical issues raised by the IPCARD Feasibility Study were common to the proposed IPCARD chest study. Nevertheless, the proposed study carefully evaluated the acceptability of the questionnaire in Study 1 for use in a higher risk population, who might have been more anxious than a chest X-ray population of previous IPCARD feasibility study. Study 1 explored and carefully monitored participants' responses to an invitation to participate in the study, and any concerns raised to indicate any changes that need to be made. Any changes to the research design would have been submitted to the ethics committee before starting the main study.

Methodology

4.11.1 Data protection and Confidentiality

In the interest of data protection and confidentiality, all patient identifiable information on questionnaires were not scanned for data entry, or stored electronically. This information (on the first page of the questionnaire) was removed from questionnaires and stored in a separate locked filing cabinet in the Faculty of Health Sciences at the University of Southampton for 10 years. Audiotape recordings and transcripts were also stored in similar way.

Chapter 5: Study 1

5.1 Introduction

The purpose of the Study 1 was to investigate the acceptability and content validity of a patient-completed symptom questionnaire in a secondary care population being investigated for lung cancer using qualitative research methodology (semi-structured and cognitive interviews). The semi-structured section of the interview explored participants' experiences of completing the IPCARD questionnaire; including any emotional impact, or identify any items that were particularly sensitive to this population. The full range of symptoms and co-morbidities experienced by those with COPD referred for lung cancer investigation were also investigated in the semi-structured interview while, the cognitive interview investigated the ease of interpretation of the questionnaire items, and reasons for any missing data.

Previous IPCARD feasibility studies have confirmed the content validity, acceptability, and test retest reliability of the IPCARD questionnaire in a population of GP referred chest X-ray attendees, which included some patients with COPD but not patients with COPD and lung cancer prior to diagnosis. Therefore, careful consideration was required to assess the IPCARD questionnaire for use in this higher risk population awaiting diagnosis by comparing the content validity in those with COPD and lung cancer, and those with COPD.

Objectives:

- (1) Explore participants' experiences of completing the questionnaire prior to diagnosis.
- (2) Determine whether any of the questionnaire items were consistently misinterpreted by this population and identify reason(s) for any inconsistencies observed.
- (3) Identify any symptoms potentially relating to lung cancer and chronic respiratory diseases not captured by the IPCARD questionnaire in this population referred to secondary care with COPD.

5.2 Methods

5.2.1 Sample size for qualitative research

We aimed to recruit approximately 10 participants to take part in a semi-structured and cognitive interview. Participants were purposefully selected for the interviews on the basis of their questionnaire responses and diagnoses; lung cancer diagnosis and COPD diagnosis, and any other common co-morbidities within this population) to obtain a maximum variation sample or heterogeneous sample (Patton 2002). Previous feasibility studies developed and evaluated the questionnaire in a population of newly diagnosed lung cancer patients, and also in a population of GP referred chest X-ray attendees, so it was anticipated that any changes required to the questionnaire would be minor.

Cognitive interview usually uses small samples due to the vast amount of rich data generated (Willis 2005). Nevertheless, sampling for the qualitative study (Study 1) was carried out until data saturation was achieved.

5.2.2 Qualitative data collection (Interviews)

There was a semi-structured and cognitive section to the interview process. The semi-structured interview explored participants' experience of completing the IPCARD questionnaire, and recent health, in order to identify any symptoms experienced by the participant that were not recorded by the questionnaire. This was followed by a cognitive interview, in which the researcher used the 'thinking aloud' method to reveal participant's interpretation of the questionnaire items, and identify the context of difficulties in questionnaire response; the interviewer read out specific questions selected from the questionnaire based on certain criteria (e.g. non-response, or inconsistency with earlier interview responses) and the participants were encouraged to think out loud their thought-process in arriving with an answer. An interview schedule had been prepared for the purpose of this study (see appendix 10).

Verbal probing was used to encourage the participant to elaborate on their responses where necessary. Content validity of the IPCARD questionnaire has been investigated in a lung cancer population and in a chest X-ray population which had included people with chronic respiratory disease. Therefore, it was

anticipated that few changes would be required to improve the content validity of the IPCARD questionnaire. However, it was possible that certain items may cause more anxiety, or be interpreted slightly differently in this population. All interviews were tape-recorded and transcribed verbatim.

5.2.3 Qualitative data analysis

The purpose of the qualitative analysis was to determine acceptability, symptom experience, and interpretation of questionnaire items to identify content validity, and reasons for missing data.

The theoretical framework underlying this study aligns itself to Hammersley's concept of 'subtle realism' (Hammersley 1992). 'Subtle realism' states that all research involves subjective perceptions and observations, and accepts that different methods produce different perspectives (Kirk and Miller 1986; Hammersley 1992). The philosophy of 'subtle realism' attempts to represent 'reality' rather than to attain 'the truth'. Therefore, Hammersley's (1992) position is one compatible with the perspective of combining research methodologies (qualitative and quantitative). This fitted in with the qualitative methodology used in Study 1; in an attempt to identify the 'reality' of symptom experiences for participants, but also taking into consideration the context within which the account of symptoms was produced, and the data collection method used.

Thematic analysis

Thematic analysis was used to analyse data from the semi-structured interviews. This involved in-depth reading and consideration of the content of the interview transcripts to obtain a detailed knowledge and level of understanding to develop a thematic framework. Analysis began by coding the data. A 'code' is a segment of the data which is relevant to the research question (Boyatzis, 1998; Braun and Clarke, 2006). The initial coding of the data was guided by responses of participants in the questionnaire based on the aims of the qualitative study. From the coded data, patterns were explored and identified to form 'themes'; these themes helped to highlight important topics using rich descriptions (Braun and Clarke, 2008).

Study 1

Key issues within the data; based on the study aims, and guided by responses of participants in the interview (experiences of taking part in the study, symptom experiences, and interpretations of questionnaire items), were identified when constructing the thematic framework. The coding frame was then applied to the interview text and text was arranged accordingly to each theme identified.

Content analysis

For the purpose of this study, content analysis of the cognitive section of the interview was also used to identify the content and context of difficulties in questionnaire response including unclear items) and any reasons for missing data, *which apply to this population* but were not apparent in earlier evaluations of the questionnaire. Any discrepancies between questionnaire and interview responses, and any ill health not captured by the questionnaire would be identified to evaluate the content validity of the questionnaire for use within this population of high rates of respiratory diseases.

Therefore, findings from this analysis would have informed any minor modifications that might be needed to the symptom questionnaire or instructions to participants in future studies.

5.3 Results

5.3.1 Recruitment to the Qualitative research

A sample of 10 patients recently diagnosed with lung cancer and/or COPD were recruited from Southampton General Hospital and interviewed about their health experiences in the previous 12 months. All participants had COPD. Earlier IPCARD studies had not investigated the content validity of the questionnaire for participants with COPD and lung cancer before diagnosis. Furthermore, an objective of the study was to identify lung cancer symptoms in a sub-group with COPD. Therefore, the qualitative study was designed to compare the content validity of the questionnaire for these two populations were recruited.

Participating lung cancer interviewees were purposefully selected to achieve equal number of lung-shadow clinic attendees diagnosed with COPD only, and

those with COPD and lung cancer; with varying degree of the disease(s) to allow for maximum variation sampling. Eligible respondents who did not agree to being interviewed during the recruitment (those that had not signed the box consenting to an audio-recorded interview) were not contacted. 11 of the 21 participants invited to take part, declined to participate in an interview. All those who agreed to participate were interviewed.

All ten respondents were interviewed between seven and 14 days after a confirmed diagnosis; to minimise anxiety levels during the period of anticipation, waiting for results of their diagnosis.

Participants in this study reflect the characteristics of the lung cancer population in the UK. Participants were aged between 52 to 76 years; mean age: 65.5 years. There were slightly more males (n=6) than females (n=4) among the participants, which is similar to the UK statistics of lung cancer incidence (3:2; male: female) (CRUK, 2009). Participants' characteristics can be found in Table 5.1

Study 1

Table 5.1 Characteristics of interview participants

	N (%)
AGE (years)	
median	64.0
mean	65.5
Gender	
Female	4 (40)
Male	6 (60)
Diagnosis	
COPD only	5 (50)
mild	1 (20)
moderate	2 (40)
severe	2 (40)
COPD + LC	5 (50)
mild	2 (40)
moderate	3 (60)
Smoking status	
Ever smoker	10 (100)
Never smoker	0
Current smoker	3 (30)

5.3.2 Characteristics of people who declined to be interviewed

11 of the 21 participants declined to be interviewed. All 11 participants in this group had a diagnosis of both COPD and lung cancer (8 inoperable lung cancer, 3 operable lung cancer- referred for surgery); aged between 68 to 90 years (mean age: 77.8). Of the 11 participants, 6 were males and 5 were females.

5.3.3 Results of Qualitative Analysis

The examples used to illustrate each theme are quoted as spoken by the respondent (*sic erat scriptum*) (Boyatzis, 1998); each respondent has also been given a unique identification number (ID) to maintain anonymity (**P01** to **P10**).

5.3.3.1 Acceptability of the questionnaire to a population of lung-shadow clinic attendees

Theme 1: Emotional impact of completing the questionnaire

In the first theme, questionnaire respondents discussed the impact of filling in the questionnaire in the lung clinic.

Participants did not report feelings of anxiety as a result of completing the questionnaire. Instead, issues raised related to the emotional implication of the context in which they were filling out the questionnaire; the apprehension of waiting for a potentially anxiety-provoking clinic appointment. One participant noted that although completing the questionnaire did not cause any anxiety, *“It didn't bother me at all but some people might be a bit **afraid of being there** etc. and it might just throw them”* suggested that *“it might be a good idea to put if you are sending a letter out for an appointment that you intend to do this- to have an explanation that you might be required to do this survey”* (Participant 08, age 68, diagnosed LC and COPD); to inform future attendees of the possibility of completing the questionnaire.

Similarly, another participant suggested that it might be helpful to fill in the questionnaire at home instead (Participant 01, age 62, diagnosed mild COPD):

“... when you are sat in that situation in hospital maybe you are thinking of other things than questionnaires, and sometimes maybe sitting at home it would be easier to answer”

“... because when you are in the hospital maybe you are worried or anxious, and so your brain doesn't always work as it should do”

These accounts suggested a level of anxiety and sensitivity as a result of being in the clinic; having to attend clinic under such uncertain circumstances where one is waiting on results of potentially worrying diagnosis/outcome.

Theme 2: Time

Time-consumption was a factor commonly mentioned by participants. One of the participants commented that the questionnaire took longer to complete than the time that was available whilst waiting for the appointment.

"I mean to fill that in correctly, you probably need about half an hour to thoroughly go through ... everything rattles through so fast though it's unbelievable ..." (Participant 02, age 65, diagnosed moderate COPD).

"when you write one of these out when you are in a hospital you are doing it quicker than you should be because you are waiting to go in to see someone" (Participant 03, age 52, diagnosed moderate COPD).

On reflection, some of the participants mentioned that they should have completed the questionnaire at home and posted it back:

"... the patient should be sending it back to give them more time to read it and understand it properly" (Participant 03, age 52, diagnosed moderate COPD).

However, all participants were given the option of filling out the questionnaire at home and pre-paid postal return.

5.3.3.2 Inconsistencies between symptom experiences reported in the interviews post-diagnosis and those recorded in the questionnaire pre-diagnosis

Table 5.1 demonstrates the symptom responses elicited in the interview (after diagnosis) and those recorded in the questionnaire (before diagnosis) for comparison, across all 10 participants who were recruited in the qualitative study. 'Discrepancy' refers to an event of inconsistency where a key symptom was reported in the questionnaire but not in the interview following diagnosis, or vice versa. Discrepancy and inconsistency was used interchangeably throughout the chapter, where appropriate.

13 individual cases of inconsistencies between the interview-elicited symptoms and questionnaire-elicited symptom responses were found in the following main symptoms; fatigue (n=4), breathlessness (n=4), pain/ache/discomfort (n=2), hot/cold sweats (n=2), and cough (n=1).

Systemic symptom, 'fatigue', and non-systemic symptom, 'breathlessness', had the highest number of cases discrepancies. Three of the four case discrepancies in 'breathlessness' related to breathing changes recounted during the interview (post-diagnosis) that had not been elicited by the questionnaire. All case discrepancies in 'fatigue' were in relation to fatigue elicited by the interview only (post-diagnosis); i.e. more participants were ready to report symptom of 'fatigue' in the interview (see Table 5.1). Generally, the discrepancies were observed to be in the same direction for 12 out of the 13 occurrences- symptoms were more likely to be elicited in the interview than by the questionnaire.

Possible reasons for these observations were discussed within the findings of the thematic analysis.

Study 1

Table 5.1 Comparison of symptoms in questionnaire responses and interview symptom data

Participant ID	Diagnosis	Pain/ache/discomfort		Cough		Breathlessness		Fatigue		Haemoptysis		Weight Loss		Sweats		Voice change		Taste change	
		Int	Qx	Int	Qx	Int	Qx	Int	Qx	Int	Qx	Int	Qx	Int	Qx	Int	Qx	Int	Qx
01	COPD	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	\	Yes	Yes	Yes	\	No	Yes	Yes
02	COPD	No	No	Yes	Yes	Yes	No	Yes	Yes	No	No	No	No	No	No	No	No	No	No
03	COPD	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	No	Yes	Yes	\	No	\	No	No	No
04	COPD	Yes	No	Yes	Yes	Yes	Yes	Yes	No	No	No	\	No	No	No	No	No	\	No
05	COPD	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	No	No	Yes	No	No	No	No	No
06	LC + COPD	Yes	No	Yes	Yes	Yes	No	Yes	No	No	No	No	No	Yes	Yes	No	No	Yes	Yes
07	LC + COPD	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	No	No	No	No	Yes	No	Yes	Yes	No	No
08	LC + COPD	No	No	Yes	Yes	No	No	No	No	No	No	Yes	Yes	No	No	No	No	No	No
09	LC + COPD	Yes	Yes	Yes	No	Yes	No	Yes	No	No	No	No	No	No	No	Yes	Yes	\	No
10	LC + COPD	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	No	No	No	Yes	Yes	No	No	\	No

Theme 1: Distinguishing between normal and abnormal bodily sensations and changes in health

Some symptoms not reported on the questionnaire tended to be normalised by participants. Often, these tended to be the less specific systemic symptoms, and were not likely to be symptoms that had led them to see their GP or prompted health seeking. The section on 'tiredness' as a systemic symptom showed the most number of inconsistencies (four participants; three diagnosed with COPD and lung cancer, and one with severe COPD) between those recorded ever feeling unexpectedly 'tired' in the questionnaire and those reported during the interview (see Table 4.1). This showed that participants were admitting to feeling more tired after a diagnosis, during the interview. It appeared difficult for participants to differentiate between "normal" tiredness and "abnormal" tiredness; where "normal" tiredness relates to the type of tiredness interpreted in light of everyday activities.

Two of the inconsistencies were prompted by spousal input:

PP 09: *"Hmm, well, I notice how, that you were more tired. I would just say you were getting more tired and when you tried to do any work. This is why we hadn't done very much this last year. Umm, very little work"*

PP= *Participant Partner (wife/husband)*

P 09: *"Yes, very little ... I didn't do a lot of work on the boat, did I, this last time?"*

Similarly, *Participant 04 (P 04)* did not recognise his tiredness as a symptom, which his wife had noticed becoming more prominent after his previous cancer treatment:

R: *"...I'd like to build a detailed picture of how your health was over the last 2 years so could you talk me through how it has been?"*

P 04: *"Over the last two years? Ahh, well, as far as, I've been getting breathless for a couple of years but umm, touch wood,..."*

PP 04: *"But you've been able to you know, be active"*

Study 1

P 04: *"Oh yea. But I've been, how do I put it ((pause))?"*

PP 04: *"a bit tired"*

P 04: *"Yea"*

PP 04: *"He's always a bit tired. But then I suppose maybe after having cancer and all that treatment maybe it's just you are that way. You know tired"*

'Tiredness' is a fairly common symptom in general and is not uncommon for a participant to experience it whether as a symptom of a pathology or part of normal everyday experiences of daily living. Participants find it difficult to distinguish increased tiredness from normal or expected tiredness caused by daily living or specific events. Accounts in which partners notice an increase in tiredness also suggests ambiguity in the term 'tiredness'.

One participant denied being tired on the basis that he wasn't doing enough to be tired, was discerned below (*Participant 09, age 76, diagnosed LC and COPD*):

"Tiredness? Well umm, well, I don't know how to actually differentiate between tiredness and sleepiness. I mean, I can get on the bus here. And I get sleep when I get off the bus at Southampton but I don't feel tiredness. I mean I can close my eyes, and nod off ((laughs)) sort of thing. So, I have to differentiate between tiredness, and sleepiness. Ahh, I don't do enough to be tired, really"

This internal conflict consistently observed in the interview could have explained his response denying any unexpected tiredness within the last 12 months in the questionnaire:

R: *"Can you in your own words describe to me what you, what do you understand from the word tiredness?"*

P 09: *"Umm, well, to me tiredness follows a period of activity which I've found sort of too energetic or prolonged and you know as a result I feel I feel I could stop umm, but that's about the only umm, if I don't undertake these activities then I don't have ah, unexpected"*

tiredness, do I? Sleepiness, yes but not tiredness. Again, again you see there's this difference between sleeping and tiredness. Umm, It could be that some of the sleepiness is due to the (ah) you know, the medication that I'm taking. That makes you sleepy. I mean, the Gabapentin one of the possible side effects is that you know, it makes you sleepy. (Participant 09, age 76, diagnosed lung cancer and COPD).

R: *"But you've taken Gabapentin before"*

P 09: *"For a long time"*

R: *"For a long time and this sleepiness only, you've only noticed it over the last, last couple of?"*

P 09: *"Yes, yes, the last few months"*

It was established that his sleepiness or need to take more naps during the day was unlikely to be caused by the effects of his long-term medication and this inconsistency was explored a little further. This ensured that the reason for the inconsistency was not due to poor interpretation or contextual difficulties but rather the questionnaire was not able to capture the relationship between exercise and 'tiredness', or distinguish between sleepiness and 'tiredness' for this participant.

Two of the inconsistencies relating to hot sweats were explained by the normalisation of sweats; everyday explanations were provided for their occurrence by **P 05** (*Participant 05, age 71, diagnosed severe COPD*):

"I just thought it's the fact that you know, I've got extra bedding and umm, I've just got too hot, dreaming about something like that you know, got myself into a bit of a state. It doesn't worry me really. It's not very often"

In addition, it was suggested that the hot sweats were episodic (*Participant 05, age 71, diagnosed severe COPD*):

"...It's not regular thing. It's not every night even once a week. Just every now and again"

Study 1

One participant (*Participant 07, age 71, diagnosed LC and COPD*) described the hot sweats in relations to menopausal hormonal changes and bouts of pneumonia:

“I was sweating, the bed was wet through, but then the doctor thought I had a touch of pneumonia so that wasn’t a sweat as much. It was because of the pneumonia”

However, more recent accounts of hot sweats were also elicited in the same interview:

“Bout six, eight weeks ago, won’t you say (addressing spouse)? Yea. When the weather started getting really cold, yes. I could sweat then some nights”

Exploring the discrepancies within the interviews confirms that there is a natural tendency to obtain more data or draw out details from an interview than a questionnaire. Generally, with interviews, details of particular symptoms of interest can be probed more easily with questioning.

Theme 2: Joint interpretation of symptom experiences between spouses

Spouses of participants were present in three interviews (**P 04, P 07, and P 09**). There was an element of joint interpretation of the symptom experiences observed and described in theme 1, where spouses were present during the interviews. For the most part, spouses supplemented additional information to the symptom accounts which would have been forgotten otherwise or left out by the participants, and on occasions these interactions elicited accounts of symptoms that had not been recorded on the questionnaire:

R: *“Did you notice anything else changed about your health then?”*

P 09: *“Not that I can remember. Do you @looking at wife@?”*

PP 09: *“No, but just that he was very lethargic” (wife of Participant 09).*

Participants with spouses also admitted to not noticing changes in their health experiences as much as their spouses might do, which could have resulted in some of the respiratory symptoms not being recorded in the questionnaire:

R: *"Was it fairly easy for you to recall experiencing any of the symptoms? Did you have any trouble recalling?"*

P 04: *"No. I mean obviously my wife notices difference in me that I don't but ah, ((laughs))"*

This suggested that the changes to the health were not perceived exactly as symptoms by participant themselves but were noticed by their spouses, which resulted in a non-affirmative response to symptoms on the questionnaire.

Theme 3: Differences between questionnaire-elicited (pre-diagnosis) and interview-elicited (post-diagnosis) symptoms

All but one case consisted of symptoms being reported in the interview (post-diagnosis) but were not recorded in the questionnaire (refer Table 5.1). The one exception in this direction of discrepancy; **P 07** (*Participant 07, age 71, diagnosed LC and COPD*) had denied experiencing any symptom of 'breathlessness' during the interview but she had reported ongoing history of chronic respiratory problems that required long-term inhalers:

R: *"Any breathlessness, or breathing changes?"*

P 07: *"No, not really. Not really, no...I've been going over to see my GP and she'd asked me about my chest. Because I've been sort of, bronchitis and all that, previous years... a couple of years ago, I had pneumonia. Only a slight pneumonia. And then, for that I had pleurisy, years ago, and then she asked me if I was still using my asthma pumps because now I've got four, asthma pumps"*

As Table 5.1 suggested, symptoms were more likely to be elicited in the interview than by the questionnaire (pre-diagnosis).

For instance, **P 06** (*Participant 06, age 65, diagnosed LC and COPD*) revealed that she was experiencing more symptoms than those previously recorded in the questionnaire. During the interview, she mentioned experiencing chest aches/ pain, breathlessness, and fatigue but these changes to her breathing was not recorded in the questionnaire:

Study 1

"I can be talking like this, and get out of breath and like ((gasp)) but I never used to be doing that. I used to be doing it automatically through my nose and everything used to be just carrying talking like that"

When asked whether participant had experienced any aches the back, shoulders or joints, participant recalled experiencing severe pain around her shoulder blade:

R: *"Any aches, or pain in the back, shoulders or joints?"*

P 06: *"Sometimes, you see, that's another thing I don't even think about. I mean, I get.. Sometimes I think it's my bed 'cause I'd wake up and I think, what you've got to remember is I'm 64. 65 coming up it's not like I'm 20. You know, you'd expect to get some sort of umm things, but, if I get a backache which I do quite frequently ... But sometimes I'd be out walking and I'd be ohhhh I've got to sit down, got to sit down, back's killing me but it's right up here"*

P 06: *"it's like in two patches like that, two distinct patches, right there and there. And I'd think ooohh. But if I could sit down, and I'd sit down 5 minutes, like this and it goes off and I get up and it's all gone... it's fine. Gone again. But that could just be **old age**"*

These discrepancies were further explored in the cognitive interviews to evaluate participants' level of understanding of the questionnaire items.

Although the episodic pain described was quite severe, the participant made several references to the pain as part of the ageing process, and related it to a prior history back pain.

It was also possible that participants were recalling more symptoms in light of the diagnosis, which would suggest an element of recall bias, a widely acknowledged limitation of retrospective data collection. This could be an implication of retrospective collection of symptoms as a result of reinterpreting accounts of symptoms following diagnosis:

P 06: *"It's right up high you see. And I've said, I've said. See, I've noticed this over the last few months. And I've said to (daughter), see that's what the back ache is I bet I've got cancer at the bottom of my lungs, that's what it is. And I think its spread. That's why I thought*

its spread all over me. But, you know I don't get it all the time. Every now and then I get it"

R: *"What made you relate to the pain in your back?"*

P 06: *"Cause I thought about my lungs, I thought I got lung cancer. And it's at the bottom of my lungs, and it feels like it's right inside, inside my ribs, it was at the bottom of my lungs it feels, I'm assuming my lungs are still right down here still. This is what I said to (daughter), it's right at the bottom of my lungs, it is. That's what it is, it's that cancer it's right at the bottom of my lungs, that's why I get the pain there,..."*

Furthermore, having gone through the process of medical investigations and clinical consultations, it is not inconceivable that participants become more perceptive of health experiences that had been earlier overlooked.

Discrepancies in breathing changes were also observed. It appears that some of the breathing changes were mild changes on exertion, and in light of an underlying condition, were not captured by the questionnaire (breathlessness was reported in the interview but not in the questionnaire following diagnosis):

R: *"When did you first notice you were breathless?"*

P 06: *"... I was **a bit breathless** I suppose and then they just automatically assumed because I was smoking it was COPD. Minor. And then I started off on the blue ventolin and when that wasn't doing as much for me"*

R: *"Could you describe this breathlessness for me?"*

P 06: *"It's **not**, well, it is more breathless when I'm out **walking the dog** and things like that, you know what I mean? Or straining, or if I'm doing the **hoovering** and things ... Maybe my breathlessness is not really breathlessness"*

5.3.3.3 Content validity of the questionnaire in a population consisting of high rates of chronic respiratory diseases referred to lung-shadow clinic

Theme 1: Interpretation of questionnaire response items

Generally, participants did not report any problems understanding the questionnaire when interviewed. The majority of them did not recall any difficulties filling out the questionnaire:

"It wasn't difficult at all" (Participant 01, age 62, diagnosed mild COPD).

"They are quite straight forward questions" (Participant 04, age 71, diagnosed COPD).

"Yes, it is straightforward" (Participant 03, age 52, diagnosed COPD).

Analyses of the cognitive interviews suggested that most participants were able to interpret the questionnaire items appropriately to suggest sufficient understanding. Verbal dialogues of interpretations of questions relating to chest, upper body or shoulder pains included:

R: *"Discomfort or pain that is not brought on by physical activity?"*
(R: Researcher)

P 10: *"Well, just like when you're sat there, you're in pain. Umm, and you're not actually doing anything physical"*

R: *"Question one is asking have you ever experienced any discomfort in chest, upper body or shoulders, and you ticked no. What does discomfort in your chest, upper body or shoulders mean to you?"*

P 06: *"To me it means whether I've had any sort of pain or whether, like my wheeziness, and my chest, all that, is caused me real, like the chest pain or discomfort because it's uncomfortable when you can't breathe properly... or sometimes you can cough really, really, and it hurts your chest. I mean I know when I have got a chest infection it hurts, it's like razor blades in my chest and when I cough, and you think ooooooh, that, you can feel that, you know"*

In one cognitive interview, *Participant 05* did interpret the statement in the questionnaire “... ache or pain in the centre of chest or ribs...” to mean heart-related problems, and resultantly, answered ‘no’:

“In the centre, well just sort of heart problem. If there was any heart problem” (Participant 05, age 71, diagnosed severe COPD).

This was the only exception, and it demonstrates the implication of how a participant’s interpretation of a symptom descriptor might influence response. Despite this, collectively, findings of the cognitive interview showed that items of the questionnaires were understood and interpreted appropriately by this population.

Overall, the symptom descriptors of the IPCARD questionnaire were able to consistently record symptom experiences described by the participants to match accounts of symptoms from the interviews.

Two of the 10 participants reported experiencing haemoptysis; the main symptom concern which led to their referral for a chest X-ray. The two participants were diagnosed with different stages of COPD (**P 01** and **P 05**):

*“Well that was one of the reasons why he sent me to the hospital because I was coughing up **just tiny bits of blood**” (Participant 01, age 62, diagnosed mild COPD).*

The other participant described her one episode of haemoptysis as “... clearing my chest, coughed up **some blood** which umm ... which really sort of made me panic quite a bit... It’s just one lot” (**P 05**). This matched her response of ‘coughed up mostly blood (blood with little or no phlegm)’ in the questionnaire:

“... Just a piece about 50 pence. Maybe a lil’ bit bigger” (Participant 05, age 71, diagnosed severe COPD).

Descriptions of haemoptysis in the interviews corresponded well with respective responses in the questionnaire thus, further validating the questionnaire in capturing the distinction between types of haemoptysis (coughing up mostly blood or little blood in sputum).

Theme 2: Reasons provided for missing data

Reasons for non-response or missing data included accidental omissions because participant simply forgot to fill in some of the items. Some participants had overlooked the questions when filling in the questionnaire:

"I forgot to tick the no" (Participant 08, age 68, diagnosed LC and COPD).

"Probably an accident more than anything" (Participant 02, age 05, diagnosed COPD).

Most of the questionnaires were filled out completely and appropriately. Aside from the far and few symptom items missed out, only one questionnaire response presented with large sections omitted or a full incomplete page especially towards the end of the questionnaire. When probed during the cognitive interview, participant 02 (**P 02**) established clear understanding of what was asked of the question items. He then explained that limiting time was most likely a contributing factor to the outcome (*Participant 02, age 65, diagnosed COPD*):

"Yea, it can't be done in 10 minutes"

Theme 3: Symptoms and descriptors not captured by questionnaire

Muscle cramp was the only symptom experience described in interviews that had not been captured by the questionnaire:

"I get cramps" (Participant 04, age 71, diagnosed COPD).

There was no medical reason or diagnosis to explain the cause of the muscle cramps as the participant had not consulted his GP or other medical profession regarding it. Another description of discomfort associated with muscles, not recorded by the questionnaire, was found in another participant relating to the back (thoracic spine region):

"... it's this two bits of stitch and you've got to hold it" (Participant 06, age 65, diagnosed LC and COPD).

Therefore, it was not a common symptom experience in this interview population. .

One of the recurring descriptions for the type of cough experienced was “*dry cough*”, which was mentioned in five of the 10 interviews. At present, the questionnaire does not contain a descriptor using the terminology ‘dry’ cough as a statement to describe recent cough(s). However, alternative cough descriptors to depict similar cough experience can be found in the questionnaire, and were used by those describing a dry cough in interviews, which had been further discussed in the following section.

Theme 4: Changes to existing cough.

All 10 participants reported having experienced some form of cough that had lasted longer than three weeks in the last three months prior to an appointment in the lung clinic.

The type of cough described in the questionnaire was found to be fairly consistent with the participant’s cough experience described in the interview. For instance, the item response ‘*A cough without phlegm*’ corresponded well with all five participants’ description of ‘*a dry cough*’ in the interview. Although the commonly used term in the interview, ‘*dry cough*’, was not recorded in the questionnaire, an alternative term ‘*a cough without phlegm*’ was ticked in the questionnaire. This suggests that the questionnaire was able to pick up dry cough.

However, for some of the participants, the occurrence of changes to their coughs over time was described during the interview but not in the questionnaire. Coughing was often associated with smoking:

“It’s what I call smoker’s cough. Umm, Cause I started smoking when I was 17” (Participant 09, age 76, diagnosed LC and COPD).

The changes in cough were also associated with changes in smoking behaviour. Four participants had reported noticeable improvements to their coughs in relations to quitting smoking:

“But since I’ve stopped smoking. Umm, I hardly cough at all now” (Participant 05, age 71, diagnosed severe COPD)

“There’s been when I’ve given up smoking; I never had a cough or anything. All the symptoms have gone away. And there was no sign of

Study 1

anything and I begin to wonder obviously, it's the smoke that's irritating my chest rather, than anything else" (Participant 06, age 65, diagnosed LC and COPD)

"Yes. I think it was less than a year ago because while I stop smoking during the summer it didn't seem so bad as I remembered it but maybe that's just wishful thinking. I don't know" (Participant 09, age 76, diagnosed LC and COPD)

Clear distinctions of the characteristics of their cough experiences over time were highlighted:

"The cough that I had then, was you could hear it, ohh, like I call it, an old man's cough. And it's makes you @pretend to cough@ You can't stop, you can't stop ... It's not like that anymore. Now's it's like a very dry. Like a (whoopany), it's a real dry cough, where you, you know, it's an entirely different cough altogether. One's a (whoopany) cough and one's like a congested cough. It was before but now it's not" (Participant 06, age 65, diagnosed early LC and COPD).

"It's a dry cough. There's nothing there. No phlegm there and it's just a little cough you know, not sort of a hacking cough I had before" (Participant 05, age 71, diagnosed severe COPD).

Despite the changes clearly expressed in the interview, analysis indicated that some of the participants proceeded to record the type of cough they used to have, in the questionnaire, rather than their more recent coughs. It is possible that the individual participant no longer interpret the current cough as a symptom, due to the perceived improvement from a productive cough to a non-productive one, as a result of stopping smoking.

The questionnaire did attempt to capture changes in the cough through an open-ended question; with two participants describing changes relating to their cough in this section.

5.4 Discussion and Conclusion of Study 1

The semi-structured interviews indicated that questionnaire completion did not appear to cause anxiety or raise concerns for participants. Furthermore, participants did not explicitly associate the study, or the questionnaire with lung cancer. The questionnaire had been designed to avoid generating unnecessary anxiety amongst participants that might have lung cancer. Some of the participants reported inadequate time to fully complete the questionnaire before the clinic appointment, but 74.2% of those eligible completed the questionnaire (reasonable completion rate for a self-completed questionnaire).

All of the 10 respondents interviewed reported multiple symptomology in the questionnaire (see Table 4.1). Four of the ten symptoms recorded by the questionnaire were sometimes reported in interviews post-diagnosis when they had not been recorded by the questionnaire pre-diagnosis (refer Table 5.1); tiredness, breathlessness, pain/ache/discomfort, hot/cold sweats, and cough, experienced before diagnosis were not always captured by the questionnaire. One possible reason for this inconsistency was due to joint interpretation of spouses; where spouses recognised changes in the health of their respective spouse, which prompted the recollection of a symptom during the interview. Some of the symptoms were not elicited by the participant themselves but by the spouses. However, these observations were limited to only three couples in the sample, and therefore, cannot be generalised to all cases where symptoms were presented in the interview but not recorded by the questionnaire.

Alternatively, participants might have been more likely to interpret a bodily sensation as a symptom in light of a diagnosis. Furthermore, considering the participants had just been through the process of receiving a new diagnosis at the time of study participation, participants could have been made more aware of symptoms such as breathing changes that may have been elicited during clinical consultations. Interviews were all carried out retrospectively. The bodily sensations not captured by the questionnaire were often mild or ambiguous.

Breathlessness and fatigue, the two symptoms most often not recorded by the questionnaire, tended to be mild changes on exertion, episodic, or 'normalised' by participants. Episodic and severe symptoms of pain/aches

Study 1

were also 'normalised', and not captured by the questionnaire in two participants. This process of 'normalisation' attributed potential lung cancer symptoms to normal processes of everyday causes, and not an indicative of ill-health (Brindle et al. 2012; Molassiotis et al. 2010; Levealahti et al. 2007; Corner et al. 2006).

It is not unusual for this population to be experiencing ongoing chest symptoms relating to their existing chronic respiratory disease, which makes it more difficult to identify any mild changes in breathing, if there were any. Qualitative interviews provide the opportunity to probe for more details, and ask follow up questions. For a minority of the participants, this led to changes in breathing being elicited.

Fatigue was commonly reported by participants in this study regardless of diagnosis. It is also a common symptom in the general population. Therefore, fatigue as a generic symptom might have low specificity and the use of the questionnaire to prospectively record the symptom 'tiredness' might be favourable. As the findings suggested, perhaps 'tiredness' could be a symptom that is more likely to be elicited in light of a diagnosis.

Hot sweats (unreported in the questionnaire pre-diagnostically) experienced by the two participants were characteristically intermittent and normalised within accounts of common explanations for the change in health experience (menopause, pneumonia, and night terrors). One of the two unreported cases of hot sweats also revealed a recent change in the pattern of hot sweats that occurred in the last three months leading up to diagnosis in the interview process. This hidden symptomology was likely elicited from further probing of the normalised accounts of everyday life experiences. Non-specific, episodic and non-progressive symptoms had been found to be normalised by patients with early-stage lung cancer (operable) who felt healthy and well (Brindle et al. 2012).

Progressive dyspnoea, dyspnoea that worsens with exercise, and/or persistent dyspnoea, and cough (chronic or intermittent; productive or non-productive) have been identified in the clinical literature as potentially suggestive features of COPD (Rabe et al. 2007). These cough descriptors were effectively recorded by the IPCARD questionnaire. Often, reports of changes in cough appear to be associated with smoking cessation. However, findings of the content analysis

suggest a possible need for the addition of closed questions to record *changes in cough*.

Overall, the questionnaire appeared to record the full range of symptoms experienced by patients with COPD and lung cancer. Most of the bodily sensations and health changes normalised by participants in interviews were recorded by the questionnaire. However, in a few cases, probing within interviews, and/or involvement of the spouse (supplementary input) further elicited normalised changes in health and bodily sensations. The qualitative analysis established the content validity and acceptability of the questionnaire in a population of lung-shadow clinic attendees with high rates of respiratory problems.

Chapter 6: Study 2

6.1 Introduction

The purpose of Study 2 was to investigate the feasibility of symptomatic diagnosis of lung cancer in a population referred to secondary care with high rates of chest and respiratory diseases. Previous studies on symptomatic diagnosis of lung cancer had been limited by methodological weaknesses such as retrospective study design and unsystematic data collection leading to potential bias and confounding (see Systematic Review Chapter; Section 3.5.3). Earlier UK-based primary care research (Hamilton et al. 2005; Jones et al. 2007; Hippenley-Cox and Coupland 2011) relied on secondary sources such as medical and GP records and databases to collect symptom data, which are subjected to recording bias by clinicians. Other quantitative studies that collected data directly from patients either through questionnaires or interviews were retrospective in study design (Hoppe 1977; Kubik et al. 2001); that is, reports of symptoms were obtained following diagnosis, potentially introducing recall bias. One of the studies was prospective and used a standardised MRC questionnaire (see Appendix 4 for MRC questionnaire). However, this was an interviewer administered questionnaire, and recorded common chest symptoms only (Kubik et al. 2001).

The current study collected prospective, symptoms, risks and co-morbidity data using a patient-completed symptom questionnaire. This method of data collection allowed a wide range of patient-elicited symptom experiences to be collected before a diagnosis is known. As such, the results of this study would not be affected by recording bias or bias resulting from retrospective data collection. Despite the strengths of the study methodology, there was a proportion of missing data because a patient-completed questionnaire was used. 'Missingness' is an inherent feature of participant-completed questionnaire data that needs to be addressed by the methods used to analyse the data. The following chapter present the methodology, methods and results of the analyses to estimate the discriminatory value of symptoms for lung cancer diagnosis, taking account of missing data.

Study 2

Objectives:

- (1) Explore the feasibility of using patient self-reported symptoms to identify lung cancer in a secondary care population (a population with high rates of chronic respiratory disease that has been referred to lung shadow clinics on suspicion of lung cancer).
- (2) Identify patient self-reported symptoms that are independently associated with a diagnosis of lung cancer in a secondary care population with high rates of chest and respiratory diseases.
- (3) Identify patient self-reported symptoms that are independently associated with the diagnosis of lung cancer in a secondary care population with COPD.

6.2 Methodology Section

Missing data

6.2.1 Introduction

Incomplete data due to participant non-response, or invalid response, is a characteristic of most studies involving surveys or participant-completed questionnaires (Allison 2002; Wood et al. 2004). Missing data may also be caused by non-coverage, participant withdrawal, and/or data entry errors. As follow-up data were extracted from patient medical data, outcome and clinical data were mostly complete in the IPCARD Chest Clinic Study dataset. However, a strategy was required for handling participant non-response and erroneous questionnaire data.

Missing data are a problem because it can result in bias and/or loss of power. There are two main approaches to handling missing data: i) complete case analysis (using a dataset of fully observed cases); and ii) analysis of partially observed data using imputation with either single or multiple imputation techniques. To use complete case analysis is to ignore missing data in the analysis, which implicitly assumes that the data are missing completely at random or MCAR, an assumption that is often hard to validate. Therefore, complete case analysis not only result in loss of power but also, possibly

Study 2

introduces bias if the missing data are not MCAR data (MCAR will be defined below). The latter method (imputation) can avoid bias and loss of power, but requires the identification of appropriate analytical methods that do not introduce bias into the dataset and are feasible with data analysis software.

There is no way of confirming the true reason for missing data in each questionnaire item. However, the extent of missingness and its possible impact on the final analysis can be considered when deciding how to handle the missing data in the primary analysis (Little and Rubin 2002).

6.2.2 Rationale for missing data methodology

The implications of missing data for the analyses are dependent upon the missing data mechanism. Missing data mechanisms have been classified according to three mutually exclusive categories (Rubin 1976; Little and Rubin 2002); Missing Completely At Random (MCAR), Missing At Random (MAR), and Missing Not At Random (MNAR). Whereas most methods of analyses would still be appropriate for MCAR data although with a potential loss of power, for MAR and MNAR data, complete case analyses might result in bias, and invalid inferences.

For data that are missing completely at random, MCAR, the only mechanism for missingness is that of randomness. When data are MAR, conditional on the fully observed variable x (covariate), the missing observations on y (variable) are missing completely at random i.e. *“if conditional on this fully observed variable, we assume the chance of seeing the partially observed variable does not depend on its values, the data are MAR”* (Carpenter and Kenward 2008). A further implication of MAR is that the distribution of potentially missing data is the same (conditionally) for all units or participants with the same observed data (Carpenter and Kenward 2008). Therefore, using appropriate analytic methods, the partially observed data can be used to produce unbiased estimates, and the mechanism defining the value of missing observations can be ignored.

Under MNAR assumptions, the reason for missingness depends upon the unseen observation; that is, the probability of missingness depends upon the

Study 2

unobserved value and the value of missing observations cannot be inferred on the basis of the observed information (Carpenter and Kenward 2008).

Therefore, the missingness mechanism is non-ignorable and valid inferences require a model of the data and missingness mechanism. However, when data are MNAR, the missingness mechanism is rarely known.

Under MAR assumptions, particular analyses that ignore the missing value mechanism are valid; a valid analysis can be achieved within a Bayesian or likelihood framework. Some frequentist methods traditionally used under MCAR assumptions can also be modified for MAR, for example, through the use of weighting with estimating equations (Molenberghs et al. 2004; Carpenter and Kenward 2008).

Even though under the assumptions of MAR, the missing data mechanism (i.e. the mechanism for the relationship between missingness and the value of missing observations) is ignorable as conditionally, the distribution of potentially missing data is the same for all participants with the same observed data, it is not possible to determine from the data whether the missing observations are MCAR, MNAR or MAR (Mallinckdrot et al. 2013; Carpenter and Kenward 2008). However, it is possible to distinguish between MCAR and MAR. For data MCAR, the proportion of participants with data missing will not vary with the observed covariates. The identification of covariates that are associated with missingness, would be consistent with MAR. However, even when observations may appear consistent with the data being MAR, MNAR cannot be ruled out (Molenberghs et al. 2004). Therefore, goodness of fit of the imputed model to observed data cannot be used to confirm the validity of the model (Molenberghs et al. 2004, Carpenter and Kenward 2008). That is, the missing value mechanism for observations that are missing not at random cannot be ignored. Furthermore, as under MAR assumptions, conditionally on fully observed covariates or independent variables, missing observations are missing **completely** at random, for most datasets the true mechanism is probably MNAR (Carpenter and Kenward 2008). As MNAR cannot be disproved, there is a general consensus that sensitivity analysis is an important part of the modelling process where imputation is involved (Molenberghs et al. 2004; Carpenter and Kenward 2008). Sensitivity analyses identify how inferences vary under assumptions of MAR, and various MNAR models (Carpenter and Kenward 2008), where several statistical models are considered simultaneously,

120

Study 2

allowing the implications of hypothesised models to be identified (Molenberghs et al. 2004; Carpenter and Kenward 2008; Carpenter et al. 2007). Sensitivity analyses might be informed by expert opinion regarding potential missing data mechanisms. For example, in this study, the researcher's experience of responding to participants' questions about the questionnaire items, and probing reasons for missing data with participants, might identify situations where the probability of missingness appears to depend upon the unobserved values of the variable of interest, and inform MNAR models. However, sensitivity analysis cannot identify the correct MNAR model; MNAR models informed by expert opinion, are a transparent way of expressing potential deviations from MAR that provide a point of departure for 'principled' sensitivity analyses (Carpenter and Kenward 2008).

Even when data are MNAR to some degree, methods valid for use with MAR data may still be of use (Carpenter and Kenward 2008; Mallinckdrot et al. 2013). For example, the information provided by the 'unseen data' might be negligible. Rubin et al. (1995) stated that "*it quite often happens that after accounting for the information about the missingness mechanism in the observed data, there is relatively little information remaining in the unseen data*". Furthermore, as the missingness mechanism for MNAR data are rarely known, a "shift" to MNAR assumptions and an MNAR analysis, does not provide greater validity; it is not possible to identify the correct MNAR model (Carpenter and Kenward 2008, p.20). Where the observed data are at least partially consistent with MAR assumptions, the use of appropriate analytic methods that ignore the missing value mechanism (e.g. multiple imputation) combined with sensitivity analyses, might be appropriate.

Although there is a consensus that sensitivity analyses are advisable when analysing data under MAR assumptions, there is also a recognition that sensitivity analyses can be beyond the scope of small applied research project (Carpenter and Kenward 2008). Carpenter and Kenward (2008) recommend that any analyses on the imputed (partially observed) data be presented alongside analyses on complete cases (those participants or units with no missing data) so that conclusions can be compared, and explanations provided for any differences. However, as complete case analysis is only valid if the

Study 2

missing data mechanism is MCAR, the analyses might differ because the missing data mechanism is MAR, and the complete case analyses produce biased results. In this case, Carpenter and Kenward (2008) recommend that the ways in which the mechanism appears to depart from MCAR, and how this departure affects the complete case analyses should be explained. Such explanations might also address the potential for bias in multiple imputation procedures, where differences between complete case and multiple imputation analyses do not appear to be explained by visible departures from MCAR.

Considering the time and budgetary constraints of applied PhD projects, in this case a small feasibility study, sensitivity analyses are not realistic. For the analyses presented in this thesis, where the partially observed data appears consistent with MAR, a likelihood based method (multiple imputation) will be used to analyse data. Prior to regression of independent variables on cancer diagnosis, the distribution of variables in the complete case and imputed datasets will be compared to check that any differences between the distributions of the observed and imputed variables are sensible within the context of this study (Stata 2013). Furthermore, models resulting from complete case analyses will be compared with analyses of the partially observed (imputed) dataset to identify any differences; where appropriate, differences will be explained in relation to visible departures from MCAR. More generally, the potential implications of missing data that do not appear consistent with MAR assumptions, for the interpretation of results, will be discussed in the results chapter.

6.2.3 Multiple imputation (MI)

MI allows inference on a complete-data statistic, by fitting a complete-data model to the observed data (Little and Rubin 2002). The following section explains the theory underpinning MI. There is always the risk of getting invalid results with any imputation process, which can be minimised through careful checks, and a good understanding of the chosen imputation method.

In a fully parametric model, it is possible to calculate maximum-likelihood estimates from the incomplete data using specialised statistical methods such as expectation maximization (EM) algorithm. It could be argued that such procedures may be somewhat more efficient than MI because they involve no

Study 2

simulation, but this current study is non-parametric. Furthermore, this is applied statistics, where missing data are nuisance rather than the primary focus of this research. Therefore, a practical, approximate solution with good properties can be preferable to one that is more efficient but complicated to implement.

Theory of MI:

Based on conditions from a Bayesian posterior distribution of missing data (Bayes' Theorem), MI combines the use of various statistical techniques (EM/ maximum likelihood estimation, propensity score estimation, and Markov Chain Monte Carlo (MCMC) model) to generate estimates of the incomplete data (Rubin 1987; Schafer 1999). Bayes' theorem is a theory of probability; the likelihood of something in the event of new evidence.

MI is a three-step approach to estimating incomplete data regression models based on Rubin's Bayesian paradigm (Rubin 1987): 1) the imputation step, 2) the complete-data analysis or estimation step, and 3) the pooling step.

It starts with a data augmentation process, where possible values for missing observations are created that reflect the variability about the non-response/ missing model. With the assumption of MAR, these values are then used to replace or impute the missing observations. This imputation step is carried out at each repetition; $t = 0, 1, 2, 3, \dots, T$, in an iterative process (part of the MCMC model) (Schafer 1999), until 'm' sets of imputed datasets (completed datasets) are generated under a chosen imputation model. This is achieved through the iteration of the EM algorithm consisting of the expectation step and the maximization step to maximise log-likelihood function.

In the second step, these parallel datasets can then be analysed using the standard methods for complete datasets. This is also called the completed-data analysis step, the primary analysis to be performed on the completed dataset.

Finally, in the pooling step, results ('M') obtained from 'm' completed-data of these analyses are combined to produce a single multiple-imputation based estimate (denoted 'e' in Figure 6.1). Figure 6.1 below illustrates the process of MI.

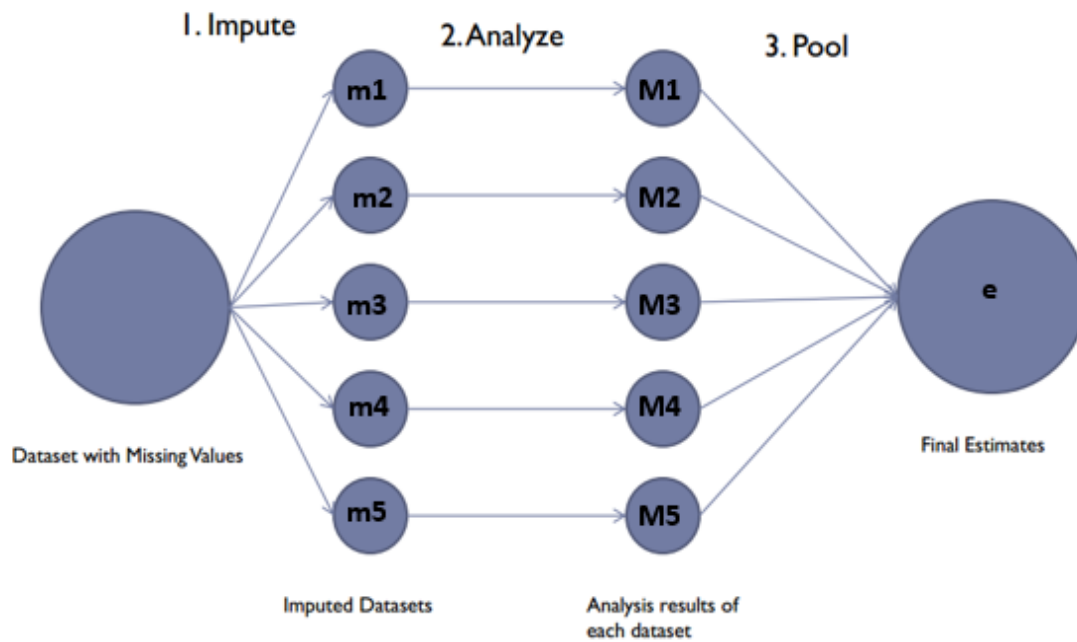


Figure 6.1: Illustration of multiple imputation process (diagram modified from Humphries 2012)

There are two main methods to imputation: multivariate normal (MVN) model and multiple imputation by chained equations (ICE / MICE). MVN uses a joint modelling approach based on a multivariate normal distribution model by Rubin (1987) (Schafer 1999). It assumes that the variables follow a normal distribution and is intended for continuous variables only. In some cases, it can be transformed to model binary and ordinal variables. However, unjustifiable assumption of normality in non-normally distributed variables may introduce bias (White et al. 2011).

MICE, also known as sequential regression multivariate imputation (Raghunathan et al. 2001), uses a Gibbs-like algorithm to impute multiple variables sequentially using chained univariate fully conditional specifications (FCS) of prediction equations (chained equations) (White et al. 2011). This variable-by-variable method of imputation specifies an imputation model per variable, makes it suitable to handle incomplete datasets with arbitrary missing data and have larger number of categorical variables.

All procedures were implemented using Stata version 13.0. MI was performed using the built-in MICE package in Stata (Royston 2005). The iterated chained

Study 2

equation process requires specifications of conditional models for each incomplete variable in relation to all other variables. Stata converges the estimates across imputations by the command `-mi estimate-` for analysis.

6.2.4 Identifying symptoms that predict lung cancer diagnosis

Presented with a large dataset of possible independent variables (113 independent symptom variables recorded using the IPCARD questionnaire), the aim of the study was to find out which combination of them are most useful for predicting LC diagnosis in this population of lung-clinic attendees.

The most natural use of multiple regression is when the outcome variable is **binary**, and when they are many independent variables (categorical, ordinal, or continuous) involved. The assumptions made in logistic regression are different to that of linear regression for this reason, and will be discussed later.

The logistic regression model is one form of a generalised linear model that uses maximum likelihood to estimate model parameters. The logistic model does not require the data to be normally distributed data, nor does it assume a linear relationship with the dependent variable (Agresti 1996). This makes the logistic regression model suited to handling real-life data and a well-established technique readily used in medical research (Bland and Altman 2000). The model uses the concept of “odds”, which makes it easier to relate within a clinical context. The odds of an event occurring is defined as:

$$\frac{\text{Prob (event occurs)}}{\text{Prob (event does not occur),}}$$

where *Prob (event does not occur)* is equal to $1 - \text{Prob (event occurs)}$.

The IPCARD questionnaire also gathers information on co-morbidities and risk factors. These epidemiological risk factors precede lung cancer and its symptoms in the causal pathway, and a relationship between the symptom and lung cancer that is independent of the risk factor would not necessarily be anticipated. In this case, the inclusion of the risk factor in a multivariate model, whilst potentially accurately predicting LC diagnosis, might obscure the relationship between the symptom and lung cancer. As the project aims to

Study 2

identify symptoms that predict lung cancer diagnosis, and there are clinical criteria for LC investigation that are satisfied by the occurrence of potential symptoms, but not by the occurrence of risk factors alone (see Systematic Review Chapter 3, page 41), two separate models will be developed to predict LC diagnosis; 1) Lung cancer and symptoms, adjusting for any confounders (age, gender, and/or comorbidities), and 2) Lung cancer, symptoms and risk factors, adjusting for any confounders.

Where risk factors are also comorbidities or lifestyle factors that might cause current symptoms similar to LC at baseline/questionnaire completion and therefore, be confounders (asbestos related illness, COPD, pneumonia within the last 3 months and current smoking) the potential for confounding is investigated and clinical confounders are adjusted for in model, where appropriate.

However, interaction terms between symptoms were not tested in the models due to the exploratory nature of the study analysis. Given the sample size, any analysis of interaction effects in this study would be underpowered, so it would be possible to miss interactions that were there. It would be difficult to make effective conclusions about interactions based on the current analyses alone. Interactions between variables in a model can cause inaccuracies in the estimates of individual coefficients and their accompany variance, an issue of collinearity (Robins and Greenland 1992). Not including interactions if they exist, would produce a less accurate statistical model. The implications of not testing for interactions will be reflected upon as part of the study limitations (Section Six).

6.2.5 Method of variable selection

The end goal of building any predictive model should be to create an accurate statistical model that can be applied to future, external data that is not too large for practical use. In order to achieve a parsimonious model, we need a subset of independent variables that completely explains the outcome (fit) with as few variables as possible because every irrelevant predictor decreases the precision of the estimated coefficients; to achieve a parsimonious model (Beale 1970; Vittinghoff et al. 2005; Faraway 2015).

Study 2

It is not uncommon for model builders to be overly concerned with creating models that only contain statistically significant predictor variables (Greenland 1989). However, the p-values should not be applied too literally as it is possible to miss out on clinically important variables that could be useful to the predictive model (Vittinghoff et al. 2005; Faraway 2015). As such, non-significant variables of clinical interest to the prediction of lung cancer could still be included in the models developed in this study. Multiple testing in regression generally adds uncertainty to the validity of the p-values. Removing less significant predictors tend to increase the significance of the remaining predictors, which could over-inflate their importance (Vittinghoff et al. 2005; Faraway 2015).

One of the drawbacks of the stepwise procedure is that the statistical significance of individual variables is not directly correlated to the predictive model's overall accuracy and so, may not really help solve the problem of interest. When variables are dropped from the model, it does not mean that they are completely unrelated to the outcome. They can still be correlated, just that they do not provide additional explanation beyond those variables that are already in the model (Faraway 2015).

Forward stepwise logistic regression was used in this study to find the best fitting model. This approach includes or discards variables based on the '*F*' statistic of models. It requires two significance levels; one threshold for entry ($p < 0.05$) and another for removing variables ($p > 0.1$). Stepwise regression adds the variable that has the highest partial correlation with the outcome variable whilst considering all the variables in the model, at each stage of the analysis. The resultant aim is to find the set of independent variables, which maximises the '*F*' statistic (Stata 2013).

A possible limitation to stepwise in this dataset relates to its use in a small sample. Logistic regression requires quite large sample sizes because it uses maximum likelihood estimations, which are less powerful than least squares estimations (used in linear regression). Generally, at least 10 cases are recommended per independent variable. With insufficient data for each variable, the stepwise procedures may fit the randomness that is inherent in most datasets and generate tenuous models (Stata 2013; Faraway 2015).

Study 2

Despite the theoretical concerns against the use of traditional stepwise regression, there is no clear consensus on which variable selection method works best and stepwise regression still appears to offer a quick and reliable means for creating reasonable statistical models.

6.2.6 Variable selection with an imputed dataset

Variables selection were performed using the `-stepwise-` estimation in Stata (version 13.0). The imputed dataset was then converted into a single stacked dataset (mlong data structure) and appropriate weighting by $1/M$ where M = number of imputations, was implemented with 'stepwise' (Woods et al. 2008). As currently, there are no guidelines on variable selection in imputed datasets, this method of weighting was intuitively simple and pragmatic. Woods et al. (2008) suggested that it is a reasonable way of performing variable selection in a multiple imputed dataset.

The stacking approach for MI variable selection results in a slightly overstated Type I error but retains more power compared to Rubin's approach (Woods et al. 2008). This is not to say that the stacking method is without limitations, but rather it was a sensible option for this study at the model-building stage.

6.2.7 Assessing fit of the model

To compare the models and decide on the best fitting model, the Akaike's Information Criterion (AIC) was applied after estimation of the imputed dataset (after each logit). The AIC is a comparative measure of fit that allows meaningful comparisons between two parallel models. The model with the lowest AIC value indicates the best fitting model (Akaike et al. 1998).

Generally, adding more variables in a model always improve the fit of the model. The AIC accounts for this by penalising larger models, which discourages over-fitting; making the AIC a better measure of goodness of fit (Akaike et al. 1998).

The performance of the set of diagnostic criteria will be assessed using the area under the receiver operator characteristic (ROC) curve. ROC analysis quantifies the accuracy of diagnostic test used to discriminate between two conditions; e.g. lung cancer or no lung cancer, using a graph of sensitivity

Study 2

plotted against (1-specificity) of the diagnostic test. The discriminatory accuracy of the test is measured by its ability to correctly classify those with the disease and those without (Pepe 2004; Altman and Bland 1994c).

6.3 Method Section

Missing data

6.3.1 Description of missing data

Missing data was described in each variable. When a variable had a high level of missing data, or relatively more missing values than the other variables, plausible reasons for the missingness were investigated; For example the possibility that it was a questionnaire design issue, or is it a particularly sensitive question, was considered. Where variables were missing a lot of data, the possibility that missing observations had anything in common was explored (MAR assumptions were investigated).

To do this empirically, the distribution of missing data by socio-demographic variables; i.e. age and gender, and by outcome variables (LC/not LC), was investigated. For this, a dummy variable was created for each variable of interest, indicating whether an observation was missing or not; 1 = missing and 0 = non-missing. Stata can automatically create a binary missing-value indicator (miss_variable) for each variable using the command, **gen(miss_)**. Cross-tabulations between 1) the miss_variable (cases of non-responses to a question or variable) and a socio-demographic variable, and 2) the miss_variable and outcome variable of interest (e.g. LC diagnosis or COPD diagnosis) were carried out to inform us if missingness in that response variable differs by age or gender, and whether the differences are large enough to affect the response variable i.e. which characteristics were more likely to be associated with missingness. Cross-tabulations allow for the inspection and comparison of differences among groups of variables with nominal or ordinal response data to be made (Field, 2009).

6.3.1.1 Missing data rates

The distribution of missing data in variables that had high missing data (10% or more of observations that were missing) were explored by demographic variables and the diagnostic outcome variables (COPD/LC). It may be that the reason for these high missing data was random and cannot be explained, which would qualify as MCAR data, as explained earlier. However, it may also

Study 2

be an indication of a lack of comprehension or clarity on that question for a particular sub-group of the sample population. Complete case analysis under the wrong assumption that the complete cases are representative of the partially observed data/missing values can lead to bias. This is an exploratory study, and dropping variables on the basis of high missingness, could risk excluding an important covariate in the predictive model.

In addition to variables with high missingness, missing data on the generic symptom variables; questions identifying whether an individual has ever experienced a particular generic symptom or not, were also descriptively analysed. Generic symptom variables form the section headings of the IPCARD questionnaire i.e. chest/shoulder pain, cough, breathing changes, tiredness, haemoptysis, chest infections, appetite change, weight loss/gain, voice, and skin changes. Within each of these sections, there are variables consisting of lay descriptors relating to each symptom, which provide more detailed information about the symptom. Analysing missingness in all the generic symptom variables would ensure the systematic coverage of all symptoms potentially indicative of lung cancer. Each of these generic symptom variables was dichotomised to either 'ever' having the symptom or 'never' (labelled: _ever). There are too many variables in the IPCARD questionnaire to descriptively analyse the missing data for every questionnaire item, where missingness was low.

The percentages of missing data in the generic symptom variables varied, ranging from 1.1% to 6.7%. A summary of the frequency and percentage of missing data for each variable is found in Appendix 11.

6.3.1.2 Missing data patterns

Statisticians discuss patterns of missing data in two forms: monotone and arbitrary or random (non-monotone) pattern. Figure 6.2 illustrates the two patterns of missingness.

Study 2

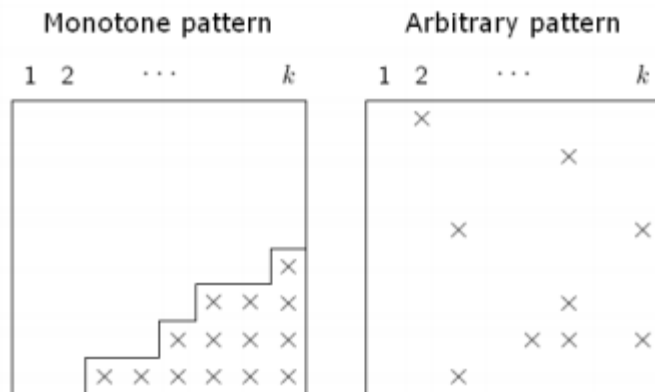


Figure 6.2: Illustration of monotonic and arbitrary (non-monotonic) patterns of missing data for ' k ' number of variables (adopted from Husmain 2008)

The importance of the pattern of missingness relate mainly to the type of approach to imputation in Stata. The process of imputation is much easier from a statistical viewpoint when the missing data are monotone. However, in practice, the pattern of missingness is hardly ever monotone (Little and Rubin 2002).

Arbitrary-patterned missing data will require more complex methods of imputation such as multivariate parametric model (Schafer 1999) and multiple imputation by chained equation approach (MICE) (Little and Rubin 2002; van Buuren et al. 1999).

Missing data were explored using the 'misstable patterns' command on Stata to identify whether the data followed a monotone or arbitrary pattern. With large number of cases and variables, it can be difficult to identify a monotone missing data pattern. As such, the 'misstable nested' commands were also used to examine the nesting structure of the missing values; patterns where a missing value on one variable is always missing on another variable.

6.3.2 Consistency check with MAR assumptions

Before building the imputation model, the assumptions for missing data must be checked. As previously explained, mechanisms of MCAR, MAR, and MNAR missing data cannot be completely proven or determined as such from the observed data. It is, however, possible to test whether MCAR assumptions hold by exploring the relationship between the observed data and the occurrence of

Study 2

missing data (Marchenko and Eddings 2011; Carpenter and Kenward 2008). Therefore, the observed data can rule out MCAR (Carpenter and Kenward 2008). However, this does not prove that the data are MAR as it is still possible that the data are MNAR. The argument is based on the interpretation of the departures of the missing data from MCAR (Carpenter and Goldstein 2004; White et al. 2011).

Logistic regression analyses were carried out for each of the chosen variables; between the missing dependent dummy variable, `miss_var` and the chosen covariates for imputation as independent variables. This analysis helps to understand the missing data mechanism, which informs the best technique to use to address the missing data. If there is any relationship between the observed covariates and the occurrence of missing data in the chosen variable of interest, data are not likely to be MCAR as the missingness could be explained by other variables, which suggests that data are consistent with MAR (Little and Rubin 1987; Marchenko and Eddings 2011). Under MAR assumptions, there might be systematic differences between the missing and observed values, but these can be generally explained by other observed variables.

Where the observed data hold up to MAR assumptions, likelihood based method of imputation such as MI would be valid. However, as explained, when considering models to impute missing data, the tenability of MAR data can, inherently, never be definitely determined from observed data. One can only test the plausibility of MAR missing data (Allison 2002; White et al. 2011).

6.3.3 Imputation model building

The aim of imputation modeling is to build a model that preserves all the main characteristics of the observed data. The quality of the imputation, and the applied analyses that follow, depend on the quality of the imputation model.

There are a number of factors to be considered when building an imputation model, including, but not limited to, the choice of: imputation method; the variables to include (inclusion of predictors of missingness, outcome variables, and interactions); and the number of imputations to generate.

Study 2

There are many different approaches to the selection of the predictor variables, each with their own justification. In general, it is recommended to use as many variables as possible in the model, as the imputation model then uses all of the information in the dataset (Rubin 1996, van Buuren et al. 1999). Rubin (1996), Taylor et al. (2002), and White et al. (2011) recommend including all variables associated with the probability of missingness, in addition to the variables to be included in the analyses (variables already contained in the dataset).

It was not practically feasible to use all the variables in our dataset to predict missing values on a given variable due to the large number of variables. This approach may be unnecessary as many of the variables in the IPCARD dataset are likely to contain redundant information, making the model less effective. Instead, a subset of relevant predictors of missingness in each variable with missing data was selected based on clinical knowledge. The predictors included in the imputation model had to be clinically associated to the variable of interest (with the missing data), which can make the MAR assumptions more plausible. The choices of the predictors were also informed by earlier IPCARD studies. There were, however, practical issues relating to this method e.g. predictors which perfectly predict binary variables (perfect prediction problem), which had to be addressed. There is also the possibility of omitting a key predictor from the imputation model in this method, which could lead to biased estimates being generated.

Some imputation models include auxiliary variables that may or may not have missing values. These variables are supplementary variables within the original data that are correlated to the variables of interest but are not part of the analysis (Little and Rubin 2002). Auxiliary variables also make the MAR assumption more plausible as they can help to keep the missing process random (Little and Rubin 2002). Gender and age were included as auxiliary variables in this study. If the study is still in the research design phase, changes to the data collection process can still be made to include auxiliary variables that might help to inform missingness in a variable. No other post-hoc auxiliary variables were planned at the time of imputation.

Another element towards building an imputation model that preserves the integrity of the observed data, is to include design variables that represent the structure of the data in the model i.e. interaction variables; e.g. those with

Study 2

COPD and without COPD, and important risk factors. Whilst it would have been preferable to carry out a separate imputation for the COPD sub-group, or impute interaction terms with COPD, as the sample size for the COPD sub-group was small, this was not practically feasible within the remit of the PhD and will be acknowledged as a limitation in this study. The intention was to build a functionally reasonable model that allows for substantive analyses, which will be compensated by other diagnostic checks.

The outcome variable was present in the imputation model to obtain valid results. There were no missing values to the outcome variable as they were extracted from the patient medical records.

Perfect prediction often occurs in the analysis of many categorical data, and may occur when variables are imputed using one of the above mentioned methods of imputation (logit, ologit, or mlogit). It can generate variables with large coefficients and large standard errors during estimation, resulting in new coefficients drawn to be largely positive or largely negative. Furthermore, in binary imputation variables, missing values may be all imputed as ones, or all zeros; presenting bias in the estimates.

Alternatively, eliminating the offending covariates that perfectly predict the outcomes could violate the theoretical underpinning of the imputation model. It may be reasonable for some research to omit the perfect predictors from the analysis, and have their inferential conclusions adjusted accordingly. However, this option is not recommended for studies that impute large number of variables, among which are many categorical variables. Therefore, perfect predictions were handled directly in the imputation process via the 'augment' option, used ad hoc for all categorical imputation methods (White et al. 2010). This approach computationally adds a few extra observations with very low weights in such a way that they have negligible effect on the results but prevents perfect predictions.

6.3.4 Imputation model and variable type

Fully conditional specification (FCS) multivariate imputation method to modeling was used in this study because it could accommodate a mixture of

Study 2

different types of variables, and preserve some of the important characteristics observed in real data (Lee and Carlin 2010).

Most of the variables in the IPCARD questionnaire were categorical variables except for five severity scale questions. Depending on the variable, one of the following univariate imputation methods were used: logistic, ordered logistic and multinomial logistic. For a binary variable, 'logit' was used. For categorical variables, 'ologit' was used to impute missing categories if they were ordered, and 'mlogit' was used to impute missing categories if they were unordered (Raghunathan et al. 2001). Data were set in wide form, where all the observations for a single subject are on the same line.

The variables included in the imputation model were chosen based on clinical understanding of the relationship between the symptom variables, and the relationship between the symptoms variables and comorbidities or risk factors, e.g. Q22d_BrChnges; breathing problems that require the use of an inhaler is thought to be related to a history of asthma, Q69c_Asthma. Therefore, using logistic regression analysis, miss_Q22d_BrChnges would be the dependent variable over Q69c_Asthma, the independent variable (predictor).

6.3.5 Number of imputations

The confidence interval of the estimate produced from MI is usually stronger than those of simple imputation methods (e.g. mean substitution) because the model would have considered variability due to sampling and variability due to imputation, which creates statistically valid inferences (Rubin 1987; Schafer 1999). Hence, the higher the number of imputations carried out, the better is the understanding of the variability introduced into the results. The recommended number of imputations (M) selected varies amongst statisticians. For small fractions of missing data, five to 20 imputations can be considered to be sufficient (Stata manual 2013). White et al. (2011) suggested more than 10 cycles are needed to achieve convergence. Rubin (1987) advocates that three to five imputations provide excellent results, Schafer (1999) also concur that no more than 10 imputations are required. With technological advancements making imputation more computationally friendlier and more feasible, statisticians are now recommending higher number of imputations. For this study, 20 sets of imputations were performed.

6.3.6 Convergence

Complex models such as ‘mlogit’ can fail to converge due to a large number of categorical variables, leading to small cell sizes. When convergence didn’t occur, the cause of the problem was identified by removing most of the variables to create a working model, and then adding variables back one at a time or in small groups until it stops working. For some variables, it might be feasible to remove the problematic variable from the model. However, where variables are potentially important to the imputation model (for example, as informed by clinical knowledge) it is difficult to justify excluding them as a result of non-convergence. That said, if a single model fails to converge, the imputation process will fail as a whole. To overcome this problem, dummy variables were used in which the missing values of the risk variables were assigned to a ‘missing’ category. The limitations to this method are acknowledged and discussed further in section 7.5 (Chapter 7).

The convergence of MICE is achieved by running multiple independent chains. It is likely that convergence will not be achieved when there are too many parameters in the model to be supported by the amount of observed data, causing instability to the imputation model.

6.3.7 Diagnostics

It is good practice to examine the sensibility of the imputed values before carrying out the primary data analysis. There are no set rules to diagnostics for imputations. Imputations can be checked using a standard of reasonability, where the differences between the observed and the imputed data values are checked to reveal any unusual patterns that might suggest problems with either methods, and to see if they make sense within the context of the study (Abayomi et al. 2006).

Any differences between the distributions of the observed and the imputed data were checked to identify any unanticipated discrepancies as recommended in Stata manual (2013). The full stepwise analyses of the observed and the imputed data were also compared.

Study 2

Complete case analyses of the original dataset with missing data were carried out, and compared with the imputed analyses, as recommended by Carpenter and Kenward (2008). Assuming that the data are MAR, complete case analyses would not be valid, as it would by theory, introduce bias. The comparison between complete case analyses and the imputed analyses allows discrepancies to be explained in relations to any visible departure from data that are MCAR, where appropriate. This process should at least flag up where the modelling may not be appropriate, if not where the missingness assumptions are not met (Abayomi et al. 2006). Essentially, the aim of the diagnostics is to identify potential problems, and fix or refine the imputation model.

6.3.8 Sample size calculation

Main model:

Power calculations indicated that, for this exploratory study, a model using 125 lung cancer cases (and approximately 2 non-cases for each case) would have a 80% power (2-sided alpha 0.05) to detect a difference in symptom frequency of 20% (from 20% in the non-LC group to 40% in the LC group). This was also justified to detect clinically significant effect sizes; odds ratio >2 or <0.5 . The identification of lower frequency symptoms strongly associated with LC (odds ratio >2 or <0.5), but with an alpha of >0.05 (not statistically significant) (Sterne and Kirkwood 2003) would inform the design of future studies. The study set out to recruit 450 participants (estimated to provide 180 participants with LC), to allow for a more flexible margin to counter the possibility of a lower than expected proportion of participants with LC (estimates of 40% of lung-shadow clinic attendees diagnosed with LC was based on estimates gathered from respiratory consultant in SUHT). This proportion (%) might vary over time.

Secondary model:

The secondary objective was to obtain predictive values of symptoms for lung cancer in a sub-group with COPD. Based on the definition of COPD used in the study, we expected approximately 60% of those attending the lung-shadow clinic to have a diagnosis of COPD, following spirometry taken at the clinic. At the projected recruitment rate, the study would have recruited 270 participants

138

Study 2

over the 10 month period (about 25 participants each month) with COPD attending lung-shadow clinic. It was estimated that 40% of participants in the COPD group would have LC (96 cases). A model using 96 LC cases had a power of 80% to detect a difference in frequency from 20% in the non-LC group to 40% in the LC group. Similar to the main model, the identification of symptoms strongly associated with LC (OR>2 or <0.5) but with an alpha of >0.5 (statistically not significant) (Sterne and Kirkwood 2003) would inform the design of future studies.

6.3.9 Recruitment strategy

The recruitment process and consent process for questionnaire administration completion were the same as used in study 1 (see 'Methodology' Chapter Four).

6.3.10 Quantitative data collection

Questionnaire: The data collection process was the same as that in the study 1. Questionnaires received after the pre-established cut-off date (after date of diagnosis) were excluded from the analysis.

Diagnoses and comorbidities: All participants' diagnoses were extracted from the electronic or paper records at the secondary care site (hospital) and identified for the presence or absence of LC and other diagnoses of chest, or respiratory disorders. Diagnoses were recorded six months after recruitment. The stage of the lung cancer and whether or not the lung cancer was operable, at the time of treatment decision was determined to distinguish between early and late stage disease.

The COPD sub-group was defined on the basis of abnormal spirometry post-bronchodilator (predicted FEV1 < 70% and FEV1:FVC ratio <0.70), **and/or** clinical diagnosis by the respiratory clinician, indicated by symptoms of COPD and pre-disposing risk factors (e.g. smoking, age and family history). This category of COPD excluded those diagnosed on the basis of CT evidence of structural damage to the lung to minimise the heterogeneity between patients diagnosed with COPD in primary care and those in lung shadow clinic. In

Study 2

practice, primary care clinicians are unlikely to diagnose COPD on the basis of CT results.

6.3.10.1 Independent symptom variables

Questionnaire items consisted of;

- 1) the generic symptom filter questions, which indicate the presence or absence of that symptom type,
- 2) other variables that refer to the generic symptom filter question; features of the symptom e.g. severity, progression, and chronicity (when symptom first indicated), and
- 3) symptom descriptor variables that are independent symptom variables, which indicate the presence or absence of the symptom indicated by the descriptor, and can be grouped within the generic symptom type.

A summary of the three attributes of the questionnaire items recorded is shown in Table 6.1.

Table 6.1 Types of response categories for questionnaire items

	Type	Description	Examples of questionnaire items corresponding to each type
1	Binary	Questionnaire items with binary response categories that indicate either the presence or absence of the symptom variable.	Q3a_Pain: A niggles, pain or ache that feels like wind or indigestion but not associated with eating Response options: Yes; No

Study 2

2	Categorical	<p>Questionnaire items with multiple response categories indicating the temporality of the symptom; for example symptom frequency, which were collapsed to indicate the presence or absence of the symptom.</p> <p>Questionnaire items with multiple response categories that indicated properties of symptoms; for example: period of onset,</p>	<p>Q13a_Cgh: An irritating cough</p> <p>Original response options: Never; Once; Occasionally; Most of the time</p> <p>Q12_Cgh_R: When did you first had a cough that lasted for more than 3 weeks</p> <p>Original response options: Within the last 3 months; 4-12 months ago; More than 12 months ago</p>
3	Ordinal (10)	Symptom severity variables (not presence or absence of the symptom)	<p>Q9_Pain: On a scale of 0-9, how much the chest pain interfered with everyday life and activities when at its worst</p>

For type 2 and 3 questions, those with multiple categories, tetrachoric correlation matrices were used to identify the optimum cut-off for distinguishing between LC and not LC.

Tetrachoric correlation is a statistical method to estimate the correlation between the two assumed continuous variables underlying the measured dichotomies (Drasgow 1988). This method assumes a dichotomous measure to be a measure of a normally distributed continuous variable. The higher the correlation coefficient, the higher the association between the variables, and the less likely it is that the variables are independent of each other. Tetrachoric correlation provides an alternative to the use of Chi-square to inform the choice of cut-offs that facilitates comparison and evaluation of the magnitude of correlation, when transforming categorical into binary variables.

Study 2

Multiple response variables with more than two responses (never, once, occasionally, most of the time) were collapsed at different response levels to form different dichotomies of the same variable. Tetrachoric correlations were then carried out and the magnitude of the correlations compared across the different possible dichotomies for the same variable. The resultant change in effect sizes were evaluated and discussed.

Similarly, generic symptom variables were collapsed at different response levels to form three binary responses; ever/never, current/non-current, and current or last three months. Results suggested that the cut-off was better for the current or in the last three months. Therefore, subsequent data analyses for this study were carried out at this level (symptoms that are current or in the last three months).

6.3.10.2 Independent socio-demographic variables

LC risk increases with age, and symptoms might also be associated with ageing. Therefore, age was adjusted for, in the multivariable models. Whereas LC risk is also associated with gender, this association would not necessarily be anticipated in an already referred population. Furthermore, a relationship between gender and potential lung cancer symptoms has not been clearly established. Therefore, the relationship between gender and LC, and the relationship between symptoms, gender and lung cancer, were explored in univariate and bivariate analysis, with a view to adjusting for gender in the multivariate analyses, where necessary. Considering the small sample size of this study, potential confounders were not adjusted for where there was no evidence of confounding, as the general rule of thumb is to have at least 10 cases per independent variable.

Whilst logistic regression does not require linearity of the dependent and independent variables, it does assume a linear relationship between the independent variables and the log odds. Therefore, the relationship between age and log odds of cancer diagnosis was investigated, and age was transformed (AGE-squared was also included into the model) to improve linearity.

Information on those not approached was not obtained as there was no permission to collect the information on the date of birth (age) or gender.

6.3.10.3 Independent epidemiological risk variables: clinical and behavioural risk factors

Table 6.2 presents patient-reported risk variables, variables that increase the likelihood of developing LC in the future. These epidemiological risk variables were based on clinical risk factors and a behavioural risk factor (e.g. smoking) identified in previous research.

The Liverpool Lung Project (LLP) risk prediction model targeted at lung cancer screening, identified history of pneumonia within the previous 5 years, exposure to asbestos, previous malignancy, family history of LC, and smoking history, as risk variables (Field et al. 2013; Cassidy et al. 2007). COPD was also included as an independent risk variable in the current study because of its positive association with the development of LC (Hamilton et al. 2005; Punturieri et al. 2009).

Risk items cancer, asbestos related illness, pneumonia and COPD had 5 response categories to indicate when the exposure/diagnosis first occurred: 'Never', 'Within last 3 months', '4-12 months ago', '1-5 years ago', or 'More than 5 years ago'. On the basis of previous research, most of these risk variables (previous cancer, diagnosis of asbestosis, diagnosis of COPD) were dichotomised (yes/no) to capture ever experienced/diagnosed or ever been exposed to these risks. For example, Q69e_COPD was dichotomised into 'yes', ever had a diagnosis of COPD, and 'no', never had a diagnosis of COPD. Pneumonia was coded as a dichotomous variable to indicate pneumonia infection within the last five years.

LC is known to be largely attributable to smoking. Field et al. (2013) have used pack years smoked and number of years smoked as the only smoking parameters to indicate lung cancer risk. They did not record information on the time of smoking cessation to distinguish between ex-smokers and current smokers (Field et al. 2013).

History of smoking variables that enable the calculation of pack years smoked (number of years smoked, and average smoked per day), and smoking status (current, former, or never-smokers) were recorded in the IPCARD questionnaire. However, there were too many missing data in the components

Study 2

of smoking variables (number of years smoked, average of cigarette smoked, and current smoker), in the current study, to calculate the pack years smoked.

The information on current smoking status (and length of time since quit smoking) was obtained to investigate confounding rather than being a risk factor; as a co-existing covariate that might explain symptoms similar to symptoms of lung cancer, and therefore, might be a confounder.

Table 6.2 Risk variables

RISK VARIABLES	
Q69a_Pneumo_last5yrs	Pneumonia in the last five years
Q69e_COPD_ever	Chronic Obstructive Pulmonary Disease (COPD)
Q69h_Cancer_ever	Previous cancer diagnosis
Q69j_Asbес_ever	Asbestos-related illness
Q71d_FamilyHx	Family history of LC
Q73_Smoke	Ever smoker
Q74_Smoke	Age started smoking (Year)
Q75_Smoke	Smoking duration (years)
Q76_Smoke	Smoking status (current/former)
Q77_Smoke	Average amount of cigarettes per day
Q79_Smoke	Average amount of cigarettes used to smoke per day

Note: Q75_Smoke, and Q77_Smoke or Q79_Smoke provide information to calculate the pack-year.

6.3.10.4 Comorbidities

Table 6.3 presents comorbidities (current medical conditions) that may co-exist along with lung cancer. Asthma, allergy, heart disease or angina, arthritis, and respiratory inflammatory-related problems such as recent pneumonia, and COPD, which have some similar symptoms to LC, were analysed for confounding with potential lung cancer symptoms.

Study 2

Symptom presentations of the comorbidities often overlap with symptoms of LC, which might confound the symptom LC relationship. An apparent association between a symptom and LC, or the apparent absence of an association, may be due to the occurrence of a comorbidity that has a relationship with the symptom and lung cancer, rather than being an independent association between the symptom and LC. The investigation of potential confounding involves controlling for the third (confounding) variable by investigating the association between the symptom and lung cancer within strata of the third variable. However, testing for confounding also involves ruling out the potential for effect modification (interaction). Where the effect of a symptom differs between strata of a third variable, interaction rather than confounding would provide an appropriate explanation for observed relationships.

Table 6.3 Baseline comorbidities

COMORBIDITIES	
Q69a_Pneumo_3mths	Pneumonia in the last three months
Q69c_Asthma_ever	Asthma
Q69d_Allergy_ever	Allergy
Q69f_HD_Angina_ever	Heart disease/ angina
Q69e_COPD_ever	COPD
Q69j_Asbestos_ever	Asbestos-related illness
Q69k_Arthritis	Arthritis

6.3.11 Data entry and data cleaning

In the data entry process, data were entered using Remark©, which is an automated questionnaire data entry system to reduce human error. Using this software, data were exported into an Excel spreadsheet.


Data were cleaned in Excel; questionnaire data were checked for consistency and naturally occurring 'non-response' were appropriately coded zero, '0' (see

Study 2

Figure 6.3). Consistency checks would have identified any data that were out of range, logically inconsistent, or had extreme values. Inconsistencies and reasons for missing data of all questionnaire responses were carefully explored (reported in chapter six). The original unedited version was kept so that any changes could be cross-referenced and verified if necessary. Responses were coded appropriately before importing the dataset in Stata for data analysis.

Section 1 - Chest and upper body aches, pain or discomfort

Q1 Have you ever experienced any discomfort in your chest, upper body or shoulders?

No ☒ Please go to Section 2, page 5 

Yes and I still have the pain/discomfort ☐ Please go to question 3

Yes but I no longer have the pain/discomfort ☐ → Q2 Have you had pain/discomfort in the last three months

Yes ☐ No ☐

Please go to question 3 Please go to Section 2, page 5

Q3 Please indicate whether the statements below accurately describe chest or upper body aches, pains or discomfort you have experienced **currently or within the last 3 months** by marking yes or no for each statement.

	Yes	No	
a) A niggle, pain or ache that feels like wind or indigestion but not associated with eating	<input type="radio"/>	<input checked="" type="radio"/>	Non-response
b) Discomfort or pain when laying/sitting in a particular position	<input type="radio"/>	<input checked="" type="radio"/>	
c) Discomfort or pain that feels like bruising	<input type="radio"/>	<input checked="" type="radio"/>	
d) Discomfort or pain that is not brought on by physical activity	<input type="radio"/>	<input checked="" type="radio"/>	

Figure 6.3 Example of questionnaire data with naturally occurring non-response

6.3.12 Quantitative data analysis

The aim of the analysis was to estimate the discriminatory value of symptoms for LC diagnosis in order to develop a model that distinguishes those with LC from those without the disease. All data were entered into Stata version 13.

6.3.12.1 Univariate analysis

Univariate associations between each of the symptom variables and the dependent outcome variable (LC diagnosis) were explored. Statistical significance was assessed using chi-squared test for dichotomous data, and

Study 2

odds ratios for the symptom variables were presented with 95% confidence intervals.

6.3.12.2 Bivariate analysis

Bivariate analyses using MH methods were carried out to explore confounding with comorbidities. Adjusted and crude ORs were compared. A $\geq 10\%$ difference between the adjusted and unadjusted OR can suggest confounding.

Although testing for interactions may contribute to the explanatory effect of the outcome and improve the fit of the model, such tests usually require large sample size. As this was an exploratory feasibility study with low-moderate power, interaction terms were not included in the model to avoid erroneous results.

6.3.12.3 Multivariate analysis

In the main and secondary analyses, multivariate logistic regression was used to model the relationships between:

1. Symptoms, potential confounders and LC status
2. Symptoms, risk factors, and LC status

The aim of the main analysis was to identify symptom variables that distinguish between lung cancer diagnosis and a non-lung cancer diagnosis in a population referred to secondary care (lung-shadow clinic) with high rates of respiratory disease. The secondary analysis, a sub-group analysis of this population, was to distinguish lung cancer diagnosis in a population a population referred to the lung-shadow clinic with COPD.

6.3.12.4 Variable selection procedure

The process of variable selection was performed using the built-in `-stepwise-` command in Stata based on hypothesis tests between nested models using Wald tests, likelihood ratio tests, and *F* statistics. Variables identified in the univariate analyses were modelled using forward stepwise regression with weighting by $1/M$, M = number of imputations (Woods et al. 2008). Two models were built, which consisted of:

Study 2

Model 1: Symptoms; adjusted for age.

Symptom variables that presented with statistically significant univariate associations ($p < 0.05$) were entered into the model.

The relationship with gender and age was explored in the univariate analyses. There was no imbalance in gender to suggest potential confounding; the distribution of male LC cases (61%) was comparable to that of the non-LC cases (59%). Similarly in females, the distribution was almost the same; 39% in the LC group, and 41% in the non-LC group. Adding gender to the model did not make any difference to the main effects model with similar variables remaining in the model (log odds ratio or β coefficient = -0.08). This suggests that gender hardly accounted for the variability of the model predicting the outcome, and was therefore, removed.

Age was entered into the logistic regression models to adjust for the difference between the LC and non-LC group (potential confounding). Ageing is a known risk factor in the development of cancer (CRUK 2009). Incidence of cancer increases with age possibly due to the cumulative exposure of different risks over time. As the variable age is rarely linear and often needs to be transformed. Previous IPCARD studies also accounted for age-squared, which was found to account for some of the variability in the outcome and the better the fit of the model. Therefore, both age variables were added to the models (age was treated as a continuous variable).

Main effects model (model1) was fitted using stepwise logistic regression. The criteria for entry into model 1 was set at significance level of $p = 0.05$ and the criteria for removal from the model was $p = 0.1$.

If no confounders were suggested from the bivariate analyses, they were not adjusted for in the model.

Model 2: Symptoms; adjusted for age, and risk variables.

Risk variables and symptoms (identified in the bivariate analyses) which might improve the predictive accuracy of the model were added to model 1. The entry and exit criteria for model 2 were set at $p = 0.05$ and $p = 0.10$, respectively.

All discarded variables were checked against the final models.

Study 2

Multivariate analysis with relaxed criteria for variable selection ($p < 0.15$): Model 1 and 2

A separate model for 1 and 2 with variable selection at a lower threshold for significance level ($p < 0.15$) in the univariate associations was also explored. This included variables that showed an unadjusted OR > 2 , or < 0.5 . Significance levels (p-values) should be interpreted within the context of each individual study design and parameters to avoid omitting clinically important variables, especially if the study is lacking power (Sterne and Kirkwood 2003). However, over fitting and under fitting of the model should be avoided. The model should be fitted correctly, where only meaningful variables are added with no extraneous variables. An approach to ensure this, was to use the forward stepwise method to estimate the logistic regression. The entry and exit criteria in the stepwise regression was also relaxed, set at $p = 0.10$ and $p = 0.15$, respectively. The Akaike's Information Criterion (AIC) was used to inform the fit of the model and parsimony.

All discarded variables were checked against the final models.

Multicollinearity

Multicollinearity or collinearity occurs when two or more independent variables are linearly related to other independent variables in the model. Stata will drop variables that are perfectly collinear to the other variables. The omitted variable should be justifiable in theory, rather than assumed to be the 'correct' variable removed (Menard 2010; Berry and Feldman 1985). Unusually large odds ratio and standard error in the model will be investigated for multicollinearity. As mentioned, goodness-of-fit tests were carried out after each logit procedure to identify observations that might have significant impact on model fit.

6.3.12.5 Developing a set of diagnostic criteria for LC population

Based on the variables identified from the main effects model, a potential set of diagnostic criteria was suggested using a simple point-scoring system that might distinguish between LC and non-LC in a population with high rates of acute or chronic respiratory disease more generally.

Study 2

The sensitivity, specificity, predictive values and likelihood odds ratios was calculated for all levels of cut-offs to determine the optimal threshold at which the set of criteria best distinguish between those with lung cancer and those without. The levels of cut-offs are usually decided on an arbitrary point on a continuum of the trade-off between sensitivity and specificity that was optimal for the diagnosis of a particular disease (Altman and Bland 1994b). This optimal cut-off between sensitivity and specificity would indicate the diagnostic accuracy of the threshold for referral. It will show the number of cases with cancer correctly referred over those who were unnecessarily referred for LC investigation. These diagnostic indicators can be represented in a 2x2 contingency table as shown below:

		Diagnosis	
		Lung cancer	No Lung cancer
Test Outcome	Symptom	True Positive (TP)	False Positive (FP)
	No Symptom	False Negative (FN)	True Negative (TN)

The formulas to calculate the diagnostic indicators are as follow (definitions can be found in the abbreviations page (refer to Page xxv):

- Sensitivity (% of true positive or true positive rate) = $TP/(TP+FN)$
- Specificity (% of true negative or true negative rate) = $TN/(TN+FP)$
- Positive Likelihood ratio (LR) = $Sensitivity / (1 - specificity)$
- Negative Likelihood ratio = $1 - (Positive Likelihood ratio)$ or $(1 - sensitivity)/specificity$
- Odds Ratio = $(TP/FN)/(FP/TN)$ or $\frac{PPV/(1-PPV)}{(1-NPV)/NPV}$ or $\frac{Positive\ LR}{Negative\ LR}$

6.3.12.6 Weighted diagnostic criteria

A potential set of weighted diagnostic criteria was also developed using a weighted scoring system. In the same way, coefficients for symptoms that were significantly associated with LC in the model were used to derive a weight for each criterion. Log odds ratio (β coefficient) was calculated for each variable used. Therefore, each item or variable was assigned a weight to reflect the relative value as a predictor.

6.3.12.7 Secondary sub-group analysis

All of the above analytical processes were also applied to the sub-group analysis. This secondary analysis was modelled within a COPD sub-group to identify symptoms that distinguish between lung cancer and non-lung cancer in those with COPD.

6.4 Results Section

Missing data

6.4.1 Frequencies of missing data

The percentages of missing data in the generic symptom variables range from 1.1% to 6.7%. A summary of the frequency and percentage of missing data for each variable can be found in Appendix 11.

Data was classified as missing data when no value was observed for a variable where there should have been a response. This does not include part of a skip question where the variable does not apply to individual respondents, as those variables would have been recoded '0' because the individual has never experienced the generic symptom variable, which covers a set of items (symptom descriptors) that all fall under that symptom category. Figure 6.4 presents a generic symptom variable with reference to the IPCARD questionnaire.

Throughout the chapter, the term, 'generic symptom variable' will be used to refer to the generic symptoms (groupings of symptoms) that are represented by the filter question at the head of each section of the questionnaire (refer to Column 1 of Table 6.4).

Study 2

Section 1 - Chest and upper body aches, pain or discomfort

Q1 Have you ever experienced any discomfort in your chest, upper body or shoulders?

No ☐ Please go to Section 2, page 5

Yes and I still have the pain/discomfort ☐ Please go to question 3

Yes but I no longer have the pain/discomfort ☐ → Q2 Have you had pain/discomfort in the last three months

Yes ☐ No ☐

Please go to question 3 Please go to Section 2, page 5

Q3 Please indicate whether the statements below accurately describe chest or upper body aches, pains or discomfort you have experienced **currently or within the last 3 months** by marking yes or no for each statement.

	Yes	No
a) A niggle, pain or ache that feels like wind or indigestion but not associated with eating	<input type="radio"/>	<input type="radio"/>
b) Discomfort or pain when laying/sitting in a particular position	<input type="radio"/>	<input type="radio"/>
c) Discomfort or pain that feels like bruising	<input type="radio"/>	<input type="radio"/>
d) Discomfort or pain that is not brought on by physical activity	<input type="radio"/>	<input type="radio"/>

Generic symptom variable

Figure 6.4 Example of generic symptom variable in questionnaire

In addition to the generic symptom variables, other symptom variables (symptom descriptors) with percentage of missing data 10%, and above, were also explored (see column 1 of Table 6.5). None of the variables with high percentage of missing values were identified as problematic to the convergence process.

The frequencies of missing data in the remaining covariates were between 0.8% and 9.1%, with an average percentage of missingness of 6.0%. The distribution of missing data in the risk variables were as follows; pneumonia in the last 5 years (6.1%), COPD (8.1%), previous cancer (7.8%), asbestos-related illnesses (6.7%), ever smoker (3.3%), and family history of lung cancer (11.1%). Although current smoking status (Q76_Smoke) is a potentially important covariate, and potential confounder, the variable had too much missing data (24% missingness) to be included in the multivariate analysis. Therefore, ever smoked variable (Q73_Smoke) was modelled instead.

Diagnoses (lung cancer and COPD) were based on patients' medical notes and radiological results. The diagnoses variables accurately represented, and were as complete as, what was recorded in the notes. COPD diagnoses were

Study 2

categorised on the basis of abnormal spirometry results and/or clinical diagnosis. However, a small number of patients (n=41) did not have spirometric tests performed in the clinic, and it is possible that in those group they could have undetected COPD, albeit it would be a very small percentage. Most of these cases clinicians deemed unnecessary for spirometry. Two of the 41, had previous history of tuberculosis, and therefore, were unable to perform the spirometry in general outpatients as part of infection control policy. For most participants except for those who passed away before a diagnosis was confirmed, their clinical diagnoses were followed up and extracted.

The following section will explore the relationship between missing data and covariates and in doing so; will be checking departures from MCAR assumptions. The identification of relationships with missingness is consistent with MAR assumptions.

6.4.2 Distribution of missing data in symptom variables by socio-demographic variables

6.4.2.1 Gender

Table 6.4 and Table 6.5 present the frequency distributions for the generic symptom variables, and variables with percentage of missing data >10% by gender.

There were more men (n=213) than women (n=146) in the study. Individual chi-square tests were used to assess associations between missing data and gender. There were no statistically significant relationships with gender for the variables examined, except for variable, Q29_Tired_ever: "*Have you experienced any unexpected tiredness within the last 12 months?*" (p=0.0369).

Study 2

Table 6.4 Investigation of non-differential classification of the proportion of missingness in each generic variable by gender

Miss_Variables	Missingness (%)	Gender (%)		Chi ² p-value
		Male n=213	Female n=146	
Q1_Pain_3mths	3 (0.84)	2 (0.9)	1 (0.7)	0.795
Q10_Cgh_3mths	12 (3.34)	9 (4.2)	3 (2.1)	0.261
Q19_BrChnges_3mths	19 (5.29)	10 (4.7)	9 (6.2)	0.541
Q29_Tired_3mths	12 (3.34)	5 (2.3)	7 (4.8)	0.205
Q38_CghBlood_3mths	10 (2.79)	4 (1.9)	6 (4.1)	0.207
Q53_HCswat_3mths	14 (3.90)	6 (2.8)	8 (5.5)	0.201
Q63_New_JPain_12mths	20 (5.57)	13 (6.1)	7 (4.8)	0.595
Q64_New_JPain_12mths	22 (6.13)	15 (7.0)	7 (4.8)	0.383
Q43_ChInfectn	17 (4.74)	9 (4.2)	8 (5.5)	0.583
Q44_ChInfectn	24 (6.69)	14 (6.6)	10 (6.8)	0.918
Q51_Weight	17 (4.74)	12 (5.6)	5 (3.4)	0.333
Q60_EatChnges	20 (5.57)	9 (4.2)	11 (7.5)	0.179
Q65_Voice	15 (4.18)	9 (4.2)	6 (4.1)	0.957
Q67_Skin	21 (5.85)	9 (4.2)	12 (8.2)	0.113

Each symptom coded 1=ever, 0=never

Study 2

Table 6.5 Investigation of non-differential classification of the proportion of missingness in variables with > 10% missing data by gender

Miss_Variables	Missingness (%)	Gender (%)		Chi ² p-value
		Male n=213	Female n=146	
Q3j_Pain	44 (12.26)	27 (12.7)	17 (11.6)	0.770
Q13c_Cgh	36 (10.00)	21 (9.9)	15 (10.3)	0.898
Q13f_Cgh_	48 (13.37)	26 (12.2)	22 (15.1)	0.434
Q13k_Cgh	42 (11.70)	26 (12.2)	16 (11.0)	0.718
Q14a_Cgh	37 (10.30)	23 (10.8)	14 (9.6)	0.711
Q14b_Cgh_	42 (11.70)	25 (11.7)	17 (11.6)	0.978
Q14c_Cgh_	41 (11.42)	26 (12.2)	15 (10.3)	0.572
Q48_ChInfectn	36 (10.03)	18 (8.5)	18 (8.5)	0.230
Q62_EatChnges	48 (13.37)	27 (12.7)	21 (14.4)	0.641

6.4.2.2 Age

Multiple logistic regression was used to test the association between the missingness of symptom variables (binary independent variable), and age (continuous variable); see Table 6.6 and Table 6.7, below. To ensure valid result, the data were checked against the assumptions of logistic regression. The dependent variable need not be normally distributed, nor linearly related to the independent variable.

Missingness in generic symptom variable for breathing changes/difficulties (Q19_BrChnges_ever) was found to be statistically associated to age (cut off at 68). According to Table 6.6, for every year increase in age, the likelihood that the data in Q19_BrChnges_ever is missing increases by 0.053. There were no observable relationships between missingness in the remaining generic symptom variables and age.

Study 2

Table 6.6 Linear regression of generic symptom variables as explanatory/predictor variables and age the dependent variable (continuous)

AGE	Coefficients	Standard error	p-value	95% Confidence Interval	
Q1_Pain_ever	-0.067	0.044	0.128	-0.154	0.019
Q10_Cgh_ever	0.012	0.024	0.629	-0.036	0.059
Q19_BrChnges_ever	0.053	0.021	0.011*	0.012	0.094
Q29_Tired_ever	0.021	0.021	0.299	-0.019	0.062
Q38_CghBlood_ever	0.016	0.026	0.557	-0.036	0.067
Q53_HCsweat_ever	0.033	0.023	0.152	-0.012	0.078
Q63_NewJPain_12mths	0.031	0.019	0.109	-0.007	0.069
Q64_NewJPain_12mths	0.033	0.019	0.074	-0.003	0.070
Q43_ChInfectn	0.015	0.020	0.473	-0.025	0.055
Q44_ChInfectn	0.019	0.018	0.269	-0.015	0.054
Q51_Weight	0.035	0.021	0.098	-0.006	0.076
Q60_EatChnges	0.021	0.019	0.262	-0.016	0.059
Q65_Voice	0.036	0.022	0.107	-0.008	0.080
Q67_Skin	0.011	0.018	0.535	-0.025	0.048

*statistically significant

Each symptom coded 1=ever missing, 0=not missing

Analysis suggested no observable relationship between missingness in the symptom variables with 10% or more missing data and age (see Table 6.7).

Study 2

Table 6.7 Linear regression of variables with 10% or more missing data as explanatory/ predictor variables and age the dependent variable (continuous)

AGE	Standard		p-value	95% Confidence	
	Coefficients	error		Interval	
Q3j_Pain	-0.018	0.013	0.160	-0.044	0.007
Q13c_Cgh	-0.013	0.014	0.366	-0.041	0.015
Q13f_Cgh_	-0.004	0.013	0.742	-0.029	0.021
Q13k_Cgh	0.020	0.014	0.147	-0.007	0.046
Q14a_Cgh	-0.000	0.014	0.972	-0.028	0.027
Q14b_Cgh	-0.003	0.013	0.828	-0.029	0.023
Q14c_Cgh	-0.003	0.014	0.828	-0.029	0.024
Q48_ChInfectn	0.021	0.015	0.143	-0.007	0.050
Q62_EatChnges	0.009	0.013	0.496	-0.016	0.034

*statistically significant

Socio-economic status

Postcode data was not available in this feasibility study dataset. Therefore it was not possible to derive socio-economic status from postcode-linked deprivation scores to be able to check the distribution of missing data by socio-demographic data.

6.4.3 Distribution of missing data in symptom variables by clinical outcome variables and clinical covariates

6.4.3.1 Cancer (Outcome variable)

Chi-square tests were used to explore the relationship between missingness in generic symptom variables and symptom variables with >10% missing data, and lung cancer diagnosis (LC/not-LC).

All of the variables in Table 6.8 had p-values above 0.05 suggesting no statistical differences in missingness between those with and without lung cancer, except for one of the generic symptom variables, for breathing changes in the last three months (Q19_BrChnges_3mths), $p=0.001$.

Study 2

Table 6.8 Investigation of non-differential classification of the proportion of missingness in each generic variable by lung cancer diagnosis

Miss_Variables	Missingness (%)	Lung cancer diagnosis (%)		Chi ² p-value
		LC n=77	No LC n=282	
Q1_Pain_3mths	3 (0.84)	0 (0)	3 (1.1)	0.363
Q10_Cgh_3mths	12 (3.34)	3 (3.9)	9 (3.2)	0.760
Q19_BrChnges_3mths	19 (5.29)	10 (13.0)	9 (3.2)	0.001*
Q29_Tired_3mths	12 (3.34)	2 (2.6)	10 (3.5)	0.681
Q38_CghBlood_3mths	10 (2.79)	1 (1.3)	9 (3.2)	0.371
Q53_HCsweat_3mths	14 (3.90)	3 (3.9)	11 (3.9)	0.999
Q63_New_JPain_12mths	20 (5.57)	3 (3.9)	17 (6.0)	0.470
Q64_New_JPain_12mths	22 (6.13)	5 (6.5)	17 (6.0)	0.880
Q43_ChInfectn	17 (4.74)	5 (6.5)	12 (4.3)	0.412
Q44_ChInfectn	24 (6.69)	7 (9.1)	17 (6.0)	0.340
Q51_Weight	17 (4.74)	5 (6.5)	12 (4.3)	0.412
Q60_EatChnges	20 (5.57)	6 (7.8)	14 (5.0)	0.338
Q65_Voice	15 (4.18)	5 (6.5)	10 (3.5)	0.252
Q67_Skin	21 (5.85)	5 (6.5)	16 (5.7)	0.786

Each symptom coded 1=ever, 0=never

No statistically significant differences were observed for variables with >10% missingness between the LC and non-LC group ($p>0.05$), see Table 6.9.

Table 6.9 Investigation of non-differential classification of the proportion of missingness in variables with missing data >10% by lung cancer diagnosis

Miss_Variables	Missingness (%)	Lung cancer diagnosis (%)		Chi ² p-value
		LC n=77	No LC n=282	
Q3j_Pain	44 (12.3)	10 (13.0)	34 (12.1)	0.825
Q13c_Cgh	36 (10.0)	7 (9.1)	29 (10.3)	0.757
Q13f_Cgh_	48 (13.4)	10 (13.0)	38 (13.5)	0.911
Q13k_Cgh	42 (11.7)	7 (9.1)	35 (12.4)	0.422
Q14a_Cgh	37 (10.3)	8 (10.4)	29 (10.3)	0.978
Q14b_Cgh_	42 (11.7)	10 (13.0)	32 (11.3)	0.692
Q14c_Cgh_	41 (11.4)	9 (11.7)	32 (11.3)	0.934
Q48_ChInfectn	36 (10.0)	6 (7.8)	30 (10.6)	0.461
Q62_EatChnges	48 (13.4)	9 (11.7)	39 (13.8)	0.625

6.4.3.2 COPD

Chi-square tests were used to explore missingness across the according to COPD diagnosis (COPD/no COPD). No statistically significant differences between COPD and non-COPD group were observed for the generic symptom variables and variables with >10% missingness ($p>0.05$); see Table 6.10 and Table 6.11.

Study 2

Table 6.10 Investigation of non-differential classification of the proportion of missingness in each generic variable by COPD diagnosis

Miss_Variables	Missingness (%)	COPD diagnosis (%)		Chi ² p-value
		COPD n=124	No COPD n=235	
Q1_Pain_3mths	3 (0.84)	1 (0.8)	2 (0.9)	0.965
Q10_Cgh_3mths	12 (3.34)	2 (1.6)	10 (4.3)	0.185
Q19_BrChnges_3mths	19 (5.29)	7 (5.6)	12 (5.1)	0.828
Q29_Tired_3mths	12 (3.34)	3 (2.4)	9 (3.8)	0.480
Q38_CghBlood_3mths	10 (2.79)	4 (3.2)	6 (2.6)	0.713
Q53_HCswat_3mths	14 (3.90)	4 (3.2)	10 (4.3)	0.632
Q63_New_JPain_12mths	20 (5.57)	7 (5.6)	13 (5.5)	0.965
Q64_New_JPain_12mths	22 (6.13)	9 (7.3)	13 (5.5)	0.517
Q43_ChInfectn	17 (4.74)	6 (4.8)	11 (4.7)	0.947
Q44_ChInfectn	24 (6.69)	7 (5.6)	17 (7.2)	0.567
Q51_Weight	17 (4.74)	7 (5.6)	10 (4.3)	0.555
Q60_EatChnges	20 (5.57)	5 (4.0)	15 (6.4)	0.356
Q65_Voice	15 (4.18)	7 (5.6)	8 (3.4)	0.313
Q67_Skin	21 (5.85)	8 (6.5)	13 (5.5)	0.724

Each symptom coded 1=ever, 0=never

Study 2

Table 6.11 Investigation of non-differential classification of the proportion of missingness in variables with missing data >10% by COPD diagnosis

Miss_Variables	Missingness (%)	COPD diagnosis (%)		Chi ² p-value
		COPD n=124	No COPD n=235	
Q3j_Pain	44 (12.3)	16 (12.9)	28 (11.9)	0.786
Q13c_Cgh	36 (10.0)	15 (12.1)	21 (8.9)	0.343
Q13f_Cgh_	48 (13.4)	19 (15.3)	48 (20.4)	0.430
Q13k_Cgh	42 (11.7)	15 (12.1)	27 (11.5)	0.865
Q14a_Cgh	37 (10.3)	17 (13.7)	20 (8.5)	0.123
Q14b_Cgh_	42 (11.7)	16 (12.9)	26 (11.1)	0.606
Q14c_Cgh_	41 (11.4)	15 (12.1)	26 (11.1)	0.770
Q48_ChInfectn	36 (10.0)	13 (10.5)	23 (9.8)	0.834
Q62_EatChnges	48 (13.4)	18 (14.5)	30 (12.8)	0.643

6.4.3.3 Epidemiological Risk variable

Chi-square tests showed no statistically significant associations between the risk variable, previous history of smoking, and missingness in symptoms variables (generic symptom variables and variables with >10% missingness).

6.4.4 Distribution of missing data in symptom variables by clinically relevant symptom covariates

Logistic regression was used to explore if there is an association between missingness in the symptom variables (generic symptom variables and variables with >10% missingness) and symptoms selected for the imputation model. The variables used in the imputation model were chosen strictly based on the clinical understanding of the associations between symptoms. Therefore, some relationships between these variables were expected. The conditional models for the imputation process are included in Appendix 12.

Results of the analyses have been tabulated (see Table 6.12 and Table 6.13). Five of the covariates (chosen for imputation) predicted missingness of four missing dummy variables at a significant level, which indicates an association

Study 2

exist between the variable with the missing data and the variables intended for imputation. The variables that explain the missingness makes the assumption of MAR more plausible (Marchenko and Eddings 2011).

Covariate, Q43_ChInfectn, predicted missingness in Q13f_Cgh (*Cough that feels as though it arises in one or other lung or side of the chest*), and Q13k_Cgh (*A hard or harsh cough without phlegm*) at a significant level ($p < 0.05$; $p = 0.006$ and $p = 0.014$, respectively). Participants who reportedly currently having a phlegmy chest or chest infection were found to be almost 4 times ($OR = 4$) more likely to have missing values in Q13f_Cgh and Q13k_Cgh.

Similarly, covariates, Q10_Cgh and Q69j_Asbes_ever, were found to predict missingness in Q48_ChInfectn; *whether one had noticeably more colds or flu within the last 12 months than the year before* (with p -values = 0.015 and 0.003, respectively). However, the confidence intervals of the association between Q69j_Asbestos_ever and missingness in Q48_ChInfectn were relatively wider (2.1 to 35.3), and should be interpreted with caution.

Covariate, Q49_Weight (*Having to eat more in order to maintain a steady weight*) was also associated to missingness in Q62_EatChnges at a statistically significant level (p -value = 0.048) with a high OR of five. The confidence interval for this variable, Q62_EatChnges, was also comparatively wider (1.01 to 23.65), which suggests less precision, and more uncertainty about the unknown parameter (see Table 6.12).

Study 2

Table 6.12 Covariates that predicted the missingness in the variables with higher missing data > 10% at a statistically significant level ($p < 0.05$)

Dependent variable	Independent Covariates	Odds Ratio (OR)	p-value	95% Confidence Intervals	
miss_Q13f_Cgh	Q43_ChInfectn	4.22	0.006*	1.52	11.71
miss_Q13k_Cgh	Q43_ChInfectn	4.04	0.014*	1.32	12.36
miss_Q48_ChInfectn	Q10_Cgh_3mths	0.30	0.015*	0.11	0.79
	Q69j_Asbess_ever	8.59	0.003*	2.08	35.50
miss_Q62_EatChnges	Q49_Weight	4.90	0.048*	1.01	23.65

* Statistical significance

Statistically significant associations were found between symptom variables in the imputation model, and three out of the nine generic symptom variables, as shown in Table 6.13.

Dependent variables, miss_Q51_Weight (OR = 6.9; p-value = 0.049), miss_Q60_EatChnges (OR = 76.6; p-value = 0.045), and miss_Q64_New_Jpain_12mths (OR=3.09; p-value=0.05) had statistically significant associations with independent covariates, Q44_ChInfectn, Q49_Weight, and Q69k_Arthritis, respectively (see Table 6.13). However, caution is needed when interpreting the effect sizes of some of the associations; Q51_Weight and Q60_EatChnges. The wide confidence interval (1 to 5285) suggests large variability and uncertainties about the unknown parameters.

Study 2

Table 6.13 Covariates that predicted the missingness in generic symptom variables at a statistically significant level ($p < 0.05$)

Dependent variable	Independent Covariates	Odds Ratio (OR)	p-value	95% Confidence Intervals	
miss_Q51_Weight	Q44_ChInfectn	6.94	0.049*	1.00	47.88
miss_Q60_EatChnges	Q49_Weight	76.63	0.045*	1.11	5286
miss_Q64_New_Jpain	Q69k_Arthritis	3.09	0.050*	1.00	5.82

* Statistical significance

The individual relationships between miss_variables and independent covariates that were not statistically significant but ORs > 2.0 or < 0.05 were also explored in Table 6.14 and Table 6.15. The lack of statistical significance does not definitely exclude all associations between the dependent variables and the independent covariates.

Table 6.14 Covariates with ORs > 2 but not statistically significant ($p > 0.05$) (variables with higher missing data $> 10\%$)

Dependent variable	Independent Covariates	Odds Ratio (OR)	p-value	95% Confidence Intervals	
miss_3j_Pain	Q69j_Asbess_ever	2.46	0.292	0.46	13.17
miss_13c_Cgh	Q19_BrChnges_ever	2.43	0.182	0.66	8.92
miss_Q13f_Cgh	Q19_BrChnges_ever	3.58	0.101	0.78	16.42
miss_Q13k_Cgh	Q69a_Pneumo_last5yrs	3.42	0.070	0.90	12.94
miss_14a_Cgh	Q19_BrChnges_ever	2.01	0.303	0.53	7.58
miss_Q14b_Cgh	Q69d_Allergy_ever	2.71	0.078	0.89	8.20
	Q69j_Asbess_ever	2.73	0.250	0.49	15.16
miss_Q14c_Cgh	Q19_BrChnges_ever	3.21	0.140	0.68	15.13
	Q43_ChInfectn	2.21	0.158	0.73	6.66
	Q42_CghBlood	2.03	0.231	0.64	6.50

Study 2

Table 6.15 Covariates with ORs > 2 but not statistically significant ($p > 0.05$)
(generic symptom variables)

Dependent variable	Independent Covariates	Odds Ratio (OR)	p-value	95% Confidence Intervals	
miss_Q1_Pain_3mths	Q69k_Arthritis	6.43	0.115	0.63	65.06
miss_Q19_BrChnges_3mths	Q69c_Asthma_ever	2.61	0.351	0.35	19.52
	Q69e_COPD_ever	2.76	0.294	0.41	18.42
miss_Q38_CghBlood_3mths	Q69e_COPD_ever	4.91	0.265	0.30	80.50
miss_Q43_ChInfectn	Q10_Cgh_ever	3.13	0.321	0.33	29.9
miss_Q44_ChInfectn	Q36_Tired	2.95	0.231	0.50	17.39
	Q69e_COPD_ever	2.03	0.477	0.29	14.41
miss_Q53_Hcsweat_3mths	Q43_ChInfectn	3.57	0.371	0.22	58.14
miss_Q63_New_JPain_12mths	Q69k_Arthritis_ever	2.32	0.278	0.51	10.65
miss_Q65_Voice	Q43_ChInfectn	8.01	0.095	0.70	92.07

The relationships between Q43_ChInfectn and miss_Q13f_Cgh; Q43_ChInfectn and miss_Q13k_Cgh; Q10_Cgh_3mths and miss_Q48_ChInfectn; Q69j_Asbes_ever and miss_Q48_ChInfectn; Q49_Weight and miss_Q62_EatChnges; Q44_ChInfectn and miss_Q51_Weight; Q49_Weight and miss_Q60_EatChnges; Q69k_Arthritis and miss_Q64_New_JPain_12mths appear to be consistent with MAR (see Table 6.14 and Table 6.15). An association indicates a visible departure from MCAR in the proposed analysis as the missing data mechanism is less likely to be working within the assumptions of ‘completely at random’, which suggests that for these variables at least, MAR assumption is plausible. There is no definitive way of assessing whether the data are MCAR, MAR or MNAR (Allison 2002). We can only state any observable departures from MCAR assumptions in the missing data (Carpenter and Kenward 2008).

Study 2

Diagnostics performed in the Section 6.4.28.5 will help to explain how this departure might affect the complete case analyses, and describe the differences between complete case and multiple imputation analyses.

6.4.5 Issues of collinearity

Two generic symptom variables, discomfort in chest, upper body, and shoulders in the last three months (miss_Q1_Pain_3mths), and coughing that lasted for more than three weeks in the last three months (miss_Q10_Cgh_3mths), had large number of variables omitted from the bivariate logistic regression analysis due to collinearity. Collinearity suggests a linear relationship between two explanatory variables (i.e. covariates) (Belsley et al. 1980). However, collinearity diagnostics performed disproved any adverse collinearity problem with a variance inflation factor (VIF) of one and a tolerance value of almost one, both suggesting that no two 'X' variables were correlated. Multicollinearity is indicated if VIF is greater than 10 and tolerance is less than 0.10 (Belsley et al. 1980; O'Brien 2007). There are many remedies to multicollinearity issues in regression, which include dropping one or two of the explanatory variables in order to produce a model with significant coefficients. However, this will result in loss of information. The current study preserved the model as it was. The occurrence of collinearity in the regression model is likely to be due to the small proportion of missing cases (1= missing) in the dummy variables for the generic symptom variables, rather than associations between the explanatory variables, and should not affect the regression model (Allison 2012). Most of the generic symptom variables had reasonably low percentage of missing observations, as little as 1% (1% to 6%). Hence, it is not unreasonable for the model to predict failure ('0') perfectly resulting in the omission of independent variables. Results of the collinearity diagnostics are included in Appendix 13.

6.4.6 Result of convergence

There were issues of non-convergence during the imputation process, which were possibly due to the large volume of categorical variables with multiple responses. In the attempt to create a working imputation model, some of the

Study 2

variables with more than two responses were dichotomised, where appropriate. The response cut-offs for these multiple response variables were carefully explored using tetrachoric correlations in Section 6.4.12. A dry run of the model was performed ‘noisily’ in Stata to diagnose potentially problematic variables. Blocks of variables that were repeatedly omitted due to collinearity in the dry run were removed, and added back in small groups until the model stopped working. The group of risk variables was found to be a problem to convergence. However, it was difficult to justify dropping the risk variables because they were potentially important to the imputation model based on clinical knowledge.

6.4.7 Pattern of missing data

The Stata output for the ‘misstable pattern’ command showed that there was no obvious monotone pattern to the missingness in the dataset, therefore, the missing data pattern was arbitrary, which cannot be modelled with MVN. MICE suitably handles both monotone and arbitrary missing data patterns. Its variable-by-variable approach to imputation, is appropriate for incomplete datasets with large number of categorical variables. Variables in the MICE approach also do not assume multivariate normality which meant that a variety of variables (binary, ordinal, and categorical variables) could be imputed (White et al. 2011). Although MICE lacks the theoretical underpinnings that MVN has (Raghunathan et al. 2001; Schafer 1999), it provides means of capturing important data characteristics such as ranges and restrictions within a subset of the data, which MVN cannot. The flexibility of MICE was needed for the practical success of imputing in this study. Stata output of the missing data pattern is found in Appendix 14.

6.4.8 Diagnostics for imputations

Table 6.16 compares the descriptive statistics of the severity variables, scaled from 0 to 9, of the complete-case and imputed data. The means and ranges (95% confidence intervals) looked equal for both data. Negligible differences were also observed in the frequency distribution of the severity variables cut off at 6 (see Table 6.17).

Study 2

Table 6.16 Means, medians, and standard deviations of the complete case and imputed for severity variables

Variable	Complete case				Imputed		
	Mean	Range			Mean	Range	
Q8_Pain	3.22	2.87	3.57		3.24	2.89	3.59
Q9_Pain	2.60	2.25	2.95		2.62	2.26	2.98
Q18_Cgh	3.68	3.34	4.03		3.66	3.30	4.01
Q28_BrChnges	3.71	3.36	4.05		3.68	3.34	4.02
Q37_Tired	2.36	2.04	2.69		2.36	2.03	2.69

Table 6.17 Frequency distribution of ordinal variables (cut-off at 6)

Variable	Complete case			Imputed	
	0-5 (%)	6-9 (%)		0-5	6-9
Q8_Pain	240 (71.9)	94 (28.1)		257 (71.7)	102 (28.3)
Q9_Pain	254 (76.7)	77 (23.3)		275 (76.6)	84 (23.4)
Q18_Cgh	225 (67.6)	108 (32.4)		243 (67.7)	116 (32.3)
Q28_BrChnges	217 (65.2)	116 (34.8)		235 (65.4)	124 (34.6)
Q37_Tired	265 (78.6)	72 (21.4)		283 (78.7)	76 (21.3)

Table 6.18, Table 6.19, and Table 6.20 present the frequencies of the different types of categorical variables (binary and multiple-response variables) for the imputed and complete dataset. The frequency distributions of the complete case and imputed variables were comparable. The frequencies of symptoms reported for the binary variables were equal or slightly higher in the imputed dataset.

Study 2

Table 6.18 Frequency distribution for binary variables (only 'yes' or '1' response presented)

Variable	Complete case		Imputed
	Yes/1 (%)		Yes/1 (%)
Q1_Pain_3mths	227 (63.8)		229 (63.8)
Q3a_Pain	97 (28.5)		102 (28.5)
Q3b_Pain	113 (33.1)		118 (32.9)
Q3c_Pain	83 (24.7)		90 (25.0)
Q3d_Pain	116 (34.3)		123 (34.2)
Q3e_Pain	162 (47.5)		168 (46.9)
Q3f_Pain	110 (32.6)		117 (32.5)
Q3g_Pain	112 (32.8)		117 (32.6)
Q3h_Pain	124 (36.2)		130 (36.3)
Q3i_Pain	69 (20.7)		74 (20.6)
Q3j_Pain	40 (12.7)		47 (12.9)
Q6_Pain	88 (26.4)		94 (26.2)
Q7_Pain	29 (8.6)		31 (8.7)
Q10_Cgh_3mths	239 (69.9)		245 (68.2)
Q12_Cgh_R	99 (28.6)		104 (29.0)
Q13a_Cgh	199 (59.2)		212 (59.1)
Q13b_Cgh	161 (48.6)		174 (48.5)
Q13c_Cgh	160 (49.5)		178 (49.6)
Q13d_Cgh	185 (56.4)		202 (56.2)
Q13e_Cgh	153 (46.8)		169 (47.1)
Q13f_Cgh	106 (34.1)		126 (35.0)
Q13g_Cgh	152 (46.5)		166 (46.3)
Q13h_Cgh	159 (48.5)		172 (47.9)
Q13i_Cgh	176 (52.5)		186 (51.9)
Q13j_Cgh	182 (54.3)		194 (54.0)
Q13k_Cgh	110 (34.7)		123 (34.4)

Study 2

Q13I_Cgh	141 (43.3)		153 (42.6)
Q15_Cgh	86 (25.6)		92 (25.5)
Q16_Cgh	58 (17.3)		63 (17.5)
Q19_BrChnges_3mths	244 (71.8)		258 (71.9)
Q21_BrChnges	91 (26.5)		95 (26.4)
Q22a_BrChnges	190 (55.9)		200 (55.8)
Q22b_BrChnges	139 (41.9)		150 (41.9)
Q22c_BrChnges	94 (28.1)		101 (28.2)
Q22d_BrChnges	73 (22.2)		81 (22.5)
Q22e_BrChnges	98 (29.5)		106 (29.5)
Q22f_BrChnges	77 (23.4)		85 (23.7)
Q23_BrChnges	41 (12.6)		47 (13.0)
Q24a_BrChnges	224 (65.9)		235 (65.3)
Q24b_BrChnges	140 (42.0)		152 (42.2)
Q24c_BrChnges	126 (37.3)		133 (36.9)
Q24d_BrChnges	145 (44.2)		160 (44.6)
Q24e_BrChnges	96 (29.1)		105 (29.4)
Q24f_BrChnges	90 (27.0)		98 (27.3)
Q26_BrChnges	126 (37.2)		133 (37.0)
Q27_BrChnges	54 (16.0)		58 (16.1)
Q29_Tired_3mths	170 (49.0)		176 (48.9)
Q31_Tired	73 (21.2)		75 (21.0)
Q32_Tired	172 (49.9)		178 (49.5)
Q33_Tired	162 (47.2)		170 (47.2)
Q34_Tired	152 (44.7)		160 (44.6)
Q35_Tired	113 (32.9)		119 (33.1)
Q36_Tired	50 (14.7)		53 (14.7)
Q38_CghBlood_3mths	92 (26.4)		95 (26.3)
Q40_CghBlood	70 (20.2)		73 (20.3)

Study 2

Q41_CghBlood	38 (11.1)		40 (11.1)
Q42_CghBlood	78 (22.7)		81 (22.6)
Q43_ChInfectn	124 (36.3)		131 (36.3)
Q44_ChInfectn	76 (22.7)		82 (22.8)
Q49_Weight	40 (11.9)		45 (12.5)
Q50_Weight	100 (29.1)		108 (30.0)
Q51_Weight	98 (28.7)		108 (30.1)
Q52_Weight	96 (28.2)		101 (28.0)
Q53_HCswat_3mths	131 (38.0)		136 (38.0)
Q55_HCswat	41 (11.9)		43 (11.9)
Q56_HCswat	137 (39.5)		141 (39.3)
Q57_HCswat	97 (28.3)		101 (28.2)
Q59_EatChnges	48 (14.0)		52 (14.5)
Q60_EatChnges	100 (29.5)		106 (29.4)
Q61_EatChnges	53 (15.7)		58 (16.2)
Q62_EatChnges	61 (19.6)		71 (19.8)
Q65_Voice	72 (20.9)		77 (21.4)
Q66_Voice	28 (8.2)		30 (8.4)

Table 6.19 Frequency distribution for three-response categorical variables

Variable	Complete case				Imputed		
	0	1	2		0	1	2
Q14a_Cgh	166 (51.6)	143 (44.4)	13 (4.0)		186 (51.7)	156 (43.5)	17 (4.8)
Q14b_Cgh	233 (73.5)	60 (18.9)	24 (7.6)		261 (72.7)	69 (19.2)	29 (8.1)
Q14c_Cgh	231 (72.6)	69 (21.7)	18 (5.7)		258 (71.8)	79 (22.0)	22 (6.2)
Q25a_BrChnges	246 (73.4)	56 (16.7)	33 (9.9)		263.5 (73.4)	60 (16.7)	36.5 (9.9)
Q25b_BrChnges	198 (58.2)	132 (38.8)	10 (2.9)		208 (58.0)	139 (38.8)	12 (3.2)
Q25c_BrChnges	206 (61.7)	116 (34.7)	12 (3.6)		221 (61.6)	124 (34.5)	14 (4.0)

Study 2

Q25d_BrChnges	240 (72.1)	64 (19.2)	29 (8.7)		256 (71.4)	74 (20.5)	29 (8.1)
Q46_ChInfectn	250 (73.3)	90 (26.4)	1 (0.3)		264 (73.4)	94 (26.3)	1 (0.3)
Q48_ChInfectn	258 (79.9)	64 (19.8)	1 (0.3)		289 (80.5)	69 (19.2)	1 (0.3)
Q58_HCswat	298 (86.9)	21 (6.1)	24 (7.0)		312 (87.0)	23 (6.3)	24 (6.7)

Table 6.20 Frequency distribution for four- response categorical variables

Variable	Complete case					Imputed			
	0	1	2	3		0	1	2	3
Q45_ChInfectn	138 (40.2)	119 (34.7)	59 (17.2)	27 (7.9)		144 (40.2)	124 (34.5)	62 (17.2)	29 (8.1)
Q47_ChInfectn	122 (36.5)	127 (38.0)	59 (17.7)	26 (7.8)		134 (37.2)	135 (37.6)	62 (17.4)	28 (7.8)

There were hardly any difference between the distributions of the observed, and the imputed data.

6.4.9 Discussion and conclusion to missing data

The aim of the imputation process is to generate a complete dataset to ensure robust analyses and valid results.

Even though we cannot actually test unobserved imputed values for agreement with an unknown true distribution (Abayomi et al. 2006), it is possible to use observed values to discern potential problem with the imputations. Imputed values that were pathologically different from expectations would then be discarded (Abayomi et al. 2006). There were no differences of consequence between the observed and imputed dataset, to indicate that removal was warranted. For the most part, our findings did not indicate obvious discrepancies in the imputations.

Abayomi et al. (2006) emphasised that differences in distribution between the imputed and the observed do not necessarily indicate violations of the missingness assumptions or problems with the imputation model. Some

Study 2

deviations between observed and missing values can be expected under MAR assumptions, but it is deviations that are not consistent with observed departures from MCAR, that require assessment for plausibility.

As at least some of the data were in keeping with MAR assumptions and the distribution of the complete cases (partially observed dataset) were similar to the imputed dataset, the use of the imputed dataset appeared justified. The models generated from the imputed dataset will be compared with those obtained from the complete case analysis, in section 6.4.28.5.

MI is appropriate to handle missing data under both MAR and MCAR assumptions (SSCC, 2013). However, to use complete case analysis when the missing mechanism is not MCAR (ignorable missing mechanism) could lead to highly biased results (Molenberghs and Kenward 2007). Given that some of the data are not MCAR, complete case analysis might not be appropriate here.

Main analysis

6.4.10 Recruitment strategy

The study identified 484 eligible participants attending the lung clinic in Southampton General Hospital who were approached between November 2012 and February 2014. 359 participants gave consent to participate in the IPCARD Chest Clinic Study and 125 eligible attendees of the clinic declined participation, resulting in a response rate of 74.2%. 45 eligible participants were missed in this pragmatic consecutive recruitment process. Recruitment rate was lower than the initial expected rate (67.9%), due to higher number of eligible participants being missed during a busy working clinic, or participants who did not return the questionnaire before leaving the clinic before the researcher can approach them.

6.4.11 Participant Descriptive Data

Table 6.21 shows the demographic details of the 359 participants who had attended the lung clinic in Southampton General Hospital. The median age at time of diagnosis of the 77 participants with LC was 71 years. The non-LC group were slightly younger at 67 years. There were more males ($\approx 60\%$) than females ($\approx 40\%$) in both the case group and the non-case group.

Study 2

Of the 359 participants, 77 were subsequently diagnosed with primary LC. This was 40% lower than the expected proportion of lung cancers diagnosed in the lung-shadow clinic in SGH. The number of people attending the lung-shadow clinic with a known CT scan result was higher than the initial estimate, which had to be excluded in the study based on the inclusion criteria. Other malignant diagnoses included 17 (6%) mesotheliomas, 3 (1.1%) lymphomas, 3 (1.1%) carcinoid tumours, and 17 other cancers (bladder, breast, colon, GI tract, skin, ovarian, pancreatic, and renal) (see Table 6.22). Mesotheliomas were not categorised as LC as they have different pathophysiology; LC develops in the lung whereas mesothelioma affects the lining of the lung.

At follow up, LC histology were obtained where possible. Non-small cell lung cancer (NSCLC) is the more commonly diagnosed lung cancer with 47 cases (61%) compared to the 7 cases (9%) of small cell lung cancer (SCLC), and 30% undefined cancer. There were 27 (35%) early-staged LC (stage I and II), 47 (61%) late-staged LC (stage III and IV), and 3 unknown staging. 16 (21%) of the lung cancer cases were resectable cancers.

Study 2

Table 6.21 Demographic details of participants (n=359)

	LC cases n=77 (%)	Non-LC cases n=282 (%)
AGE (years)		
median	71.0	67.0
mean	70.5	67.5
median (sample population)		68.0
40-49	4 (5.2)	23 (8.2)
50-59	9 (11.7)	56 (19.9)
60-69	18 (23.4)	87 (30.9)
70-79	31 (40.3)	63 (22.3)
80-89	12 (15.6)	43 (15.2)
90-99	3 (3.9)	7 (2.5)
100-	0 (0)	3 (1.1)
Gender		
Female	30 (39.0)	116 (41.1)
Male	47 (61.0)	166 (58.9)
Smoking status		
Ever smoker	69 (89.6)	202 (71.6)
Never smoker	4 (5.2)	72 (25.5)
Current smoker	25 (32.5)	65 (23.0)

Study 2

Table 6.22 Clinical characteristics of participants (n=359); follow up of participant's diagnoses


Other types of cancer	
Bladder cancer	4 (1.4)
Breast cancer	3 (1.1)
Carcinoid tumour	3 (1.1)
Colon carcinoma	1 (0.4)
GI tract cancer	1 (0.4)
Lymphoma	3 (1.1)
Mesothelioma	17 (6.0)
Melanoma	3 (1.1)
Ovarian cancer	2 (0.7)
Pancreatic cancer	1 (0.4)
Renal cancer	2 (0.7)
Non-malignant respiratory and chest-related problems	
Asbestos- related pleural plaques	15 (5.3)
Asthma	24 (8.5)
Bronchitis (Acute)	2 (0.7)
Bronchiectasis	18 (6.4)
Interstitial lung disease	3 (1.1)
Heart disease	13 ()
COPD	124 (34.5)
Pneumonia	10 (35.5)
Types of lung cancer	
<i>NSCLC</i>	47 (61.0)
<i>SCLC</i>	7 (9.1)
<i>undefined LC</i>	23 (29.9)
Tumour Stage	
I	11 (14.3)
II	16 (20.8)
III	20 (26.0)
IV	27 (35.1)
Operable LC	16 (20.8)
Inoperable LC	61 (79.2)

6.4.12 Dichotomising multiple-response variables

Most of the questionnaire items that consisted of more than two response categories have been collapsed into dichotomous variables. With reference to Table 6.1 (Page 140), type 2 and type 3 questions with multiple categories, were dichotomised. Tetrachoric correlation was used to identify the optimum cut-off for distinguishing between LC and not LC.

Section 2 - Cough

Q10 Have you ever had a cough that lasted for more than 3 weeks?

No ☐ Please go to Section 3, page 7 

Yes and I still have ☐ Please go to question 12

Yes but I no longer ☐ → Q11 Have you had a cough in the last three months Yes ☐ No ☐

Please go to question 12 Please go to Section 3, page 7

Q12 Please indicate when you **first** had a cough that lasted for more than 3 weeks.

Within the last 3 months 4-12 months ago More than 12 months ago

☐ ☐ ☐

Q13 Please indicate whether the statements below accurately describe your most recent cough/coughs (**that lasted for more than 3 weeks**) and how often you have had the type of cough described by that statement.

	Never	Once	Occasionally	Most of the time
a) An irritating cough (feels like an irritation in the throat or chest)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
b) A tickly cough	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
c) A cough that starts in the throat	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
d) A cough that feels like clearing the throat	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
e) A wheezy cough	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Figure 6.5 Example of questionnaire structure for multiple response questions

Tetrachoric correlations were performed for all of the symptom variables to explore the optimum cut offs for each variable, see Appendix 15. The values of the correlation are interpreted in the same way as the Pearson correlation; the

Study 2

closer the value is to 1.0, the higher is the correlation and the values close to zero indicate little association between variables (Stata manual 2013).

The tetrachoric correlations of three symptom variables from each generic symptom variable (*an irritating cough, breathlessness after walking a short distance, and feeling tired more easily than used to*) were reported and discussed in detail. Results of the remaining variables can be found in Appendix 15. Figure 6.6 to Figure 6.8 show the Stata outputs of the following tetrachoric correlations. Three possible variables with different response cut-offs, labelled Q13a_Cgh_1, Q13a_Cgh_2, and Q13a_Cgh_3, were generated for the variable, Q13a_Cgh. Q13a_Cgh_1 collapsed responses, 'once', 'occasionally', and 'most of the time' into one category to indicate presence of that symptom. Q13a_Cgh_2 collapsed 'occasionally', and 'most of the time' to form the positive response (coded 1), whilst Q13a_Cgh_3 only captured those experiencing highest frequency (most of the time) of that symptom.

Results showed that effect sizes of the varying categorisation of the response variables varied slightly. Although in all three examples, the third cut-off appears to have the largest correlation with lung cancer diagnosis. At this cut-off, the sensitivity will be improved, and ORs will be consequently lower. This might restrict the number of variables that enter the final model. As this is an exploratory study, it might not be beneficial to lose too many symptom variables at this stage, as they could be potentially important variables for lung cancer diagnosis. Therefore, the cut-off with the slightly poorer correlation, and sensitivity was applied (variable_1).

Figure 6.6: Tetrachoric correlation for cough symptom, Q13a_Cgh

```
. tetrachoric Q13a_Cgh_1 Q13a_Cgh_2 Q13a_Cgh_3 LCdiagnosis
(obs=336)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0069
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q13a_C~1	Q13a_C~2	Q13a_C~3	LCdiag~s
Q13a_Cgh_1	1.0000			
Q13a_Cgh_2	1.0000	1.0000		
Q13a_Cgh_3	1.0000	1.0000	1.0000	
LCdiagnosis	0.0045	-0.0582	0.0601	1.0000

Study 2

Figure 6.7: Tetrachoric correlation for breathing changes symptom, Q22a_BrChnges

```
. tetrachoric Q22a_BrChnges_1 Q22a_BrChnges_2 Q22a_BrChnges_3 LCdiagnosis  
(obs=340)
```

```
matrix with tetrachoric correlations is not positive semidefinite;  
it has 1 negative eigenvalue  
maxdiff(corr,adj-corr) = 0.0006  
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q22a_B~1	Q22a_B~2	Q22a_B~3	LCdiag~s
Q22a_BrChn~1	1.0000			
Q22a_BrChn~2	1.0000	1.0000		
Q22a_BrChn~3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0811	-0.0855	-0.1161	1.0000

Figure 6.8: Tetrachoric correlation for tiredness symptom, Q32_Tired

```
. tetrachoric Q32_Tired_1 Q32_Tired_2 Q32_Tired_3 LCdiagnosis  
(obs=345)
```

```
matrix with tetrachoric correlations is not positive semidefinite;  
it has 1 negative eigenvalue  
maxdiff(corr,adj-corr) = 0.0018  
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q32_Ti~1	Q32_Ti~2	Q32_Ti~3	LCdiag~s
Q32_Tired_1	1.0000			
Q32_Tired_2	1.0000	1.0000		
Q32_Tired_3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0222	-0.0116	-0.0722	1.0000

Results of the tetrachoric correlations also determined the better cut-off for the generic symptoms variables to be 'current or in the last three months' (refer to Appendix 15). Therefore, the following analyses (univariate, bivariate, and multivariate) were carried out at this level.

Full population

6.4.13 Univariate analysis of the relationship between symptoms and lung cancer

Variables that were significantly associated with the outcome (lung cancer) at the $p\text{-value} < 0.05$ and $p < 0.15$ were tabulated, see Table 6.23 and Table 6.24. The frequencies of each symptom and the cross-tabulation (2x2 tables) of these symptoms with the binary dependent variable (LC case/non-LC case) are shown in the following tables. Effect estimates were presented as odds ratios.

Study 2

Generic symptom variables experienced in the last three months prior to clinic appointment were assessed. Univariate analyses of all the symptom variables can be found in Table 6.52.

Two variables, cough that lasted for more than 3 weeks first indicated in the last 3 months (Q12_Cgh) and breathing difficulties or changes first indicated in the last 3 months (Q21_BrChnges) were associated with LC at a statistically significant level, $p < 0.05$. 39.5% of participants with LC had reported first experiencing their coughs in the last three months compared with 26% of those without LC. This difference was statistically significant ($p=0.026$) with an OR of almost 2.0. Breathing changes that were first indicated in the last three months, (Q21_BrChnges), also predicted lung cancer diagnosis with OR of approximately 2.0 ($p\text{-value} = 0.021$).

A total of 12 variables were selected at $p\text{-value} < 0.15$, or had $OR > 2.0$ or $OR < 0.5$, to explore potential associations with lung cancer.

Table 6.23 Univariate analysis of the relationship between symptoms and lung cancer (p<0.05)

Questionnaire items	Variable	(% presented with this variable)		Odds Ratio	p-value
		LC cases (n=77)	Non-LC cases (n=282)		
When did you first have a cough that lasted for more than 3 weeks (Within the last 3 months/Not)	Q12_Cgh_R	39.5	26.1	1.85	0.026*
When did you first have breathing difficulties/changes? (Within the last 3 months/Not)	Q21_BrChnges	37.1	23.5	1.92	0.021*

Study 2

Table 6.24 Univariate analysis of the relationship between symptoms and lung cancer (OR>2.0 or <0.5, or p<0.15)

Questionnaire items	Variable	Number (%) presented with this variable		Odds Ratio	p-value
		LC cases	Non-LC cases		
Ache or pain in the side of chest or ribs	Q3g_Pain	25.1	34.6	0.6	0.119
When did you first had a cough that lasted for more than 3 weeks (Within the last 3 months/Not)	Q12_Cgh_R	39.5	26.1	1.8	0.026
A tickly cough	Q13b_Cgh	56.0	46.4	1.5	0.149
A hard or harsh cough without phlegm (Yes/No)	Q13k_Cgh	25.1	36.9	0.6	0.064
When did you first had breathing difficulties/changes? (Within the last 3 months/Not)	Q21_BrChnges	37.1	23.5	1.9	0.021

Study 2

Have you experienced breathing problems that are only present or get worse at certain times of the year?	Q23_BrChnges	18.4	11.5	1.7	0.132
Is your breathlessness worse than it was 3 months ago?	Q26_BrChnges	44.7	34.9	1.5	0.131
Have you had noticeably more chest infections within the last 12 months than the year before?	Q46_ChInfectn	32.5	24.6	1.6	0.077
Have you gained weight within the last 12 months?	Q52_Weight	20.4	30.1	0.6	0.111
Have you experienced hot or cold sweats in the day?	Q57_HCswat	21.4	30	0.6	0.147
Has your appetite increased within the last 12 months?	Q59_EatChnges	8.8	16.0	0.5	0.126
Have you currently gone off certain foods you used to eat?	Q62_EatChnges	26.6	18.0	1.7	0.106

Study 2

Accepting a lower level of significance ($p < 0.15$) in the univariate analysis allowed for systemic symptom variables, such as eating changes, and weight, to be explored in the multivariate model. Although these variables did not reach traditionally accepted significance level of $p < 0.05$, they could still be potentially important predictors for lung cancer, and would inform future investigations when there is a larger sample. As this analysis was exploratory rather than explanatory, excluding too many variables for further study might not be very informative.

Univariate analysis of the relationship between risk variables and lung cancer showed statistically significant associations between previous cancer and lung cancer, and previous smoker and lung cancer (see Table 6.25). All of the risk factors, except for pneumonia (in the last five years) were higher percentage in the LC group than the non-LC group, $OR < 1.0$.

Table 6.25 Univariate associations between risk factors and lung cancer diagnosis

Risk Factor	Variable	(% presented with this variable)		Odds Ratio	p-value
		LC cases (n=77)	Non-LC cases (n=282)		
Pneumonia in the last 5 years	Q69a_Pneumo_5yrs	5 (6.5)	33 (11.47)	0.7	0.325
COPD	Q69e_COPD_ever	17 (22.1)	43 (15.2)	1.4	0.245
Cancer	Q69h_Cancer_ever	19 (24.7)	29 (10.3)	2.4	0.003*
Asbestos	Q69j_Asbес_ever	6 (7.8)	19 (6.7)	1.1	0.849
Ever smoker	Q73_Smoke	69 (89.6)	202 (71.6)	6.3	0.001*
Family history of lung cancer	Q71d_FamHx	16 (20.8)	50 (17.7)	1.3	0.469

6.4.14 Bivariate analyses (Comorbidities)

Comorbidities were investigated in the bivariate analyses if there were reasons to suspect confounding. Distributions of recent pneumonia (in the last three months), asthma, COPD, and arthritis differed between the LC and non-LC group, and therefore, were analysed for potential confounding (see Table 6.26). There was more asthma, arthritis and recent pneumonia, in the non-LC population.

Table 6.26 shows the distributions of patient-reported comorbidities in each group; LC and non-LC. Relative differences in the LC and non-LC group were observed in comorbidities; Q9a_Pneumo_last3mths, Q69c_Asthma_ever, Q69e_COPD_ever, and Q69k_Arthritis_ever. Only 14% of individuals in the LC group had asthma in comparison to the 20% of the non-LC group.

Bivariate analyses were carried out using the Mantel-Haenszel (M-H) analysis. The Mantel-Haenszel method is a technique that generates an estimate of an association between an exposure and an outcome after adjusting for or taking into account confounding (Mantel and Haenszel 1959; McDonald 2014). Data were stratified into two or more levels of a potential confounding factor; a series of 2x2 tables showing the association between the risk factor and outcome at two or more levels of the confounding factor were created. A weighted average of the odds ratios across the strata was then computed (McDonald 2014).

Table 6.26 Patient-reported comorbidities (full population data)

	LC cases n=77 (%)	Non-LC cases n=282 (%)	p-value
Patient-reported comorbidities			
Q69a_Pneumo_last3mths	1 (1.3)	9 (3.2)	0.373
Q69c_Asthma_ever	11 (14.3)	57 (20.2)	0.209
Q69d_Allergy_ever	13 (16.9)	46 (16.3)	0.943
Q69e_COPD_ever	17 (22.1)	43 (15.2)	0.155
Q69f_HD_Angina_ever	16 (20.8)	60 (21.3)	0.847
Q69k_Arthritis_ever	20 (26.0)	106 (37.6)	0.039*

Study 2

Asthma

The American College of Asthma, Allergy and Immunology [ACAAI] (2010) described coughing, breathing difficulties, tightness in the chest, and/or wheezing as the most common symptoms of asthma. Similar symptom variables for LC (cough that lasted for three weeks in the last three months, breathing changes in the last three months, and breathing symptom descriptors such as tightness in the chest, and wheezing) were investigated for possible confounding. There were no observable differences between the crude and adjusted ORs in the bivariate analyses for asthma.

Table 6.27 M-H χ^2 test for asthma

	Asthma	OR	Confidence Interval	
			95%	
Q10_Cgh_3mths	Crude	1.481	0.787	2.893
	M-H combined	1.495	0.815	2.743
Q19_BrChnges_3mths	Crude	1.105	0.578	2.188
	M-H combined	1.167	0.625	2.179
Q24c_BrChnges (<i>Tightness in chest</i>)	Crude	0.932	0.507	1.687
	M-H combined	0.990	0.557	1.763
Q25b_BrChnges (<i>Wheezing noise breathing in</i>)	Crude	0.926	0.507	1.667
	M-H combined	0.954	0.544	1.678
Q25c_BrChnges (<i>Wheezing noise breathing out</i>)	Crude	0.990	0.532	1.807
	M-H combined	1.000	0.567	1.766

COPD

There was an uneven distribution between COPD in the LC and non-LC group in the direction that would be expected for a risk factor (i.e. more COPD in the LC group), which suggested that COPD could be a confounder as well as a risk factor.

Study 2

According to the International Primary Care Respiratory Group guidelines, symptoms of COPD include persistent coughing, productive cough, breathing difficulties, and wheezing (Levy et al. 2006). Similar symptoms for LC (coughing that lasted for more than three weeks, cough that usually produces phlegm in the morning, cough that usually produces phlegm at any time of the day, breathing changes in the last three months, wheezing noise when breathing out, and when breathing in) were explored for confounding (see Table 6.28). No observable differences between the crude and adjusted ORs in the bivariate analyses for COPD were found to indicate any potential confounding.

Study 2

Table 6.28 M-H χ^2 test for COPD

	COPD	OR	Confidence Interval 95%	
Q10_Cgh_3mths	Crude	1.506	0.799	2.945
	M-H combined	1.512	0.826	2.766
Q13i_Cgh (A cough that usually produces phlegm in the morning)	Crude	1.122	0.635	1.991
	M-H combined	1.123	0.657	1.920
13j_Cgh (A cough that usually produces phlegm at any time of day)	Crude	1.360	0.768	2.431
	M-H combined	1.341	0.782	2.300
Q19_BrChnges_3mths	Crude	1.100	0.574	2.182
	M-H combined	1.046	0.556	1.967
Q25b_BrChnges (Wheezing noise when breathing out)	Crude	0.932	0.515	1.668
	M-H combined	0.878	0.502	1.539
Q25c_BrChnges (Wheezing noise when breathing in)	Crude	0.959	0.523	1.734
	M-H combined	0.904	0.510	1.600

Arthritis

Arthritis is generally indicated by aches, and pain in the joints. Symptoms such as upper body discomfort or pain in the upper body, chest, or shoulders, any new joint pain in the last 12 months, and any new unusual sensations in the last 12 months, were investigated for potential confounding. Bivariate analyses for arthritis showed no observable differences between the unadjusted and adjusted ORs that might suggest confounding (see Table 6.29)

Study 2

Table 6.29 M-H χ^2 test for arthritis

	Arthritis	OR	Confidence Interval	
			95%	
Q1_Pain_3mths	Crude	0.786	0.445	1.400
	M-H combined	0.764	0.445	1.313
Q63_New_JPain_12mths	Crude	1.079	0.559	2.021
	M-H combined	1.106	0.606	2.017
Q64_New_JPain_12mths	Crude	1.019	0.505	1.975
	M-H combined	1.059	0.560	2.001

Pneumonia

The numbers of cases with current pneumonia (pneumonia in the last three months) were too small (n=10) to draw meaningful conclusions from the bivariate analysis.

6.4.15 Multivariate analysis: Full population data

Variables were selected for the multivariate models using two separate entry criteria; a conservative criteria ($p < 0.05$), and a relaxed criteria ($p < 0.15$), to explore potentially important symptoms that might predict lung cancer.

Variable selection at $p < 0.05$

6.4.15.1 Model 1; Symptoms, adjusted for age ($p < 0.05$)

The two variables that were selected at the significance level of $p < 0.05$, were entered into the first main effects model using forward stepwise regression,. Only breathing changes/difficulties that were first indicated in the last three months (Q21_BrChnges) remained in the model (see Table 6.30). Adding the discarded variable, cough that was first indicated in the last three months (Q12_Cgh), back into the model improved the fit by at two (the minimum difference) based on Akaike's information criterion (AIC).

Study 2

Table 6.30 Model 1: Main effects model ($p < 0.05$), adjusted for age; using forward stepwise regression at $p(e) = 0.05$ and $p(r) = 0.10$

LCdiagnosis	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
AGE	1.130838	.1550448	0.90	0.370	.8643616	1.479467
AGEsq	.9993325	.0009759	-0.68	0.494	.9974216	1.001247
Q12_Cgh_R	1.735646	.5740815	1.67	0.096	.9076496	3.318977
Q21_BrChnges_R	1.858637	.6195869	1.86	0.063	.9670344	3.572295
_cons	.0010706	.0050853	-1.44	0.150	9.69e-08	11.82613

Akaike's information criterion and Bayesian information criterion

Model	Obs	ll (null)	ll (model)	df	AIC	BIC
.	5155	-138.1075	-131.7718	5	273.5436	306.2822

6.4.15.2 Model 2; Symptoms, adjusted for age, and risk variables ($p < 0.05$)

In addition to the variables for breathing changes and cough that were entered into model 1, two risk variables significantly associated with LC in the univariate analysis were added, in model 2. The risk variables were previous cancer (Q69h_Cancer_ever) and ever smoked (Q73_Smoke). The risk variable, ever smoked (Q73_Smoke) remained in model 2 (Table 6.31). The lower AIC score suggested that the resultant model 2 had better fit compared to model 1.

Table 6.31 Model 2: Main effects model with risk variables ($p < 0.05$); using forward stepwise regression; $p(e) = 0.05$ and $p(r) = 0.10$

LCdiagnosis	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
AGE	1.134789	.1573333	0.91	0.362	.8647684	1.489122
AGEsq	.9992794	.0009887	-0.73	0.466	.9973434	1.001219
Q12_Cgh_R	1.790579	.6045582	1.73	0.084	.9238471	3.47046
Q21_BrChnges_R	1.884903	.6441928	1.85	0.064	.9646698	3.682981
Q73_Smoke_	3.372939	1.263673	3.25	0.001	1.618473	7.029291
_cons	.0003641	.0017482	-1.65	0.099	2.98e-08	4.450023

Study 2

Akaike's information criterion and Bayesian information criterion

Model	Obs	ll (null)	ll (model)	df	AIC	BIC
.	5155	-138.1075	-125.6127	6	263.2254	302.5117

The discarded variable, previous history of cancer (Q69h_Cancer_ever) was checked against the final model but did not improve goodness of fit.

6.4.16 Developing a set of diagnostic criteria ($p < 0.05$)

In a similar study on ovarian cancer, Bankhead (2005) used a simple cumulative scoring system for a set of diagnostic factors to predict ovarian cancer. The same approach was used to develop a score,

The set of diagnostic criteria was developed for the best fit model. Therefore, results from the second modelling procedure with the risk variable (model 2) were used to develop this set of diagnostic criteria, which was as follows:

- Age was dichotomised at 71 years old, such that those aged 71 and above scored one point. This cut-off age was decided on the basis of the median age of the LC cases and the non-LC cases, which were 68 and 73, respectively. The mean age of the two groups is 71 years.
- Any one of the two symptoms; breathing changes/difficulties that was first indicated within the last three months (Q21_BrChnges), and coughing (for more than three weeks) that were first indicated within the last three months (Q12_Cgh) present contributed a score of one to the cumulative score.
- Those who had previously smoked also scored one point, so anyone who responded 'yes' in the ever smoked variable (Q73_Smoke_).

The total number of criteria satisfied in the case and non-case group would indicate the total number of factors that were present in each group. A maximum score of four was obtainable. A score of two out of four (2/4) represent participants who experienced at least two of the maximum four possible factors.

Study 2

Sensitivity, specificity and likelihood ratios were calculated for all levels of cut-offs, see Table 6.32. The distribution of cases and non-cases were larger than the actual sample size of 359 because it included imputed data that had been converted to 'mlong' (marginal long) data structure. Technically, this should not affect the ratio of cases and non-cases, which is needed in the calculation of diagnostic accuracy for each cut-off level.

Table 6.32 Analysis of the diagnostic performance of this set of criteria
($p < 0.05$)

Criteria		LC cases	Non-LC cases	Sensitivity	Specificity	Youden index	1-Specificity	+ve Likelihood Ratio	-ve Likelihood Ratio
No. of criteria (max 4)	0.0	1108	3821	1.00	0.0	0.0	1.00	1.0	-
	1.0	1087	3455	0.981	0.096	0.08	0.904	1.08	0.20
	2.0	912	2108	0.823	0.448	0.32	0.552	1.49	0.39
	3.0	526	605	0.475	0.842	0.28	0.158	3.00	0.62
	4.0	115	65	0.104	0.983	0.09	0.017	6.10	0.91

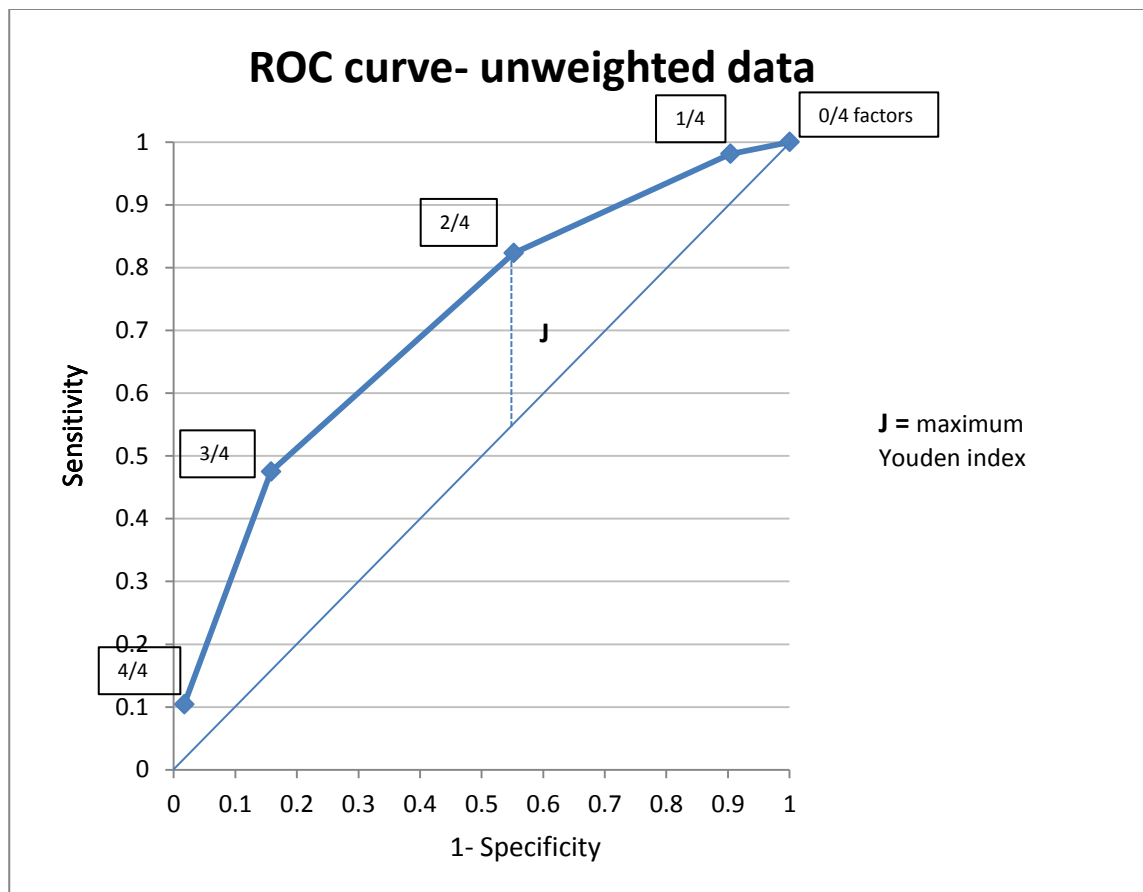


Figure 6.9 Receiver operating characteristic curve for the un-weighted set of criteria ($p < 0.05$)

ROC curve was used to gauge how well the test performs at specific cut-off points in a population (Altman and Bland 1994c). The determination of the optimum trade-off point will depend on the situation, requirement, and implications of misclassification. Youden's J statistic or Youden's index is often used in conjunction with ROC curves as a good indicator of the performance of the diagnostic test at each cut-off point (Youden 1950). It may be used as a criterion for the **optimum cut-off point** to be selected, also indicated by the cut-off value with the shortest Euclidean distance between the ROC curve and the upper left corner of the graph (Youden 1950). By this method, the optimum cut-off point for the ROC curve in Figure 6.9 was suggested at 2/4.

The use of likelihood ratios (positive and negative) is another way to express diagnostic accuracy or measure the power of a diagnostic test in increasing or decreasing the likelihood of a disease (CEBM 2014). Positive likelihood ratio or likelihood ratio of a positive test is the ratio of the probability that the positive

Study 2

result is correct to the probability that the positive test result is incorrect, and conversely, the negative likelihood ratio or likelihood ratio of a negative test is the ratio of the probability that the negative result is incorrect to the probability that the negative result is correct (Attia 2003). The larger the likelihood ratio, the greater is the likelihood of disease, and similarly, the smaller the negative likelihood ratio, the lesser the likelihood of disease (Attia 2003).

The positive likelihood ratio at the optimum cut-off level (2/4) was 1.49. This means that of those referred to the clinic with any two of the four criteria for lung cancer investigation, 82.3% of the lung cancers would be correctly identified with lung cancer (sensitivity) but 55.2% would be incorrectly ruled in. The area under curve (AUC) statistic used to measure the overall discriminatory power was 0.663, which would be considered relatively poor (Fan et al.2006), according to standard assessments.

At the cut-off of 4/4, the highest positive likelihood ratio was achieved at 6.10, above the recommended ≥ 5.0 for strong diagnostic evidence. Sensitivity at this cut-off was only 10%.

6.4.17 Weighted set of diagnostic criteria

The variables were forced into a regression model to obtain the log odds ratio, which were then rounded up to produce the weights shown below. These weights were applied to each criteria to form a weighted symptom score.

Variables	Log Odds Ratio (β coefficient)	Rounded up weight for each criterion
Age (>71)	+0.980	+1
Q12_Cgh_R	+0.690	+1
Q21_BrChnges_R	+0.620	+1
Q73_Smoke_	+1.278	+1

Adding the scores, the minimum obtainable score for the weighted criteria is 0 and the maximum is +4, which is the same as the un-weighted set of

Study 2

diagnostic criteria. Their diagnostic performances would therefore, remain the same, and there was no advantage of a weighted score.

Variable selection at $p < 0.15$

6.4.17.1 Model 1; Symptoms, adjusted for age ($p < 0.15$)

12 symptom variables were identified at a significance level of $p < 0.15$ in the univariate analysis, and added into the forward stepwise regression model, adjusted for age variables (age and age square). The entry and exit criteria for the forward stepwise model were relaxed to $p = 0.1$ and $p = 0.15$, respectively.

The resultant model 1 ($p < 0.15$) contained five symptom variables; coughing (for more than three weeks) that were first indicated within the last three months (Q12_Cgh), breathing changes that were first indicated within the last three months (Q21_BrChnges), a hard or harsh cough without phlegm (Q13k_Cgh), experience noticeably more chest infections in the last 12 months than the years before (Q46_ChInfectn), and weight gain in the last 12 months (Q52_Weight); see Table 6.33. Using the AIC score to determine the fit of the model, discarded variables ($p < 0.15$) were checked against the final model.

Table 6.33 Model 1: Main effects model ($p < 0.15$), adjusted for age; using forward stepwise regression; $p(e) = 0.10$ and $p(r) = 0.15$

LCdiagnosis	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
AGE	1.152799	.1658485	0.99	0.323	.86955	1.528314
AGEsq	.9991425	.0010235	-0.84	0.402	.9971385	1.001151
Q21_BrChnges_R	1.793639	.6251553	1.68	0.094	.9058513	3.55151
Q13k_Cgh_R	.4367246	.1570751	-2.30	0.021	.2158039	.8838042
Q12_Cgh_R	2.072139	.7261332	2.08	0.038	1.042642	4.118156
Q46_ChInfectn	1.950321	.6630441	1.96	0.049	1.001673	3.797399
Q52_Weight	.4771828	.1849239	-1.91	0.056	.2232628	1.019889
_cons	.0009118	.0045627	-1.40	0.162	5.02e-08	16.56485

Study 2

Akaike's information criterion and Bayesian information criterion

Model	Obs	ll (null)	ll (model)	df	AIC	BIC
.	5055	-135.5945	-123.2063	8	262.4126	314.6377

6.4.17.2 Model 2; Symptoms; adjusted for age and risk variables ($p < 0.15$)

Similarly, risk variables that were associated with lung cancer; previous history of cancer (Q69h_Cancer_ever) and previous smoking variable (Q73_Smoke_), and symptom variables from the univariate associations with $p < 0.15$, were added into the forward stepwise regression model.

The resultant model 2 ($p < 0.15$) contained five symptom variables; breathing changes/difficulties that was first indicated within the last three months (Q21_BrChnges), coughing (for more than three weeks) that were first indicated within the last three months (Q12_Cgh), a hard or harsh cough without phlegm (Q13k_Cgh), experience noticeably more chest infections in the last 12 months than the years before (Q46_ChInfectn), and increased appetite (Q59_EatChnges), and one risk variable, previous history of smoking (Q73_Smoke), see Table 6.34). The model was also checked against the variables that were dropped.

The variable, weight gain in the last 12 months (Q52_Weight) in Model 1 was replaced by increased appetite (Q59_EatChnges) in Model 2. Also, Model 2 had better fit than Model 1 (without the risk variables).

Study 2

Table 6.34 Model 2: Main effects model ($p < 0.15$) with risk variables; using forward stepwise regression at $p(e) = 0.10$ and $p(r) = 0.15$

LCdiagnosis	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
AGE	1.119502	.1637095	0.77	0.440	.8405235	1.491075
AGEsq	.9993315	.0010392	-0.64	0.520	.9972968	1.00137
Q73_Smoke_	3.685254	1.499058	3.21	0.001	1.660437	8.179227
Q21_BrChnges_R	1.814218	.6444097	1.68	0.094	.9043637	3.639453
Q59_EatChnges	.3081691	.1852356	-1.96	0.050	.0948736	1.000997
Q46_ChInfectn	1.820837	.6270715	1.74	0.082	.9271044	3.57613
Q12_Cgh_R	1.944038	.6887329	1.88	0.061	.9708257	3.892853
Q13k_Cgh_R	.5162156	.1874102	-1.82	0.069	.2533994	1.051615
_cons	.0007875	.0040107	-1.40	0.161	3.64e-08	17.02962

Akaike's information criterion and Bayesian information criterion

Model	Obs	ll (null)	ll (model)	df	AIC	BIC
.	5055	-135.5945	-117.5937	9	253.1873	311.9405

6.4.18 Developing a set of diagnostic criteria ($p < 0.15$)

Both model 1 and 2 developed at the relaxed entry and exit criteria, included five symptom variables, but in the risk model (Model 2), the absence of increased appetite replaced the absence of weight gain in Model 1 with the addition of previous smoking. For the purpose of developing a set of diagnostic criteria, the model with the better fit (risk model, Model 2) was chosen; with Akaike's score difference of 9. The set of criteria was as follows:

- Age was dichotomised at 71 years old, such that those aged 71 and above scored one point. This cut-off age was decided on the basis of the median age of the LC cases and the non-LC cases, which were 68 and 73, respectively. The mean age of the two groups is 71 years.
- A hard or harsh cough without phlegm (Q13k_Cgh), and increased appetite (Q59_EatChnges), which occurred significantly more frequently in those without lung cancer, were recoded so that the absence of either symptom presented a score of 1.

Study 2

- The remaining three symptoms; breathing changes/difficulties that was first indicated within the last three months (Q21_BrChnges), coughing (for more than three weeks) that were first indicated within the last three months (Q12_Cgh), and experience noticeably more chest infections in the last 12 months than the years before (Q46_ChInfectn), present contributed a score of one to the cumulative score.
- Those who had previously smoked also scored one point, so anyone who responded 'yes' in the ever smoked variable (Q73_Smoke_).

A maximum score of seven was obtainable. Table 6.35 shows the diagnostic performance (sensitivity, specificity, and likelihood ratios) of the set of criteria based on the model developed to predict LC diagnosis.

Table 6.35 Analysis of the diagnostic performance of this set of criteria (p<0.15)

Criteria		LC cases	Non-LC cases	Sensitivity	Specificity	Youden index	1-Specificity	+ve Likelihood Ratio	-ve Likelihood Ratio
No. of criteria (max 4)	0.0	1162	3954	1.00	0.0	0.0	1.00	1.0	-
	1.0	1162	3912	1.00	0.011	0.01	0.989	1.06	0.0
	2.0	1162	3757	1.00	0.050	0.05	0.95	1.05	0.0
	3.0	1104	3010	0.950	0.239	0.19	0.761	1.25	0.21
	4.0	929	1801	0.800	0.545	0.35	0.455	1.76	0.37
	5.0	520	561	0.448	0.858	0.31	0.142	3.15	0.64
	6.0	250	91	0.215	0.977	0.19	0.023	9.35	0.80
	7.0	58	35	0.005	0.991	0.0	0.009	5.64	0.96

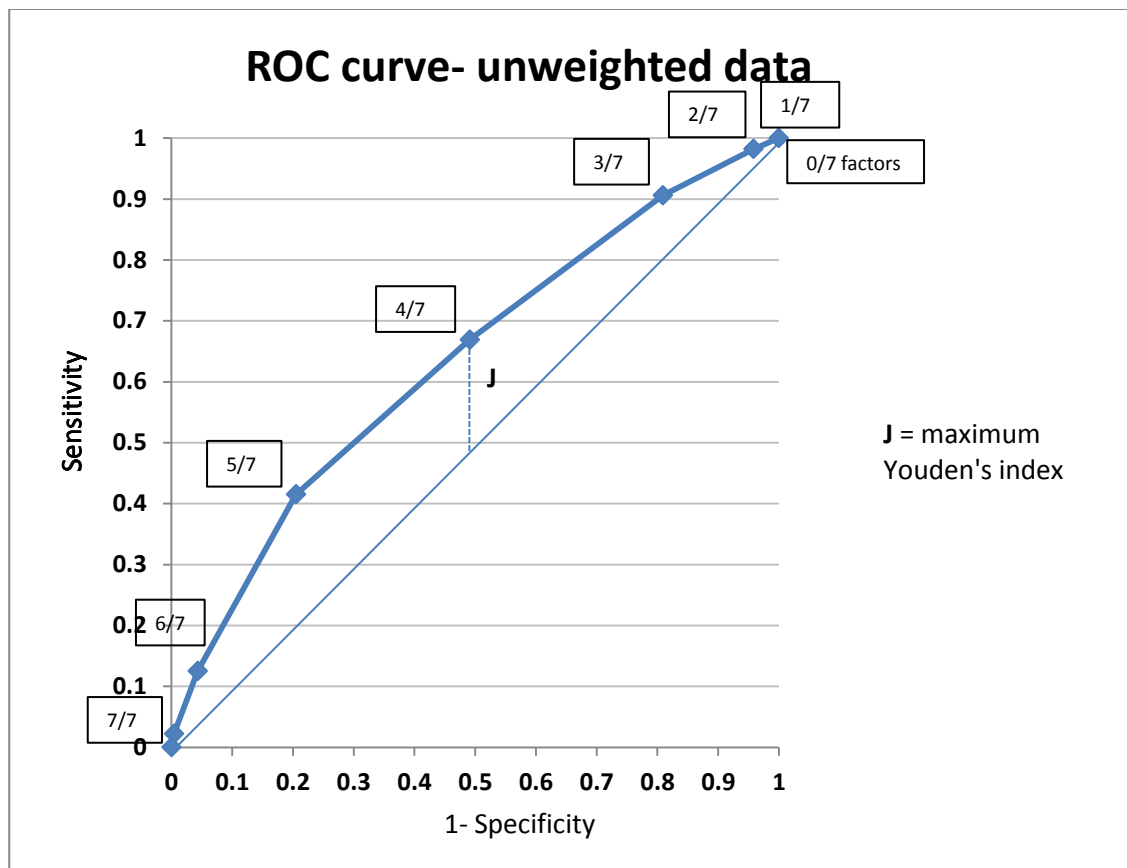


Figure 6.10: Receiver operating characteristic curve for the un-weighted set of criteria

Similarly, the Youden's index was used to determine the optimum level for cut-off on the ROC curve, also indicated by the cut-off value with the shortest Euclidean distance marked 'J' (Youden 1950). It can be seen from the ROC curve (Figure 6.10) that the optimal cut-off is at 4/7, where those referred to the clinic with four criteria were investigated, 80% of the lung cancers would be correctly identified. Sensitivity was reasonably high but specificity was lower, 54.5%. Likewise, positive likelihood ratios were used to measure and express diagnostic accuracy of the test. At the optimal cut-off level, a positive likelihood ratio of 1.76 was obtained. The AUC statistic was 0.775, indicating fairly good discriminatory power.

The highest positive likelihood ratio was achieved at a cut off of 6/7 with positive likelihood ratio of 9.35, above the recommended threshold for strong diagnostic evidence. However, the sensitivity is lower at this cut-off (21.5%).

6.4.19 Weighted set of diagnostic criteria

A weighted symptom score was also generated from the same logistic regression model. Weights obtained from the log odds ratios of the six variables (age remained dichotomous) were applied to form this weighted set of diagnostic criteria. The weights were rounded up to the next whole number and applied to each criterion.

Variables	Log Odds Ratio (β coefficient)	Rounded up weight for each criterion
Age (>71)	+0.867	+1
Q12_Cgh_R	+0.764	+1
Q13k_Cgh_R	+0.643	+1
Q21_BrChnges_R	+0.594	+1
Q46_ChInfectn	+0.560	+1
Q59_EatChnges	+1.205	+1
Q73_Smoke	+1.377	+1

The minimum obtainable score of the weighted set of criteria is again 0 and the maximum is +7. The diagnostic performance (sensitivity, specificity, and likelihood ratios) of the weighted set of diagnostic criteria would be the same as the un-weighted set of diagnostic criteria. Therefore, there is no additional advantage to using the weighted set of criteria over the un-weighted criteria.

COPD

6.4.20 Sub-group analysis: COPD population

There were 124 participants with COPD in the population recruited (35% of those referred to the chest clinic). Approximately 31% (38 participants) of this group had LC. It was initially anticipated that approximately 40% of this group would have LC. The number of participants recruited with COPD and the proportion of those with COPD found to have lung cancer, were lower than predicted. The sample size previously calculated 96 lung cancer cases would

Study 2

have a power of 80% to detect a change in symptom frequency from 20% in the non-LC group to 40% in the LC group. Therefore, 38 lung cancers in 124 participants with COPD, provides less power than previously planned. Nevertheless, the identification of symptoms strongly associated with LC in this sub-population would still inform the design of future studies.

An exploratory sub-group analysis was carried out to examine whether there were any differences in symptom distribution between LC cases and non-LC cases with COPD, and to identify symptoms that distinguish between LC and non-LC cases in a sub-population with COPD. The same modelling strategy used in the full population dataset was used.

6.4.21 Univariate analyses: COPD sub-population

Three variables were statistically significant at $p < 0.05$; unable to get enough air (Q24b_BrChnges), wheezing sensation when in a particular position (Q25d_BrChnges), and breathlessness worse than in the last three months (Q26_BrChnges), see Table 6.36.

In addition, variables that met a relaxed significance level of $p < 0.15$, were also identified to be further investigated in the multivariate analysis (refer to Table 6.37). The identification of symptoms strongly associated with LC ($OR > 2$ or < 0.5) but with an alpha of > 0.05 (statistically not significant) would still inform the design of future studies (Sterne and Kirkwood 2003).

Study 2

Table 6.36 Univariate analysis of the relationship between symptoms and lung cancer in a population with COPD (OR>2.0, or <0.5, or p<0.05)

Questionnaire items	Variable	Number (%) presented with this variable		Odds Ratio	p-value
		LC cases	Non-LC cases		
Unable to get enough air	Q24b_BrChnges	61.8	40.6	2.37	0.033*
Wheezing sensation when in a particular position	Q25d_BrChnges	30.0	17.9	2.43	0.002*
In general is your breathlessness worse than it was 3 months ago?	Q26_BrChnges	54.3	35.2	2.19	0.050*

*p<0.05 statistically significant

Study 2

Table 6.37 Univariate analysis of the relationship between symptoms and lung cancer in a population with COPD (OR>2.0 or <0.5, or p<0.15)

Questionnaire items	Variable	Number (%) presented with this variable		Odds Ratio	p-value
		LC cases	Non-LC cases		
A tickly cough	Q13b_Cgh	65.4	48.5	2.00	0.089
Have you ever experienced breathing difficulties/changes in the last 3 months?	Q19_BrChnges_3mths	85.7	74.9	2.03	0.215
When did you first had breathing difficulties/changes? (Within the last 3 months/Not)	Q21_BrChnges	29.7	16.2	2.19	0.094
Wheezing noise when breathing in	Q25b_BrChnges	56.8	42.2	1.70	0.142
Wheezing noise when breathing out	Q25c_BrChnges	50.9	37.2	1.89	0.068
Is your tiredness worse than it was 3 months ago?	Q36_Tired	10.5	20.6	0.45	0.183

Study 2

Have you currently gone off certain foods you used to eat?	Q62_EatChnges	28.0	12.7	2.71	0.056
---	----------------------	------	------	------	-------

Table 6.38 Univariate analysis of the relationship between risk variables and lung cancer in a population with COPD

Risk Factor		Number (%) presented with this variable		Odds Ratio	p-value
		LC cases	Non-LC cases		
Pneumonia in the last 5 years	Q69a_Pneumo_5yrs	7.9	17.4	0.68	0.327
Cancer	Q69h_Cancer_ever	28.9	5.8	1.27	0.408
Asbestos	Q69j_Asbes_ever	7.9	11.6	0.533	0.142
Ever smoker	Q73_Smoke	100.0	95.2	2.49	0.231
Family history of lung cancer	Q71d_FamHx	27.3	17.6	1.76	0.255

6.4.22 Bivariate analysis: COPD sub-population

Similarly to the bivariate analysis for the full population data, comorbidities were investigated in the bivariate analyses if there were reasons to suggest confounding; i.e. if distributions of number of comorbidities differed between the LC and non-LC group. Asthma, allergy, angina, and arthritis were investigated for potential confounding. The distribution of patient-reported comorbidities can be found in Table 6.39.

Table 6.39 Patient-reported comorbidities (sub-population data)

	LC cases n=38 (%)	Non-LC cases n=86 (%)	p-value
Patient-reported comorbidities			
Q69a_Pneumo_last3mths	1 (2.6)	2 (2.3)	0.922
Q69c_Asthma_ever	11 (28.9)	21 (24.4)	0.785
Q69d_Allergy_ever	8 (21.1)	15 (17.4)	0.736
Q69f_HD_Angina_ever	9 (23.7)	16 (18.6)	0.697
Q69k_Arthritis_ever	10 (26.3)	32 (37.2)	0.122

Asthma

Symptom variables for LC that were similar to symptoms of asthma (cough that lasted for three weeks in the last three months, breathing changes in the last three months, and breathing symptom descriptors such as tightness in the chest, and wheezing) were investigated for possible confounding. The differences between the two measures of association were less than 10%, and therefore, possibility of confounding was small, if any.

Table 6.40 M-H χ^2 test for asthma

	Asthma	OR	Confidence Interval	
			95%	
Q10_Cgh_3mths	Crude	1.27	0.46	3.75
	M-H combined	1.25	0.50	3.17
Q19_BrChnges_3mths	Crude	2.21	0.64	9.74
	M-H combined	2.18	0.68	7.06
Q24c_BrChnges (<i>Tightness in chest</i>)	Crude	1.07	0.44	2.58
	M-H combined	1.04	0.45	2.41
Q25b_BrChnges (<i>Wheezing noise breathing in</i>)	Crude	1.55	0.65	3.76
	M-H combined	1.53	0.69	3.43
Q25c_BrChnges (<i>Wheezing noise breathing out</i>)	Crude	1.73	0.72	4.19
	M-H combined	1.72	0.77	3.85

Angina

Symptoms associated with angina include chest pain or discomfort, pain in the arms, neck, or shoulder accompanying chest pain, fatigue, breathlessness, and sweats (NHS Choices 2013). Similar symptom variables for LC (pain or discomfort in the chest, upper body or shoulders, ache or pain in the centre of chest or ribs, pain started in shoulder blade, breathing changes in the last three months, unexpected tiredness in the last three months, and hot or cold sweat in the last three months) were investigated for possible confounding. There were no observable differences between the crude and adjusted ORs (<10% difference) in the bivariate analyses for angina.

Table 6.41 M-H χ^2 test for angina

	Angina	OR	Confidence Interval	
			95%	
Q1_Pain_3mths	Crude	0.98	0.41	2.33
	M-H combined	0.96	0.44	2.12
Q3g_Pain (Ache or pain in centre of chest or ribs)	Crude	0.86	0.32	2.22
	M-H combined	0.84	0.35	2.03
Q3i_Pain (Pain started in shoulder blade)	Crude	0.56	0.15	1.79
	M-H combined	0.54	0.18	1.60
Q19_BrChnges_3mths	Crude	2.33	0.68	10.23
	M-H combined	2.33	0.73	7.46
Q29_Tired_3mths	Crude	0.93	0.40	2.16
	M-H combined	0.89	0.40	1.96
Q53_HCswat_3mths	Crude	1.16	0.49	2.74
	M-H combined	1.20	0.54	2.67

Arthritis

Bivariate analyses for arthritis showed possible confounding; a difference of 10% was observed between the unadjusted and adjusted ORs that might suggest confounding with pain variables (see Table 6.29). Furthermore, recent evidence has suggested an independent relationship between arthritis and lung cancer diagnosis in secondary care, and therefore the potential for confounding involving symptoms of arthritis and lung cancer (Walter et al. 2015). In light of this, arthritis (comorbidity) was added to the model with the other main effects to adjust for potential confounding.

Table 6.42 M-H χ^2 test for arthritis

	Arthritis	OR	Confidence Interval	
			95%	
Q1_Pain_3mths	Crude	0.88	0.37	2.16
	M-H combined	0.98	0.44	2.24
Q63_New_JPain_12mths	Crude	1.60	0.58	4.27
	M-H combined	1.72	0.69	4.30
Q64_New_JPain_12mths	Crude	2.94	1.02	8.44
	M-H combined	2.86	1.10	7.46

Allergy

Symptoms of seasonal allergies usually relate to breathing changes such as sneezing, and nasal congestion. Therefore, breathing changes variables (breathing changes/difficulties experienced in the last three months, and breathing problems that are only present or get worse at certain times of the year), were investigated for potential confounding. There were hardly any differences between the unadjusted and adjusted ORs to suggest confounding.

Table 6.43 M-H χ^2 test for allergy

	Allergy	OR	Confidence Interval	
			95%	
Q19_BrChnges_3mths	Crude	2.21	0.64	9.74
	M-H combined	2.21	0.69	7.11
Q23_BrChnges	Crude	1.99	0.64	6.03
<i>Have you experienced breathing problems that are only present or get worse at certain times of the year</i>	M-H combined	1.97	0.73	5.28

6.4.23 Multivariate analysis: COPD sub-population

Variable selection at $p < 0.05$

6.4.23.1 Model 1; Symptoms, adjusted for age ($p < 0.05$): COPD sub-population

Age, and age-squared were fitted to the model with the symptom variables that were statistically significant at $p < 0.05$ and comorbidity, arthritis (potential confounder) using forward stepwise regression. Symptom descriptor, unable to get enough air (Q24b_BrChnges) and arthritis (potential confounder) remained in the model Table 6.44.

Table 6.44 Model 1: Main effects model ($p < 0.05$), adjusted for age; using forward stepwise regression at $p(e) = 0.05$ and $p(r) = 0.10$

LCdiagnosis	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
AGE	1.198576	.2873255	0.76	0.450	.7492272	1.917422
AGEsq	.999011	.0017056	-0.58	0.562	.9956737	1.002359
Q69k_Arthritis_ever	.3881103	.1715376	-2.14	0.032	.1632074	.9229341
Q24b_BrChnges	3.171138	1.606268	2.28	0.023	1.175055	8.557995
_cons	.0001689	.0014248	-1.03	0.303	1.11e-11	2560.03

Akaike's information criterion and Bayesian information criterion

Model	Obs	ll (null)	ll (model)	df	AIC	BIC
.	1833	-55.2956	-49.07677	5	108.1535	135.7221

All discarded variables were checked against the model to assess the fit of the model, but none improved model fit.

6.4.23.2 Model 2; Symptoms, adjusted for age and risk variables ($p < 0.05$): COPD sub-population

According to the univariate analysis, none of the risk variables; pneumonia in the last five years, previous cancer, history of asbestos-related illnesses, ever smoked, and family history of LC, reached the statistical level of significance at $p < 0.05$ in this secondary population with COPD. Therefore, the same model to Model 1 ($p < 0.05$) would be obtained.

6.4.24 A set of diagnostic criteria for a sub-population with COPD ($p < 0.05$)

Using the results from the first modelling procedure with the stricter entry criteria for variable selection ($p < 0.05$), three variables were significantly associated with a diagnosis of lung cancer (age, a symptom variable, and a co-morbidity to adjust for confounding) as follows:

- Age was dichotomised at the level of 69 years such that those aged 69 years and above would score one. This cut-off was decided on the basis of the median age between the cases and non-cases, which were 67 and 71 years, respectively.
- Symptom descriptor, unable to get enough air (Q24b_BrChnges), contributed a score of one to the cumulative score.
- The absence of arthritis (Q69k_Arthritis_ever) would contribute to a score of one

A maximum score of three was obtainable. Table 6.45 shows the diagnostic performance of the set of criteria based on Model 1 developed to predict LC diagnosis in this secondary care population with COPD.

Table 6.45 Analysis of the diagnostic performance of this set of criteria

Criteria		LC cases	Non-LC cases	Sensitivity	Specificity	Youden index	1-Specificity	+ve Likelihood Ratio	-ve Likelihood Ratio
No. of criteria (max 4)	0.0	515	1094	1.00	0.0	0.0	1.00	1.0	-
	1.0	514	1021	0.998	0.067	0.07	0.933	1.07	0.03
	2.0	376	419	0.730	0.617	0.35	0.383	1.91	0.44
	3.0	149	65	0.289	0.941	0.23	0.059	4.87	0.76

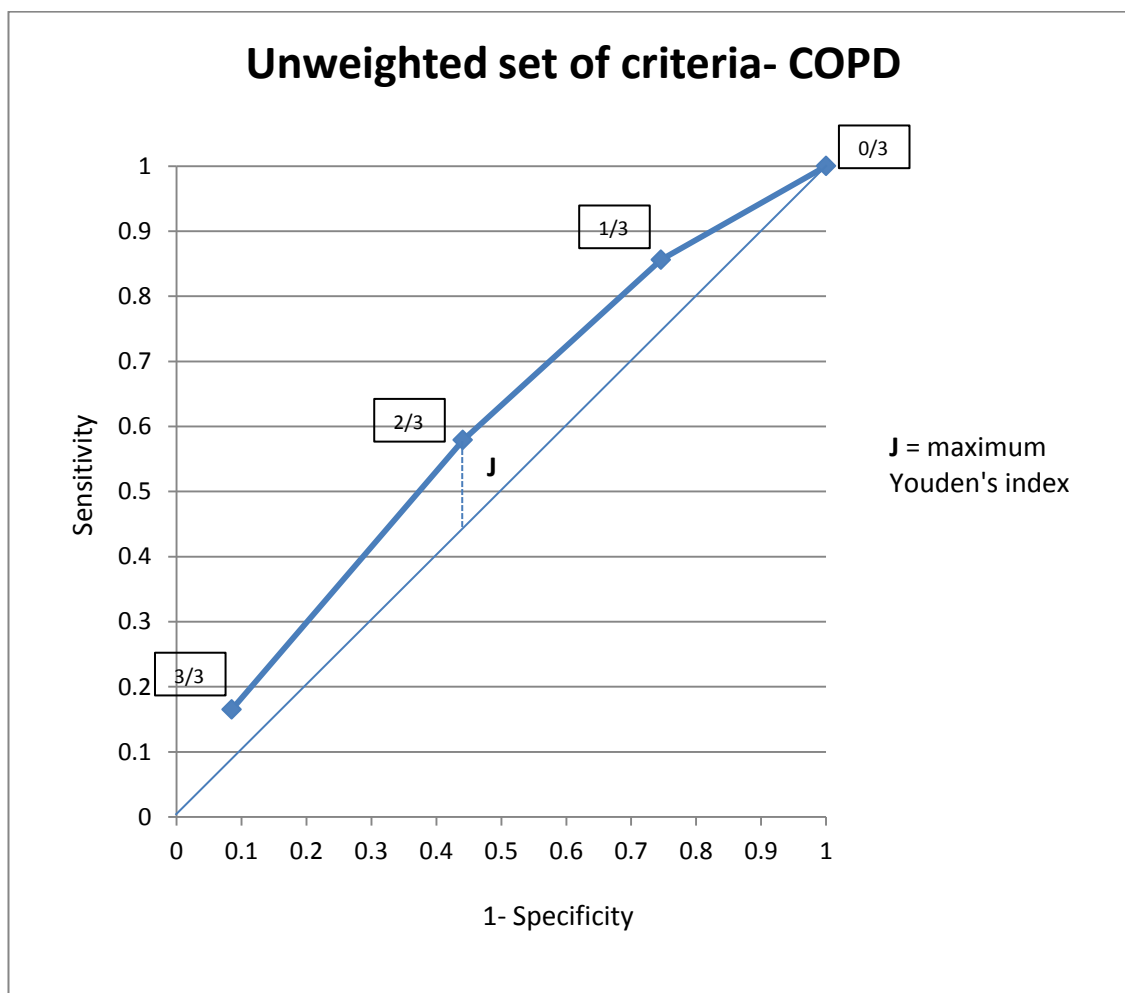


Figure 6.11 Receiver operating curve for the un-weighted set of diagnostic criteria for sub-group with COPD ($p < 0.05$)

Again using the Youden's index, the optimum cut-off level was determined and marked 'J' on the ROC curve for the un-weighted criteria (Youden 1950). As shown on the ROC curve in Figure 6.11, at the optimum cut-off level of two out of the three symptoms (2/3), 73% of the cases of lung cancer will be correctly identified (27% will be missed), and 61.7% of those investigated do not have lung cancer. Similarly, using the positive likelihood ratio, it informs us of the probability of a positive test in those with disease over the probability of positive test in those without disease (Attia 2003). This resulted in a positive likelihood ratio of 1.91 (see Figure 6.11). The discriminatory power (as measured by the AUC) for this criteria was 0.739, also considered to be slightly poor (Fan et al. 2006).

6.4.25 Weighted set of diagnostic criteria for a sub-population with COPD ($p < 0.05$)

In a similar method to the full population analysis, three variables were forced into a logistic regression model to obtain the log odds ratios, which were then rounded to obtain the weights shown in Table 6.46.

Table 6.46 Log odds ratios and derived weights from variables identified in the COPD model

Variables	Log Odds Ratio (β coefficient)	Rounded up weight for each criterion
Age (>69)	+1.280	+1
Q24b_BrChnges	+1.041	+1
Q69k_Arthritis_ever	+1.045	+1

The minimum obtainable score of the weighted set of criteria is 0, and the maximum is +3. The diagnostic performance of this weighted set of diagnostic criteria within the sub-population with COPD would be the same as the un-weighted criteria

Variable selection at $p < 0.15$

6.4.25.1 Model 1; Symptoms, adjusted for age ($p < 0.15$): COPD sub-population

Similarly, symptoms strongly associated with LC ($OR > 2$ or < 0.5 , or $p < 0.15$) were added into Model 1, adjusted for age, and potential confounder, arthritis. The p-values at which variables are considered for entry and exit from the model were also relaxed to 0.1 and 0.15, respectively (as explained in Methods).

The resultant Model 1 ($p < 0.15$) consisted of four symptom variables; breathing changes/difficulties experienced within the last three months (Q19_BrChnges_3mths), breathing changes/difficulties first indicated within the last three months (Q21_BrChnges), unable to get enough air (Q24b_BrChnges), and wheezing sensation when in a particular position (Q25d_BrChnges). The model presented in Table 6.47 was observed to have the optimal goodness of fit according to the AIC score (minimum AIC score). Discarded variables were checked against model 1.

Table 6.47 Model 1: Main effects model ($p < 0.15$), adjusted for age; using forward stepwise regression at $p(e) = 0.10$ and $p(r) = 0.15$

LCdiagnosis	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
AGE	1.248944	.3107821	0.89	0.372	.7668903	2.034007
AGEsq	.9987093	.0017509	-0.74	0.461	.9952834	1.002147
Q19_BrChnges_3mths	.1993067	.1761758	-1.82	0.068	.0352458	1.12703
Q69k_Arthritis_ever	.3883079	.1834741	-2.00	0.045	.1538107	.9803154
Q25d_BrChnges	2.128408	.8827125	1.82	0.069	.9441415	4.798139
Q21_BrChnges_R	4.766613	3.849852	1.93	0.053	.9788578	23.21133
Q24b_BrChnges	5.983991	4.392354	2.44	0.015	1.419692	25.22247
_cons	.0000562	.0004983	-1.10	0.269	1.61e-12	1960.94

Study 2

Akaike's information criterion and Bayesian information criterion

Model	Obs	ll (null)	ll (model)	df	AIC	BIC
.	1809	-54.57124	-44.60828	8	105.2166	149.2208

6.4.25.2 Model 2; Symptoms, adjusted for age and risk variables ($p < 0.15$): COPD sub-population

Risk variables, asbestos-related illness, and ever smoked, were entered into model 1 (COPD sub-group). Forward stepwise regression was also carried out at the relaxed entry criteria, $p = 0.10$ and exit criteria, $p = 0.15$. None of the risk variables entered the model. Therefore, the same model to Model 1 ($p < 0.15$) was obtained for Model 2 ($p < 0.15$) in this referred population with COPD.

6.4.26 A set of diagnostic criteria for a sub-population with COPD ($p < 0.15$)

In a similar way to the full population data, a set of diagnostic criteria was created to distinguish those with lung cancer and non-lung cancer for a population with COPD. Model 1 ($p < 0.15$) was used because none of the risk variables entered Model 2. For each case and non-case, the number of criteria that were satisfied was calculated in the following manner:

- Age was dichotomised at the level of 69 years such that those aged 69 years and above would score one. This cut-off was decided on the basis of the median age between the cases and non-cases, which were 67 and 71 years, respectively.
- Breathing changes/difficulties experienced within the last three months (Q19_BrChnges_3mths), which occurred more frequently in those without lung cancer, was recoded so that the absence of the breathing variant resulted in a score of one.
- Any one of the three remaining symptoms present breathing changes/difficulties first indicated within the last three months (Q21_BrChnges), unable to get enough air (Q24b_BrChnges), and

wheezing sensation when in a particular position (Q25d_BrChnges), contributed a score of one to the cumulative score.

- Those who do not have a history of arthritis (Q69k_Arthritis_ever) would contribute to a score of one.

The maximum score was six. Therefore, a score out of six was produced to indicate the number of criteria that were satisfied. The sensitivity, specificity, and likelihood odds ratios was calculated for all levels of cut-offs (Table 6.48).

Table 6.48 Analysis of the performance of the un-weighted set of diagnostic criteria for sub-group with COPD

Criteria		LC cases	Non-LC cases	Sensitivity	Specificity	Youden index	1-Specificity	+ve Likelihood Ratio	-ve Likelihood Ratio
1.35	0.0	531	1299	1.00	0.0	0.0	1.00	1.0	-
	1.0	531	1275	1.00	0.019	0.02	0.982	1.02	0.0
	2.0	502	1089	0.945	0.162	0.11	0.838	1.13	0.34
	3.0	360	468	0.678	0.640	0.32	0.360	1.88	0.50
	4.0	177	125	0.333	0.904	0.24	0.096	3.46	0.74
	5.0	102	33	0.192	0.975	0.17	0.025	7.56	0.83
	6.0	1	23	0.002	0.982	0.02	0.018	0.11	1.02

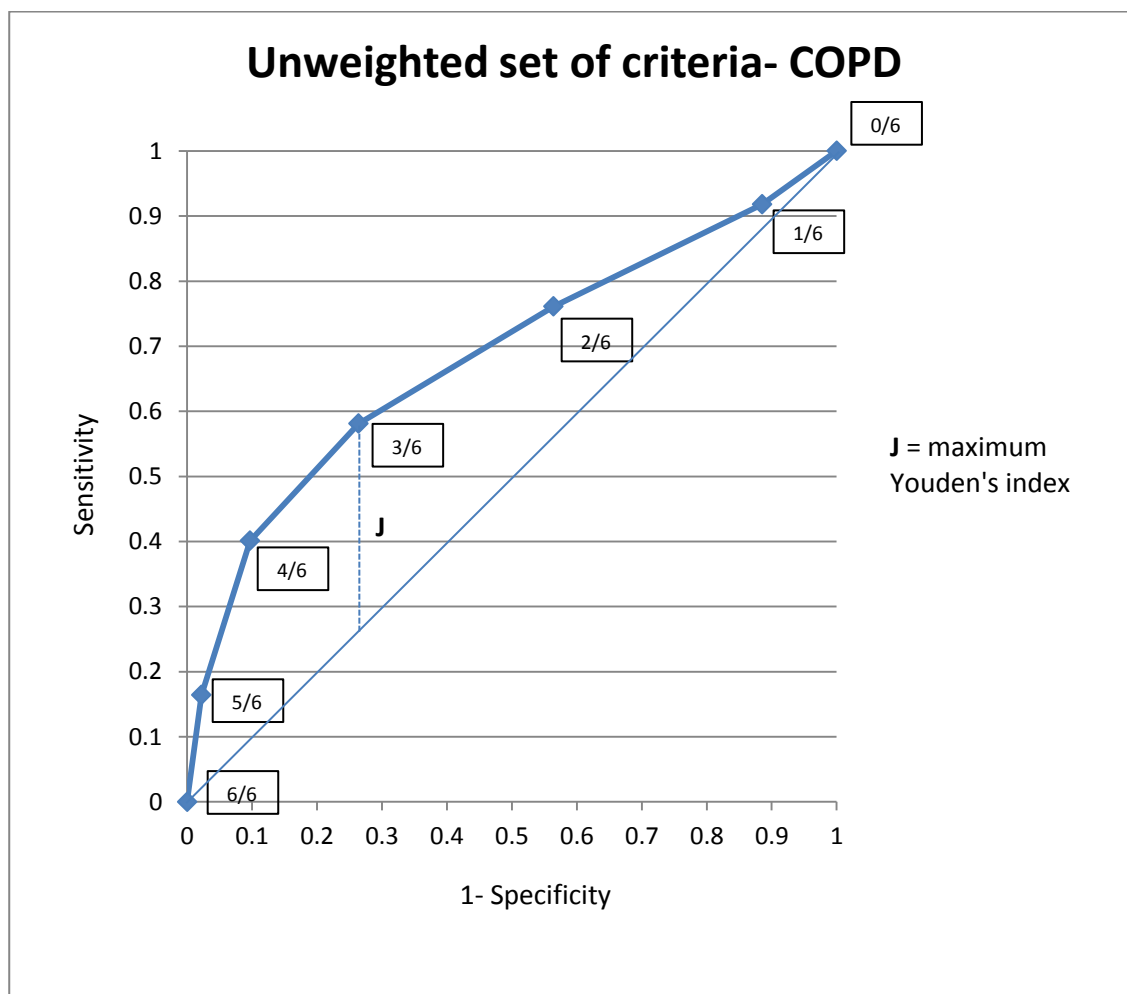


Figure 6.12 Receiver operating curve for the unweighted set of diagnostic criteria for sub-group with COPD ($p < 0.15$)

Figure 6.12 shows the receiver operating curve for the un-weighted set of criteria for this sub-group analysis. As previously explained, the interpretation of the cut-off point in a ROC curve is dependent on situation and the population it is intended for; which allow for multiple cut-off options. Again, the Youden's index was used to determine the optimum level for cut-off on the ROC curve, as it is considered a good summary measure of the ROC curve. This is indicated by the cut-off value with the shortest Euclidean distance marked 'J' on the ROC curve (Youden 1950).

Where sensitivity is equally important as specificity, the optimum cut-off for this criteria is at 3/6. This suggests that if an individual with COPD referred to the lung-shadow clinic, presented with any three out of the six criteria were

investigated, 67.8% of the cases will be identified with lung cancer and 32.2% of those investigated would have been without the disease. The positive likelihood ratio at this cut-off level is 1.88. The criteria produced a discriminatory power (determined by AUC) of 0.790.

6.4.27 Weighted set of diagnostic criteria for a sub-population with COPD ($p < 0.15$)

Individual weights derived for each criterion were applied to this set of diagnostic criteria for the sub-group with COPD (see Table 6.49). This resulted in a cumulative score of minimum 0 and maximum of +9 obtainable. The performance of this weighted set of criteria was evaluated (Table 6.50) and presented in a ROC curve (Figure 6.14).

Table 6.49 Log odds ratios and derived weights from variables identified in the model

Variables	Log Odds Ratio (β coefficient)	Rounded up weight for each criterion
Age (>69)	+1.39	+1
Q21_BrChnges	+2.12	+2
Q24b_BrChnges	+2.12	+2
Q25d_BrChnges	+1.02	+1
Q19_BrChnges_3mths	+2.22	+2
Q69k_Arthritis_ever	+1.07	+1

Table 6.50 Analysis of the diagnostic performance of the weighted set of criteria in the sub-group analysis

Criteria		LC cases	Non-LC cases	Sensitivity	Specificity	Youden index	1-Specificity	+ve Likelihood Ratio	-ve Likelihood Ratio
No. of criteria (max 5)	0.0	531	1299	1.00	0.0	0.0	1.0	1.0	-
	1.0	531	1275	1.00	0.018	0.02	0.982	1.02	0.0
	2.0	503	1186	0.947	0.087	0.03	0.913	1.04	0.61
	3.0	501	963	0.944	0.259	0.20	0.741	1.27	0.22
	4.0	352	368	0.663	0.717	0.38	0.283	2.34	0.47
	5.0	198	117	0.373	0.910	0.28	0.09	4.14	0.69
	6.0	126	45	0.237	0.965	0.20	0.035	6.85	0.79
	7.0	39	29	0.073	0.978	0.05	0.022	3.29	0.95
	8.0	1	2	0.002	0.999	0.0	0.001	1.22	1.0
	9.0	0	2	0.0	0.999	0.0	0.001	0.0	1.0

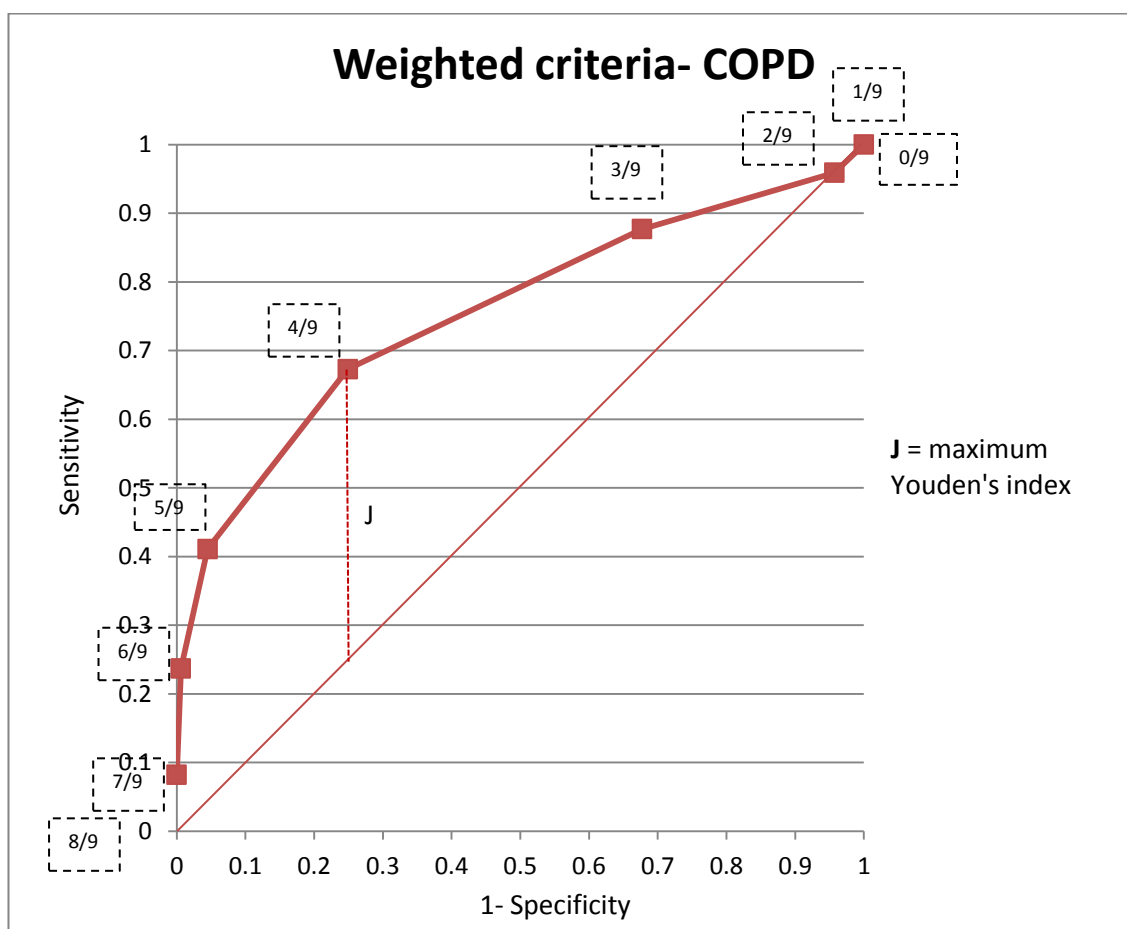


Figure 6.13: Receiver operating curve of the weighted set of criteria for the COPD sub-group

As shown in the ROC curve in Figure 6.13, the optimum cut-off point on the weighted criteria is 4/9 using the Youden's index (Youden 1950). According to the data in Table 6.50, the positive likelihood ratio at this level of cut-off is 2.34 with a sensitivity and specificity of 66.3% and 71.7%, respectively. 28.3% will be false positives.

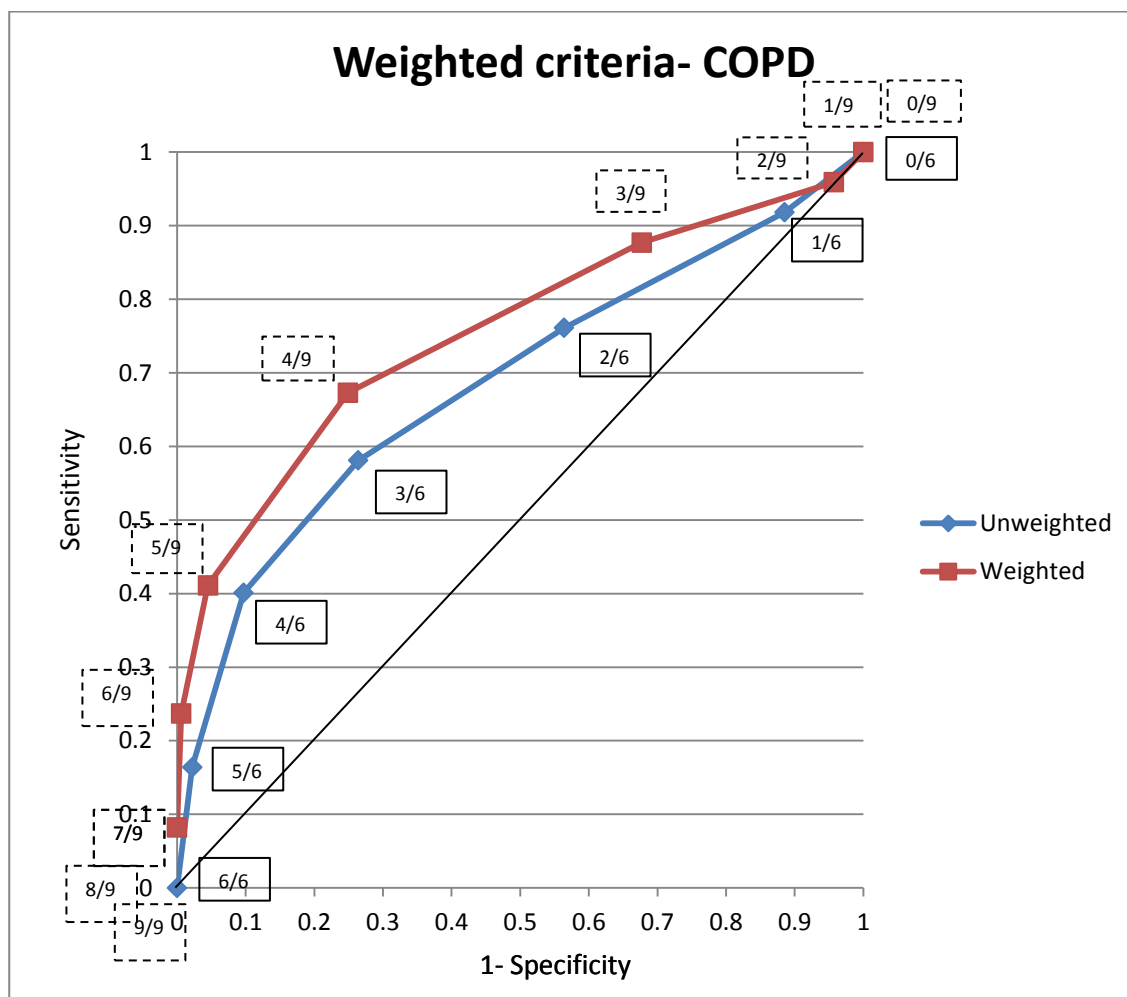


Figure 6.14 Receiver operating curves of the weighted and un-weighted set of criteria for the COPD sub-group

A comparison of the two ROC curves shown in Figure 6.14 demonstrates that the weighted criteria performed slightly better than the un-weighted criteria. The positive likelihood ratio for the weighted criteria (4/9) was slightly higher, 2.34 compared to the 1.88 in the un-weighted criteria (3/6). Overall, both versions of the symptom criteria did not demonstrate particularly strong diagnostic performance in distinguishing lung cancer in this referred population with COPD.

At the higher cut-off levels (6/9), the weighted criteria produced a maximum positive likelihood ratio of 6.85 but sensitivity is low. Although this likelihood ratio is above the recommended positive likelihood ratio (>5) to produce strong diagnostic evidence, the calculations were based on a very small sample

size using a relaxed model criteria. Therefore, these findings are highly exploratory, and should be interpreted with caution.

The diagnostic performance (AUC, optimum cut-off and +LR) at the optimum cut-off for each version of the model (relaxed and strict criteria, COPD and whole population) is presented in Table 6.51 for ease. The models cannot be directly compared against one another without performing statistical analyses to confirm, but there are indications of possible benefits to working with the COPD sub-group.

Table 6.51 A table summary of the performance of the criteria for each version of the model

	Strict criteria				Relaxed criteria			
	Full population		COPD		Full population		COPD	
	UW	W	UW	W	UW	W	UW	W
Optimum cut-off	2/4	-	2/3	-	4/7	-	3/6	4/9
+ve LR	1.49	-	1.91	-	1.76	-	1.88	2.34
AUC	0.66	-	0.74	-	0.78	-	0.79	

UW= Un-weighted criteria

W= Weighted criteria

6.4.28 Results of diagnostic models: complete case vs. imputed

Variations in effect sizes in both directions; increase and decrease of ORs, were observed, as shown in Table 6.52.

Table 6.52 Univariate analyses of symptoms and lung cancer as outcome for complete case and imputed data.

Variable	Complete case		Imputed	
	Odds Ratio	p-value	Odds Ratio	p-value
Q1_Pain_3mths	0.8	0.41	0.8	0.40
Q3a_Pain	1.1	0.69	1.1	0.68
Q3b_Pain	1.0	0.97	1.0	0.99
Q3c_Pain	0.8	0.59	0.8	0.54
Q3d_Pain	1.1	0.73	1.1	0.72
Q3e_Pain	0.7	0.23	0.7	0.25
Q3f_Pain	1.1	0.67	1.1	0.70
Q3g_Pain	0.6	0.11	0.6	0.12
Q3h_Pain	0.9	0.57	0.9	0.60
Q3i_Pain	1.0	0.97	1.1	0.83
Q3j_Pain	1.1	0.84	1.2	0.70
Q6_Pain	1.3	0.36	1.3	0.35
Q7_Pain	0.7	0.57	0.7	0.56
Q8_Pain	1.0	0.51	1.0	0.51
Q9_Pain	1.0	0.59	1.0	0.59
Q10_Cgh_3mths	1.4	0.26	1.4	0.27
Q12_Cgh_R	1.9	0.02**	1.8	0.03**
Q13a_Cgh	1.0	0.96	1.0	0.92
Q13b_Cgh	1.5	0.15	1.5	0.14
Q13c_Cgh	1.3	0.37	1.3	0.32
Q13d_Cgh	1.4	0.20	1.5	0.16
Q13e_Cgh	0.9	0.74	0.9	0.73
Q13f_Cgh	0.9	0.59	0.8	0.51
Q13g_Cgh	0.8	0.44	0.8	0.44

Q13h_Cgh	1.1	0.67	1.1	0.66
Q13i_Cgh	1.0	0.96	1.0	0.99
Q13j_Cgh	1.3	0.38	1.2	0.44
Q13k_Cgh	0.6	0.08*	0.6	0.10
Q13l_Cgh	1.1	0.62	1.1	0.70
Q14a_Cgh	1.1	0.71	1.1	0.75
Q14b_Cgh	1.0	0.91	1.1	0.77
Q14c_Cgh	1.6	0.16	1.6	0.14
Q15_Cgh	0.6	0.18	0.7	0.23
Q16_Cgh	0.2	0.58	1.2	0.59
Q18_Cgh	1.0	0.47	1.0	0.47
Q19_BrChnges_3mths	1.2	0.56	1.2	0.53
Q21_BrChnges	1.95	0.02**	1.9	0.02**
Q22a_BrChnges	0.8	0.40	0.8	0.48
Q22b_BrChnges	0.9	0.68	0.9	0.78
Q22c_BrChnges	1.0	1.0	1.0	0.91
Q22d_BrChnges	1.1	0.82	1.1	0.71
Q22e_BrChnges	1.0	0.91	1.0	0.88
Q22f_BrChnges	0.9	0.71	0.9	0.72
Q23_BrChnges	1.7	0.14	1.7	0.13
Q24a_BrChnges	1.1	0.80	1.1	0.75
Q24b_BrChnges	1.4	0.23	1.4	0.21
Q24c_BrChnges	0.9	0.76	0.9	0.75
Q24d_BrChnges	0.9	0.77	0.9	0.75
Q24e_BrChnges	1.1	0.65	1.1	0.78
Q24f_BrChnges	0.8	0.57	0.9	0.64
Q25a_BrChnges	0.8	0.62	0.9	0.69
Q25b_BrChnges	1.0	0.96	1.0	0.99

Study 2

Q25c_BrChnges	1.1	0.77	1.1	0.78
Q25d_BrChnges	1.2	0.51	1.3	0.47
Q26_BrChnges	1.6	0.1	1.5	0.12
Q27_BrChnges	0.5	0.14	0.5	0.14
Q28_BrChnges	1.0	0.44	1.0	0.44
Q29_Tired_3mths	1.0	0.95	1.0	0.96
Q31_Tired	0.9	0.78	0.9	0.80
Q32_Tired	0.9	0.82	0.9	0.78
Q33_Tired	1.0	0.89	1.0	0.92
Q34_Tired	0.9	0.78	0.9	0.77
Q35_Tired	1.5	0.16	1.5	0.16
Q36_Tired	0.8	0.49	0.8	0.47
Q37_Tired	1.0	0.79	1.0	0.79
Q38_CghBlood_3mths	1.1	0.78	1.1	0.82
Q40_CghBlood	1.3	0.35	1.3	0.39
Q41_CghBlood	1.3	0.51	1.3	0.51
Q42_CghBlood	1.1	0.81	1.1	0.84
Q43_ChInfectn	1.0	0.98	1.0	0.90
Q44_ChInfectn	0.8	0.55	0.8	0.50
Q45_ChInfectn	0.9	0.61	0.9	0.60
Q46_ChInfectn	1.7	0.07*	1.7	0.06*
Q47_ChInfectn	1.0	0.74	1.0	0.84
Q48_ChInfectn	0.8	0.41	0.7	0.36
Q49_Weight	1.1	0.81	1.1	0.87
Q50_Weight	0.8	0.57	0.9	0.66
Q51_Weight	1.1	0.69	1.1	0.69
Q52_Weight	0.6	0.09*	0.6	0.09*
Q53_HCswat_3mths	1.1	0.81	1.0	0.87
Q55_HCswat	1.2	0.67	1.2	0.70

Q56_HCsweat	0.9	0.87		1.0	0.85
Q57_HCsweat	0.6	0.13		0.7	0.16
Q58_HCsweat	0.7	0.23		0.7	0.23
Q59_EatChnges	0.5	0.14		0.5	0.14
Q60_EatChnges	1.3	0.37		1.3	0.42
Q61_EatChnges	0.6	0.28		0.6	0.22
Q62_EatChnges	1.9	0.05**		1.7	0.12
Q63_New_JPain_3mths	1.1	0.64		1.2	0.57
Q65_Voice	0.7	0.32		0.7	0.31
Q66_Voice	0.8	0.65		0.8	0.71
Risk variables					
Q69a_Pneumo_last5yrs	0.52	0.197		0.7	0.325
Q69e_COPD_ever	1.58	0.158		1.4	0.245
Q69h_Cancer_ever	2.98	0.001*		2.4	0.003*
Q69j_Asbes_ever	1.08	0.869		0.9	0.849
Q73_Smoke	6.15	0.001*		6.3	0.001*
Q71d_FamHx	1.27	0.469		1.3	0.469

Slight differences can be observed in the effect estimates and p-values between the complete case (original dataset with missing values) and the imputed univariate analysis. Although p-values were better (smaller p-value; more statistically significant) in the complete case analysis compared to the imputed analysis, higher effect sizes could be observed in the imputed dataset. This could affect the modelling process, where variable selection relies on the significance level. Models created from the complete case and the imputed dataset were compared in the following sections.

Variable selection at $p < 0.05$

6.4.28.1 Model 1; Symptoms, adjusted for age ($p < 0.05$): Complete case analysis on full population data

Using similar forward stepwise regression, variables were added and removed at $p < 0.05$ and $p < 0.10$, respectively, for the complete case model 1 and 2. Results are presented in Table 6.53 and Table 6.54. At the more stringent significance level ($p < 0.05$), variables that enter Model 1 were similar for both observed and imputed data (see Table 6.53).

Table 6.53 Model 1 including variables at $p < 0.05$: Complete case

LCdiagnosis	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
AGE	1.372605	.1986538	2.19	0.029	1.033601	1.822795
AGEsq	.9978669	.0010439	-2.04	0.041	.9958229	.999915
Q12_Cgh_R	2.280336	.7145777	2.63	0.009	1.233849	4.214399
Q21_BrChnges_R	1.755822	.5506557	1.79	0.073	.9495748	3.246622
_cons	1.95e-06	9.68e-06	-2.65	0.008	1.19e-10	.0322151

Model	Obs	ll (null)	ll (model)	df	AIC	BIC
.	288	-150.0094	-140.1782	5	290.3564	308.6712

6.4.28.2 Model 2; Symptoms, adjusted for age and risk variables ($p < 0.05$): Complete case analysis on full population data

In Model 2, both risk variables; previous cancer (Q69h_Cancer), and ever smoked (Q73_Smoke), stayed in the model for the complete case (see Table 0.51). Previous cancer was removed from the **imputed** Model 1 ($p < 0.05$).

Table 6.54 Model 2 including variables at $p < 0.05$: Complete case

LcDiagnosis	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
AGE	1.282123	.1780915	1.79	0.074	.9765501	1.683312
AGEsq	.9983134	.0010032	-1.68	0.093	.9963491	1.000282
Q12_Cgh_R	1.879609	.6503229	1.82	0.068	.9540288	3.703167
Q21_BrChnges_R	1.940496	.6645608	1.94	0.053	.9917498	3.79685
Q69h_Cancer_ever	2.778377	1.159498	2.45	0.014	1.226196	6.295388
Q73_Smoke	9.179546	6.852035	2.97	0.003	2.125423	39.64579
_cons	2.82e-06	.0000136	-2.65	0.008	2.24e-10	.035451

Akaike's information criterion and Bayesian information criterion

Model	Obs	ll (null)	ll (model)	df	AIC	BIC
.	272	-138.2983	-119.3598	7	252.7195	277.9601

6.4.28.3 Model 1; Symptoms, adjusted for age ($p < 0.15$): Complete case analysis on full population data

Table 6.55 shows the larger model using the more relaxed significance level ($p < 0.15$), or $OR > 2.0 < 0.5$, generated in the complete, un-imputed dataset. All the discarded variables were checked against the model to see if they might improve fit. The model looked slightly different to the imputed model. The model included pain in the side of the chest or ribs (Q3g_Pain), worsening breathlessness in the last three months (Q26_BrChnges), coughing (for more than three weeks) that were first indicated within the last three months (Q12_Cgh), a hard or harsh cough without phlegm (Q13k_Cgh), and a tickly cough (Q13b_Cgh).

Study 2

Table 6.55 Model 1 including variables p<0.15: Complete case analysis

LCdiagnosis	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
AGE	1.228962	.1937975	1.31	0.191	.9022162	1.674041
AGEsq	.9985678	.0011478	-1.25	0.212	.9963208	1.00082
Q12_Cgh_R	2.719028	1.021203	2.66	0.008	1.302335	5.676817
Q3g_Pain	.3316024	.1441139	-2.54	0.011	.1414773	.777228
Q26_BrChnges	2.421181	.894144	2.39	0.017	1.174026	4.993173
Q13k_Cgh_R	.3273647	.1386304	-2.64	0.008	.1427481	.7507466
Q13b_Cgh_R	2.526963	.9412453	2.49	0.013	1.217706	5.243909
_cons	.000119	.0006359	-1.69	0.091	3.38e-09	4.196968

Akaike's information criterion and Bayesian information criterion

Model	Obs	ll (null)	ll (model)	df	AIC	BIC
.	233	-121.1548	-103.7174	8	223.4348	251.0431

6.4.28.4 Model 2; Symptoms, adjusted for age and risk variables (p<0.15): Complete case analysis on full population data

Model 2 (the model in which risk variables were also entered) appeared to have better fit in the complete case analysis when compared to the rest of the models generated using complete case (minimum AIC), at p<0.15. Both risk variables; ever smoked (Q73_Smoke) and previous cancer (Q69h_Cancer_ever), entered the model. The model using the complete case analysis will be compared against the imputed model to discuss the discrepancies between the two, and suggested reasons for their departures.

Table 6.56 Model 2 including variables p<0.15: Complete case analysis

LCdiagnosis	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
AGE	1.183185	.1803617	1.10	0.270	.8776018	1.595172
AGEsq	.9988285	.0011109	-1.05	0.292	.9966535	1.001008
Q73_Smoke	7.055642	5.406728	2.55	0.011	1.571311	31.68188
Q12_Cgh_R	2.143884	.8707496	1.88	0.060	.9671226	4.75249
Q3g_Pain	.3078439	.1439559	-2.52	0.012	.123108	.7697945
Q13b_Cgh_R	2.559087	1.020329	2.36	0.018	1.171397	5.590697
Q13k_Cgh_R	.382832	.1701007	-2.16	0.031	.1602509	.9145678
Q26_BrChnges	2.420866	.9685438	2.21	0.027	1.105149	5.30299
Q69h_Cancer_ever	2.428456	1.22346	1.76	0.078	.9046767	6.518789
_cons	.0000714	.0003729	-1.83	0.068	2.55e-09	2.000702

Akaike's information criterion and Bayesian information criterion

Model	Obs	ll(null)	ll(model)	df	AIC	BIC
.	222	-113.2717	-92.4477	10	204.8954	238.9222

6.4.28.5 Comparison between complete case model and imputed model

Following the recommendations of Carpenter and Kenward (2008), a comparison of the imputed (partially observed) data analysis and the complete case (those participants or units with no missing data) analysis showed no differences in the model using the traditional entry criteria ($p < 0.05$). However, when a relaxed criteria for variable selection was used ($p < 0.15$), some disparities were observed between the imputed analysis and the complete case analysis, see Table 6.57.

Table 6.57 Comparison of Model 1 ($p < 0.15$) between the complete case and imputed data

	Complete case		Imputed
Variables that remained in the model	Q3g_Pain		Q12_Cgh
	Q12_Cgh		Q13k_Cgh
	Q13b_Cgh		Q21_BrChnges
	Q13k_Cgh		Q46_ChInfectn
	Q26_BrChnges		Q59_EatChnges

Variables; pain in the side of the chest or ribs (Q3g_Pain), a tickly cough (Q13b_Cgh), and worsening breathlessness in the last three months (Q26_BrChnges), remained in the complete case model but was dropped in the imputed model, which suggests that relative to the other symptom variables in the model, these variable might not be as strongly associated to lung cancer in the imputed dataset. However, these were variables dropped at the more conservative significance level ($p < 0.05$), which might not be a relevant variable in the fully powered study. It should be noted that the findings of the relaxed criteria were not be as robust considering the differences between the two models. The estimated regression parameters for these variables did not vary much after the imputation but rather, the variables in the imputed model had smaller standard errors compared to the complete case analysis.

Further attempts were made to explore the reason for this disparity by looking at the descriptive characteristics of the missing data in symptom variables; pain in the side of the chest or ribs (Q3g_Pain), a tickly cough (Q13b_Cgh), and worsening breathlessness in the last three months (Q26_BrChnges), by socio-demographic variables, clinical outcome variables, and other symptom covariates. Looking at Table 6.58, there appears to be more missing data in variables; Q3g_Pain, Q13b_Cgh, and Q26_BrChnges, in those who had never smoked or were ex-smokers. No other remarkable findings were observed. It is noted that only 'ever-smoked' variable (Q73_Smoke) was analysed in the

multivariate models, and that particular variable only had 3.3% missingness in the whole data.

Table 6.58 Descriptive characteristics of the three miss_variables;
Q21_BrChnges, Q46_ChInfectn, and Q52_Weight, by socio-demographic variable (age), outcome variable (lung cancer diagnosis), clinical covariates (COPD, smoking, family history)

	Q3g_Pain		Q13b_Cgh		Q26_BrChnges	
	Missing (%)	Non-missing (%)	Missing (%)	Non-missing (%)	Missing (%)	Non-missing (%)
AGE>68 *	8 (47)	182 (53)	17 (61)	173 (52)	15 (75)	175 (52)
AGE<68	9 (53)	160 (47)	11 (39)	158 (48)	5 (25)	164 (48)
LC	1 (6)	76 (22)	4 (14)	73 (22)	7 (35)	70 (21)
NO LC	16 (94)	266 (78)	24 (86)	258 (78)	13 (65)	269 (79)
COPD	3 (18)	121 (35)	9 (32)	115 (35)	5 (33)	119 (35)
NO COPD	14 (82)	221 (65)	19 (68)	216 (65)	15 (67)	220 (65)
SMOKER *	9 (56)	262 (79)	18 (75)	253 (78)	10 (63)	261 (79)
NON-SMOKER	7 (44)	69 (21)	6 (25)	70 (22)	6 (37)	70 (21)
CURRENT SMOKER	2 (22)	88(33)	3 (17)	87 (34)	0 (0)	90 (34)
EX-SMOKER	7 (78)	177 (67)	15 (83)	169 (66)	10 (100)	174 (66)
FAMILY HISTORY LC	2 (15)	64 (21)	4 (21)	62 (21)	3 (23)	66 (21)
NO FAMILY HISTORY	11 (85)	242 (79)	15 (79)	238 (79)	10 (77)	243 (79)

Mean age (full population): 68 years

*variables that entered the model (selected)

As the complete case analysis is only valid if the missing data mechanism is MCAR, it is also possible that the missing data mechanism was MAR, and the

results of complete case analyses had some degree of bias. Therefore, the variations observed might be a result of correcting for bias.

Variations were to be expected between the imputed and complete case analysis particularly in the model with the larger number of independent variables considering that the study was under-powered. Differences in the variables selection could be due to the slightly improved p-value in the complete case analysis, which might be an effect of bias. It is possible that MI was correcting for bias, which should be investigated in a larger, fully-powered study.

6.5 Discussion

The previous IPCARD Feasibility Study identified patient-elicited symptoms that predicted chest X-rays suspicious for LC in a GP-referred chest X-rays population. In the chest X-ray population, response rates (>70%) and data completion were high (>80%). Feasibility of using the questionnaire in this population needed to be established in the current population referred to secondary care with high rates of chronic respiratory disease, which was a higher risk population than the chest X-ray population. The recruitment rate was relatively good at 67.9%, and better than a recent prospective secondary care study in the UK, which reported recruitment rates of 19.5% (Walter et al. 2015). Kubik et al. (2001) achieved response rates of 83.3% in their secondary care study in Czech Republic.

6.5.1 Missing data

Unlike MI methods, complete case analysis would only be valid if the missing data mechanism was MCAR, otherwise complete case analysis would produce biased results. For MI to be appropriate data would need to be MAR.

MAR missing data could only be argued on the basis of departures of the missing data from MCAR assumptions; such departures were evident in the associations between some of the observed covariates and the missingness in the generic symptom variables, and variables with missingness > 10%.

Although some observations may appear consistent with the data being MAR, MNAR cannot be ruled out (Molenberghs et al. 2004).

One limitation of the missing data analysis is that interactions were not included in the imputation model due to the small sample size. The lack of power could mean that clinically important interactions were missed. However, Wood et al. (2008) stated that accounting for all possible interactions would make the imputation model impractically large, and add very little information to the model (Wood et al. 2008) Therefore, selection of non-linear and interaction terms would present further difficulties.

The use of dummy variables to account for missing data in the risk variables was necessary to cause the imputation model to converge. However, there is the potential for the use of dummy variables to introduce bias into the model. The implications of the use of dummy variables in the risk variables were also investigated. Complete case analyses of the original dataset with missing data were comparable to the imputed analyses, at least at the more conservative $p < 0.05$, which could suggest reasonability in the imputed data. The purpose of the diagnostic checks performed throughout the process was to flag up where the imputation modelling may not be appropriate, and check the robustness of the data. Any discrepancy observed could be further explored in the larger study. MI is not the only principled method for handling missing values, nor is it necessarily more suitable for any given problem. However, in real datasets, where missing data are a nuisance rather than the primary focus, a convenient, approximate solution with reasonable properties is preferable. Within the limitations of a PhD project, sensitivity analyses were not feasible, but this could be recommended for future research.

6.5.2 Discussion of main findings

At the traditional threshold for variable selection ($p < 0.05$), there were few symptoms that predict lung cancer even in this secondary care population referred to chest-clinics. This confirmed the findings of previous studies and systematic review (see Chapter Three). Coughing that was first indicated within

Study 2

the last three months (OR=1.74; CI=0.91-3.32), and breathing changes that was first indicated within the last three months (OR=1.86; CI=0.97-3.57) were symptoms associated with lung cancer in this study population ($p<0.05$).

The risk model (Model 2), with previous smoking variable added had a better fit compared to Model 1. However, the analysis was limited to only previous smoking data. There were many missing data in the 'current smoking' variable and variables needed to calculate pack-years, and therefore, pack-years could not be used in the modelling process.

Use of a relaxed criteria for variable selection ($p<0.15$): full population

For exploratory purposes, any symptoms that were significant at p -value <0.15 , or ORs >2.0 or <0.5 were entered into a model. This step can be helpful in identifying clinically important variables that in the univariate analyses, are not significantly related to the outcome but make an important contribution in the presence of other variables (Sterne and Kirkwood 2003). Furthermore, the current study was under-powered; therefore, in an adequately-powered study, these symptoms could significantly contribute to the model. These findings were interpreted with caution, with consideration of the confidence intervals.

The use of relaxed criteria for variable selection ($p<0.15$) generated larger models. The presence of breathing changes/difficulties that was first indicated within the last three months (Q21_BrChnges), coughing (for more than three weeks) that were first indicated within the last three months (Q12_Cgh) (from the more stringent model), and noticeably more chest infections in the last 12 months than the years before (Q46_ChInfectn) (OR=1.95; CI=1.00-3.80) were significantly associated with the diagnosis of lung cancer in the full population attending lung-shadow clinic. The absence of a hard or harsh cough without phlegm (Q13k_Cgh) (OR=0.44; CI=0.22-0.88), and the absence of weight gain (Q52_Weight) (OR=0.48; CI=0.22-1.02) were also included in the model ($p<0.15$).

Under the relaxed criteria (full population models), the absence of weight gain added to symptom Model 1 ($p<0.15$) was however, dropped in Model 2 in favour of the absence of increased appetite with the addition of previous smoking (risk variable). It is possible that the absence of weight gain is

correlated to smoking, or the absence of increased appetite explains the variability in the model more significantly than weight gain in smokers. However, the potential for interactions between these symptoms and smoking were not tested in any of the models due to the small sample size and exploratory nature of this analysis. The directionality of the effect of the symptoms, weight gain (OR= 0.48) and eating changes (increased appetite) (OR=0.31), were as expected (OR<1). Weight loss qualifies for urgent referral for chest X-ray under NICE guidelines (NICE 2014), with reported positive predictive value of 1.1% and 6.1% (see Systematic review, Chapter Three). Appetite loss is also frequently reported to be associated with lung cancer (Ades et al. 2014; Hippenley-Cox 2011; Hamilton et al. 2005). One primary care study reported higher appetite loss in the more advanced stage of diagnosis (Ades et al. 2014).

It is highlighted that these variables; noticeably more chest infections in the last 12 months than the years before (Q46_ChInfectn), the absence of a hard or harsh cough without phlegm (Q13k_Cgh), and the absence of increased appetite (Q59_EatChnges), were not significantly associated in the univariate analysis at the usual significance level of 5%. Therefore, it is important to avoid over interpreting the results of the quantitative analysis as they were highly exploratory, and were based on a small sample size.

6.5.3 Sub-population COPD analysis

Based on the evidence in the bivariate analyses with the symptom variables in the COPD sub-group, arthritis (comorbidity) was added into the model to explore the effect of potential confounding in this COPD sub-population. Using the stricter criteria for variable selection ($p < 0.15$), a variant of breathing change; unable to get enough air (OR=3.17, CI=1.18-8.56), significantly entered Model 1. A risk model (Model 2) was not created because none of the risk variables achieved the statistical threshold at $p < 0.05$ in the univariate analysis. The resultant discriminatory power of the set of diagnostic, determined by the AUC statistic, was 0.739.

Use of a relaxed criteria for variable selection ($p < 0.15$): COPD

When the criteria was relaxed to $p < 0.15$, a larger model was obtained. Four symptoms of breathing changes/difficulties; generic breathing changes/difficulties experienced within the last three months ($OR = 0.20$; $CI = 0.04-1.13$), breathing changes/difficulties first indicated within the last three months ($OR = 5.38$, $CI = 1.16-24.91$), unable to get enough air ($OR = 5.98$; $CI = 1.42-25.22$), and wheezing sensation when in a particular position ($OR = 2.13$; $CI = 0.94-4.80$) were identified. As a rule of thumb, there should really be at least 10 variables per case in the logistic regression model, which require large sample sizes. Our reasoning for the relaxed model is mainly for exploratory purposes to investigate the potential use of symptoms to discriminate between those with and without lung cancer in a homogeneous COPD group. The standard errors in the COPD-specific models were big with wide confidence intervals, which is likely to be a reflection of the small sample size. Therefore, it is recommended that findings in this sub-population should be interpreted with caution. There is also the possibility of over-fitting as a result of being under-powered. Further limitations of the study are extensively discussed in the following Chapter Seven.

6.5.4 Development of a set of diagnostic criteria

Single symptoms are not very useful to clinical practice, and are likely to have limited sensitivity. The idea behind the development of a simple symptoms score was to take an exploratory approach with no intention of over-interpreting the findings. The diagnostic performance of the criteria was broadly interpreted merely to give an indication of the potential for a symptoms score to predict lung cancer in those who had been referred to the lung-shadow clinic.

In the full population dataset, Model 1 differed from Model 2 (modelled with the risk variables), and the better fitting model (Akaike's Information Criterion lower by ≥ 2) was used to develop the criteria, which was the risk model (Model 2). At both levels of significance; $p < 0.05$ and $p < 0.15$, Model 2 had better fit. Previous primary care studies have developed models incorporating the combinations of symptoms and risk factors to estimate the absolute risk of

having lung cancer (Hippesley-Cox and Coupland 2011; Iyen-Omofoman et al. 2013). Comparing their models to the NICE criteria model (validation cohort), the derived models (with smoking, age, COPD, pneumonia, and family history of cancer) performed better than the current NICE referral guidelines.

Based on the analysis for the full population, the optimal cut-off between sensitivity and unnecessary referrals of the stricter model (model consisted of cough, breathing changes/difficulties that was indicated in the last three months, age > 71, and smoking history) was at 2/4. The sensitivity and specificity at this level was 82% and 45%, respectively. This indicates that 82% of those with lung cancer that would be detected on the basis of any two of the four symptoms will be correctly diagnosed. However, 18% of those based on this criteria would have been missed. The positive likelihood ratio was 1.5. The un-weighted diagnostic criteria took no account of the effect size of individual symptoms as each symptom was given equal weight, and therefore, a weighted score was derived where appropriate. In the case of the model with the stricter criteria ($p < 0.05$), no difference in weights (+1) was observed between the weighted and un-weighted criteria.

A simple symptom score was also developed for the referred population with COPD. In the COPD sub-group, the addition of risk variables in Model 2 produced a similar model to Model 1. Therefore, Model 1 was used. At the optimum cut-off, the positive likelihood ratio is 1.88. The AUC statistic was reasonably good at 0.790. Similar to the full population, there was hardly any improvement in the weighted set compared to the un-weighted set.

Although this test does not necessarily inform symptomatic diagnosis of early lung cancer as it was a referred population that was going to be investigated for lung cancer, exploratory analysis on a sub-population with a chronic respiratory disease, such as COPD, might provide more information relevant to primary care.

6.6 Conclusion

Common chest symptoms; cough and breathing changes/difficulties significantly predict lung cancer in this secondary care population at the conservative significance level of 0.05. The findings also concurred with current literature, and supported NICE recommendation for urgent referrals in the UK. However, it needs to be emphasised that the results of the quantitative analysis particularly the analysis of the sub-population was highly exploratory. The hypotheses generated by this research will be further investigated in the larger fully-powered study.

The use of the weighted criteria sets did not improve the discriminatory value for lung cancer, suggesting that they are not likely to be very useful in practice. Differences observed between the full population analysis and the COPD sub-group analysis, however, warrants further investigation in diagnostic studies that control for common co-morbidities such as COPD. Furthermore, COPD diagnosed by spirometry and/or clinical diagnoses by respiratory physicians based on symptomology, might be comparable to primary care-defined COPD, and therefore, the positive predictive values of symptoms in this population using an adequately powered prospective study may inform symptomatic diagnosis in primary care.

Chapter 7: Discussion and Conclusion

7.1 Introduction

Despite the increasing interest in symptomatic diagnosis in lung cancer, research efforts have not identified symptoms that consistently predict lung cancer, except for haemoptysis in primary care populations. In the UK, 86% of patients are still diagnosed at the advanced stage of lung cancer when survival is poor, due to late diagnosis (NHSIC 2011). As most cancers are detected following symptomatic presentation, the best way to improve survival rates in lung cancer patients remains with earlier symptom recognition.

The main aim of this study was to assess the feasibility of using the self-completed IPCARD questionnaire to prospectively collect patient-elicited symptoms, and inform the design of a larger fully-powered study in a secondary care population, with COPD, that had been referred for lung cancer investigation, using both qualitative and quantitative methods.

The feasibility study used qualitative methods to establish acceptability and validity of the IPCARD questionnaire in a secondary care pre-diagnosis population, before carrying out a quantitative analysis to identify symptoms that predict lung cancer in that population (with varying chest and respiratory diseases), and then in a homogeneous sub-population with COPD.

7.2 Study findings

The previous IPCARD Feasibility Study with GP-referred CXR attendees, established content validity and test retest reliability in that population. Although the current population is higher risk, and may be more anxious, it was not anticipated that the content validity would differ considerably between this and the IPCARD Feasibility Study population. The qualitative component of this study (Study 1, Chapter Five) established the acceptability of the IPCARD questionnaire in this population. Completion of the questionnaire did not raise anxiety levels in attendees of the lung-shadow clinic; however, inadequate time to complete the questionnaire was identified by some participants, and may

Discussion

result in some missing data. The IPCARD questionnaire was able to fully capture the range of symptoms experienced in this COPD population, with most normalised symptoms (bodily sensations and health changes normalised by participants) elicited. However, the qualitative study suggested that the symptom experiences of changes in cough could be better recorded using closed questions, rather than the current open-question format.

The two symptoms found to discriminate between patients with and without lung cancer in this referred secondary care population (at the significance level $p < 0.05$) were cough that was first indicated in the last three months (OR=1.74, CI=0.91-3.32), and breathing changes/difficulties that were first indicated in the last three months (OR=1.86, CI=0.97-3.57). This suggests that only 'new' respiratory symptoms predict lung cancer in this secondary care population with high rates of chest and respiratory disease. This finding is consistent with the findings of the systematic review (Chapter Three), that few of the symptoms identified in primary care studies appear to distinguish between those with and without lung cancer in secondary care.

Comparing current findings to an earlier secondary care study, Kubik et al. (2001) reported an increased lung cancer risk with chronic cough (ORs=1.99; CI=0.8-4.9), and shortness of breath (ORs=1.48; CI=0.8-2.7). However, symptoms indicated were chronic rather than recent (in the last three months). The prospective study only investigated common chest symptoms included in the standardised MRC questionnaire, and recruited only women in the study (Kubik et al. 2001).

In a primary care population, Hamilton et al. (2005) found clinical features presentation of a second attendance with cough (OR=2.7; CI=1.7-4.4), and dyspnoea (OR=4.7; CI=2.7-8.0) to be independently associated with lung cancer. Hippenley-Cox and Coupland (2011) identified a new onset cough in the last 12 months as a predictor in both males (Hazard ratio (HR) =1.47; CI=1.23-1.75) and females (HR=1.90; CI=1.56-2.32).

Many of the symptoms used in the current national guidance for lung cancer referrals are non-specific, and can be attributed to other non-malignant comorbidities. The National Institute for Health and Care Excellence (NICE) recommends urgent referral for presentations of cough that lasted for more than three weeks, or dyspnoea (NICE 2011), while national campaigns such as

the 'Be Clear On Cancer' have been carried out to raise public awareness on 'persistent cough' (CRUK 2014). It is interesting to note that the two symptoms identified in the current prospective study were first noticed or indicated within the last three months, rather than being generic symptoms of cough/chronic cough and breathing changes/difficulties. Therefore, in a population with high rates of chronic respiratory diseases, patients' reports of the onset of symptoms within the last three months might be worth exploring.

The symptoms associated with lung cancer in this study differ from those that predicted abnormal chest x-ray suspicious of lung cancer in the earlier IPCARD Feasibility Study in a lower risk chest X-ray population. Brindle et al. (2014) did not find common chest symptoms (cough for longer than three weeks, generic chest aches/pains, and breathlessness) to predict chest X-ray suspicious of lung cancer. Instead, weight loss, and less common variants of pain (pain in side of chest/ribs, severe pain, and pain that feels like indigestion not associated to eating in patients with less severe pain) predicted suspicious chest X-ray in this lower risk population (relative to the current population). The rates of those with COPD in the chest X-ray referred population were also lower.

Haemoptysis also did not appear to be an independent predictor of lung cancer in this lung-shadow-clinic referred population, or in the IPCARD Feasibility Study, despite evidence suggested in the Systematic Review from primary but not secondary care studies (see Chapter Three). A recent prospective secondary care study, Walter et al. (2015), did identify haemoptysis as a predictor of lung cancer diagnosis. However, the population in this study was at a lower risk than the population in the current study (21.5% compared to 15.9 % diagnosed with LC), and recruitment rates were lower. Although haemoptysis is highly specific, it is only reported in 20% of the patients (less than 10 cases) (Walter et al. 2015; Hamilton et al. 2005). Furthermore, Walter et al. (2015) screened referral letters, and only recruited patients with symptoms suspicious of lung cancer as opposed to the current study population, which included patients referred not on suspicion of lung cancer using continuous sample selection. It should be noted that the spectrums of disease between primary care studies (Hamilton et al. 2005) will be different to

Discussion

secondary studies; the current study and Walter et al. (2015), and therefore findings are not expected to be the same.

One study comparing the frequency of individual symptoms in lung cancer cases and non-cases, found higher records of more suggestive, non-specific symptoms in patients with undetected lung cancer than those without. They also found that the rates of symptom presentations were the same across the disease spectrum (stages I-IV), which suggests that the presentation of most symptoms six months before diagnosis was late; supporting the case for standard targeted screening (smoking) over symptom-based strategies (Ades et al. 2014).

7.2.1 Diagnostic symptoms in COPD population (sub-group)

There is a high prevalence of COPD in lung cancer populations, and overlap between symptoms of COPD and lung cancer; qualitative studies indicate that patient recognition of lung cancer symptoms might be delayed by pre-existing chronic respiratory disease (Kiri et al. 2010). Therefore, prospectively collected symptom data in a COPD population is needed to identify symptoms that would better distinguish between COPD and LC. The current study identified symptoms that distinguish between LC and COPD in a COPD subgroup of the secondary care population. However, this exploratory sub-group analysis was underpowered to identify symptoms with low or moderate effect sizes.

At the more stringent entry criteria (statistical significance at $p < 0.05$ for variable selection), only one breathing change/difficulty variable, 'unable to get enough air in', was found to be significantly associated with lung cancer in those with COPD (OR=3.17, CI=1.18-8.56).

However, at the relaxed entry criteria, arthritis (OR=0.39; CI=0.15-0.98) and four breathing variables; generic breathing changes/difficulties experienced within the last three months (OR=0.20; CI=0.04-1.13), breathing changes/difficulties first indicated within the last three months (OR=5.38, CI=1.16-24.91), unable to get enough air (OR=5.98; CI=1.42-25.22), and wheezing sensation when in a particular position (OR=2.13; CI=0.94-4.80), remained in the COPD-specific model. Three of the four variables are different to the variables in the full population model. However, it is noted that the

variables in this model have high standard errors, and wide confidence intervals, and so should not be over-interpreted.

Nevertheless, a comparison of the full population model and COPD model, suggesting that there could be a difference in symptoms that predict lung cancer in those with COPD, justifies further investigation into stratifying by a specific respiratory disease within a heterogeneous population, as it might be possible to predict lung cancer in homogenous groups.

7.3 Development of a set of diagnostic criteria

Full population model: Creating a set of criteria from the four variables in the logistic regression model with risk variables ($p < 0.05$) resulted in an optimum cut-off of 2/4. The positive likelihood ratio for this cut-off was 1.49.

COPD sub-population: For the sub-population COPD model ($p < 0.05$), the optimum cut-off is at 2/3, which would give a positive likelihood ratio of 1.91. 'Weighting' the criteria so that variables with larger effects sizes were given more weight, performed equally well as the un-weighted criteria. Therefore, there is little advantage to using the weighted set over the simpler, un-weighted set of criteria.

Adding epidemiological risk factors into the symptom models: In the full population dataset, Model 2 (with the addition of the risk variables) differed from Model 1 (symptoms model) at both levels of significance; $p < 0.05$ and $p < 0.15$. Thus, the better fitting Model 2 (risk model) was used to develop the criteria (Akaike's Information Criterion lower by ≥ 2). In contrast, in the COPD sub-group, Model 2 (with addition of risk variables) was not different to Model 1, and therefore Model 1 was used at both levels of significance ($p < 0.05$ and $p < 0.15$).

7.4 Methodological issues

Symptoms are highly subjective, and therefore, research involving symptomology is usually complex, which requires a carefully considered study design. This is reflected in the methodological factors and limitations

Discussion

associated with symptomatic studies in lung cancer; as explained in the systematic review (see Section 3.5.3, Chapter Three). The systematic review highlighted several key methodological weaknesses which included:

- Lack of standardised data collection
- Retrospective study design
- Recording bias
- Selection bias
- Likely potential confounders
- Limited generalisability

The strengths and weaknesses of this study will be discussed in relations to these methodological considerations.

7.5 Strengths and limitations of study

Prospective data collection: Prospective data are less subject to recall bias than retrospective re-interpretation of symptoms in the light of a diagnosis. Although, prospective studies that record patient-elicited symptoms in primary care would be ideal, they are very costly to conduct in order to achieve sufficient numbers of cancer cases. Therefore, a robust, and methodologically sound prospective study in secondary care, recruiting from attendees of a lung-shadow clinic that includes those referred under the two week wait referral, will be useful to provide evidence that informs the justification for future investments in primary care studies.

Patient-elicited symptoms: The study collected patient-elicited data, which are likely to be more complete and accurate than data drawn from medical records (clinician-reported). There is a potential for recording bias in clinician-reported symptoms, which tend to be diagnosis-driven rather than symptom-focused (Kroenke 2001). Therefore, the risk of under-representation of patients' experiences exists, as only those deemed important by the clinician are documented. However, even when symptoms are collected systematically using clinician-elicited questionnaires, there might still be a higher threshold for reporting symptoms to a clinician or GP.

The systematic collection of detailed and comprehensive symptom data: A further strength of the study is in the continuous, systematic method of data collection using a validated symptom questionnaire (IPCARD) that extensively records patient-elicited symptoms. This minimises the issue of recording bias observed in earlier studies, and proves invaluable for addressing some of these inconsistencies in symptom data collection. The IPCARD questionnaire was specifically designed to address some of the methodological issues of eliciting lung cancer symptomology, and included lay descriptors of symptoms experienced by those in late-stage lung cancer and early-stage lung cancer. In this, and previous research, the questionnaire was found to also elicit symptoms normalised by patients (Brindle et al. 2015). Previous research suggested that symptoms described using non-disease terms such as ‘aches’ or ‘discomfort’ were better at eliciting non-specific LC symptoms in interviews than questions that used disease-related labels such as ‘pain’ or ‘breathlessness’ (Brindle et al. 2012). Furthermore, in this study, qualitative research indicated that the questionnaire captured the full range of health-related changes in those with COPD and/or lung cancer. The IPCARD questionnaire has been validated in a GP-referred CXR population, and in the current secondary care population, that had been referred to a lung-shadow clinic with high rates of respiratory problems (Study 1).

Missing data: Where a strength of the study lies in its prospective design, collecting data within an operational clinic makes the questionnaire response prone to missing data with no way of knowing for certain its mechanism of missingness. Multiple imputation (MI) was used to handle missing data. Unfortunately, sensitivity analyses were not achievable within the constraints of a PhD, as they are too time-consuming. That said, Woods et al. (2008) also stated that sensitivity analyses require large sample sizes to be informative. Furthermore, no indications of bias in the multiple imputation model suggests that the model is sound.

Although the study only recruited 359 participants within the period of data collection, the **response rate** was higher than the 60% generally regarded acceptable for surveys, and was therefore considered to be good for a population known to be especially difficult to recruit. The study population is

Discussion

heterogeneous with referrals from GPs, secondary-care, and other tertiary centres to reflect a community-based population. However, the same heterogeneity might also have made it more difficult to identify significant symptoms that predict lung cancer; which are known to include non-specific symptoms that can be attributed to other comorbidities. Therefore, the plan in a future study, to estimate the predictive value of symptoms stratified by a common respiratory disease (e.g. COPD) to improve homogeneity within a heterogeneous population, appears to be a justifiable one.

The biggest limitation of this study is the **small sample size**, i.e. lack of power. Multiple testing can increase the occurrences of Type I error, and one of the ways to overcome this is to increase the significance threshold in the analysis (alpha is reduced, $\alpha < 0.05$) to restrict the number of inferences made. However, one could also argue that the small sample size of the study (the study is underpowered) would serve to counterbalance these errors in inference, as only large effect sizes would be detected at $p < 0.05$. This is a feasibility study, and the validity of findings can be evaluated in the larger study.

Effect modification: Due to the exploratory nature of the analysis and the lack of power, potential interactions were not investigated or included in the models, nor were they included in the imputation process. Interaction studies require large sample sizes in order to detect interactions of realistic magnitude, and therefore caution was taken to avoid spurious analysis from the overestimation of the data's robustness. The inclusion of interactions between symptoms, and between symptoms and epidemiological risk factors in the larger study, might improve the model's predictive accuracy.

Despite the methodological strengths of this study; the prospective, systematic data collection, and the collection of detailed information about a broad range of symptoms, the results from this secondary care study (full population dataset) differed little from previous secondary care studies in the systematic review (see Chapter Three), as only two (respiratory) symptoms that predicted lung cancer were identified. The ability of this model to discriminate between patients with, and without, lung cancer was low. However, the COPD model, whilst having limited power at a more stringent alpha, at a relaxed alpha identified symptoms worthy of further investigation in a fully powered study

that would also include interaction terms. At the optimum cut-off of 3/6, the positive likelihood ratio for the COPD un-weighted model was 1.88, producing an Area Under Curve (AUC) of 0.79, which is considered to be a good level of discriminatory power (Fan et al. 2006).

7.6 Implications for clinical practice and future research

Early detection of LC is of clinical interest due to the favourable relationship between earlier diagnosis and better prognosis (lower mortality rate). Previous studies (Hamilton et al. 2005; Corner et al. 2005) and the present study have shown that individuals with lung cancer experience symptoms before diagnosis. They also found that majority of these symptoms first appeared within three months of diagnosis, with some detectable even four to 12 months before diagnosis. Results from this study are not transferable to primary care due to spectrum effects. However, further investigation is warranted in the COPD subgroup. Symptoms that are found to discriminate between those with, and those without, LC in a population with COPD in secondary care, might also have some diagnostic value in primary care COPD populations with a similar spectrum of disease.

In line with the aims of policy initiatives, such as urgent chest X-ray referrals for higher-risk patients (extra-NICE, Hurt et al. 2013) and the increasing justification for clinical decision support tools (Hamilton et al. 2013; Iyen-Omofoman et al. 2012), future research in the development of a symptom-risk algorithm, based upon symptoms elicited from patients (minimised risk of recording bias) rather than symptoms recorded in GP notes, that assists GPs in primary care to make appropriate referrals for LC investigations, is indicated.

7.7 Conclusion

Findings of the present study are in keeping with the current literature on secondary care research; that there are few symptoms that predict lung cancer. Nevertheless, the evidence presented in this study is more robust methodologically. Furthermore, this feasibility study currently lacks power; a

Discussion

model built from a fully-powered study could include a greater number of predictors and interaction terms and, therefore, have improved predictive accuracy and be more robust. Sensitivity analyses could also be performed, to inform interpretations of the implications of missing data, if required. This current research can, however, establish the feasibility of the study design in a secondary population that had been referred to lung-shadow clinic.

Although, the heterogeneity across patient subgroups, within the full population, limits the comparability and transferability of the findings to primary care, the identification of symptoms that predict lung cancer diagnosis in the homogeneous COPD population, within the wider lung-shadow clinic population, is more likely to add to the limited evidence in primary care. The superior performance of the COPD-specific diagnostic criteria further supports the need for an adequately powered study to investigate the predictive values of LC symptoms in homogeneous populations with specific respiratory diseases.

Appendices

Appendix 1 Revised TNM staging system (2010) (American Joint Committee on Cancer (AJCC) 2002)

Tumor (T) stage	Clinical description
T1 a	Tumour in lung <2cm
T1 b	Tumour in lung 2-3 cm
T2	Tumour 3-7cm (or) grown into the main bronchus (or) tumour has grown into the inner lining of the visceral pleura (or) the tumour has part of the lung collapse. T2a ≤ 5cm ≤ T2b
T3	Tumour >7cm (or) tumour has made the whole lung collapse (or) tumour grown into the chest wall, the central lining of the mediastinal pleura, the diaphragm or the pericardium (or) there is more than one tumour nodule in the same lobe of lung.
T4	Tumour has grown into either the mediastinum, the heart, a major blood vessel, the trachea, the oesophagus, the spine, the vagus nerve (or) tumour nodules in more than one lobe of the same lung.
Nodes (N) stage	
N0	No cancer in any lymph nodes
N1	Cancer in the lymph nodes nearest to the affected lung
N2	Cancer in the lymph nodes in the mediastinum or cancer in the lymph nodes just under the trachea.
N3	Cancer in the lymph nodes on the contralateral side of the chest of the affected lung or lymph nodes above the clavicle or lymph nodes at the apex of the lung
Metastases (M)	
M0	No signs that cancer has spread to another lobe of the lung or organs
M1 a	Tumours in both lungs or malignant pleural effusion
M1 b	Lung cancer cells in distant organs of the body

Appendix 1

Stage	Tumour (T)	Node (N)	Metastases (M)
I a	T1a	N0	M0
	T1 b	N0	M0
I b	T2 a	N0	M0
II a	T1 a	N1	M0
	T1 b	N1	M0
	T2 a	N1	M0
	T2 b	N0	M0
II b	T2 b	N1	M0
	T3	N0	M0
III a	T1 a	N2	M0
	T1 b	N2	M0
	T2 a	N2	M0
	T2 b	N2	M0
	T3	N1	M0
	T3	N2	M0
	T4	N0	M0
	T4	N1	M0
III b	Any T	N3	M0
	T4	N2	M0
	T4	N3	M0
IV	Any T	Any N	M1 a
	Any T	Any N	M1 b

Stage	Description
I	Small tumour, no nodes or metastasis
II	Small tumour, nodes infiltrated but not mediastinal nodes; no metastasis
III a	Large tumour, ipsilateral mediastinal nodes present; no metastasis
III b	Large tumour, Contralateral mediastinal or any scalene or supraclavical nodes; no metastasis
IV	Metastatic disease

Appendix 2 Published paper

Family Practice, 2014, Vol. 31, No. 2, 137–148
doi:10.1093/fampra/cmt076
Advance Access publication 17 December 2013



A systematic review of symptomatic diagnosis of lung cancer

Joanna Shim^{a,*}, Lucy Brindle^a, Michael Simon^a and Steve George^b

^aFaculty of Health Sciences and ^bFaculty of Medicine, Southampton General Hospital, University of Southampton, Southampton, UK.

*Correspondence to Joanna Shim, Faculty of Health Sciences, University of Southampton, Building 45, Highfield Campus, University of Southampton, Southampton SO17 1BJ, UK; E-mail: js1g08@soton.ac.uk

Received August 16 2013; revised October 18 2013; Accepted October 26 2013.

Abstract

Background. Lung cancer (LC) is often diagnosed late when curative intervention is no longer viable. However, current referral guidelines (e.g. UK National Institute for Health and Care Excellence guidelines) for suspected LC are based on a weak evidence base.

Aim. The purpose of this systematic review is to identify symptoms that are independently associated with LC and to identify the key methodological issues relating to symptomatic diagnosis research in LC.

Methods. Medline, Ovid and Cumulative Index to Nursing and Allied Health Literature were searched for the period between 1946 and 2012 using the MeSH terms 'lung cancer' and 'symptom*'. Quality of each paper was assessed using Scottish Intercollegiate Guidelines Network and Consolidated Criteria for Reporting Qualitative Research Checklists and checked by a second and third reviewer.

Results. Evidence regarding the diagnostic values of most symptoms was inconclusive; haemoptysis was the only symptom consistently indicated as a predictor of LC. Generally, evidence was weakened by methodological issues such as the lack of standardized data collection (recording bias) and the lack of comparability of findings across the different studies that extend beyond the spectrum of disease. Qualitative studies indicated that patients with LC experienced symptoms months before diagnosis but did not interpret them as serious enough to seek health care. Therefore, early LC symptoms might be under-represented in primary care clinical notes.

Conclusion. Current evidence is insufficient to suggest a symptom profile for LC across the disease stages, nor can it be concluded that classical LC symptoms are predictors of LC apart from, perhaps, haemoptysis. Prospective studies are now needed that systematically record symptoms and explore their predictive values for LC diagnosis.

Key words: Diagnostic accuracy, epidemiology, lung cancer, lung neoplasm, predictive value, symptoms.

Introduction

Lung cancer (LC) continues to be the leading cause of cancer mortality worldwide, accounting for approximately 1.4 million deaths each year (1). Over the last four decades, survival rates have only improved slightly with most LC being diagnosed at late stages when curative intervention is no longer viable in both developing and developed regions of the world (1,2). Evidence from large population-based studies has since

associated relatively lower survival in some regions with delays in diagnosis (3,4).

In the UK, efforts to improve the survival rate of LC included National Institute for Health and Care Excellence (NICE) guidelines, which recommended urgent chest X-ray referral for patients experiencing any persistent symptoms that might indicate LC (5). However, most of these symptoms are non-specific and

could also suggest other differential diagnosis of lung and respiratory diseases, such as chronic obstructive pulmonary disease (COPD). The most common methods found in symptom-based research in LC are retrospective cohort studies and case-control studies that use clinical records and databases; these studies are generally limited by methodological issues regarding symptom data collection. For example, these studies are subject to the possibility of recording bias by clinicians that might have implications for the predictive values obtained. The absence of prospective studies of the predictive value of symptoms mainly reflects the large study sizes and costs involved in longitudinal studies that commence prior to referral for investigation.

Two systematic reviews have addressed the diagnostic value of symptoms in LC: Hamilton and Sharp (6) and Shapley *et al.* (7). Hamilton and Sharp (6) reviewed features of symptomatic LC across studies and estimated the likelihood ratios (LRs) of some of the symptoms in LC diagnosis. The estimated LRs reported were based on referred populations in secondary care settings with hardly any research to be found in primary care populations. The study of Shapley *et al.* (7) identified symptoms, signs and non-diagnostic test results that were highly predictive of specific cancers [where positive predictive values (PPVs) $\geq 5\%$ were reported]. The review analysed all higher quality evidence of symptoms that predicted LC in an unselected primary care population and reported two studies for LC (8,9). Only haemoptysis was identified as having high PPVs ($\geq 5\%$) in LC diagnosis. Owing to the lack of high-quality research in primary care populations in the most recent systematic review identified, this systematic review will investigate the diagnostic value of symptoms for LC, regardless of national health care system or spectrum of disease, and identify any new primary care evidence since 2010. The review will also include qualitative studies to explore the symptom experience of people diagnosed with LC and identify factors associated with patient reporting of symptoms that might have implications for the design of future diagnostic studies. This qualitative component could also reveal any non-classical symptoms or characteristics of symptoms experienced before LC diagnosis not investigated in diagnostic studies.

Methods

Search strategy

Electronic databases were searched from their commencement to July 2012 using search terms 'lung cancer*', and 'symptom*' for title abstracts. The key terms were exploded to include alternative MeSH descriptors such as 'Lung Neoplasms', 'Signs and Symptoms', and 'Differential Diagnosis'. Search hits were filtered using qualifying restrictions: diagnosis, epidemiology and aetiology for 'Lung Cancer', 'Symptoms' and 'Differential, Diagnosis'. Details of the complete search strategies of all the databases performed are included in Appendix I (see [supplementary data](#)).

Similar search methods were applied for all the electronic databases listed below:

Electronic databases
<ul style="list-style-type: none"> • MEDLINE (from 1946 to July Week 1 2012) • Embase (1946 to Week 28 2012) • Cumulative Index to Nursing and Allied Health Literature (CINAHL) (1981 to 16th July 2012) • Multi-database: Embase (1980 to 2012 Week 2); Ovid MEDLINE (1946 to July Week 1 2012); Ovid MEDLINE Daily Update (July 13 2012); Ovid MEDLINE In-Process & Other Non-Indexed Citations (July 13 2012) • Cochrane Library

In addition, the contents pages of four journals (two for quantitative and two for qualitative studies) between 1 January 2009 and 31 December 2011 were hand searched: Thorax and the British Journal of General Practice for quantitative studies and the Psychooncology and the European Journal of Cancer Care for qualitative studies. Based on the search strategy, journals with the highest number of relevant papers for quantitative and qualitative studies were selected. This generated a total of 3017 papers (1830 quantitative and 1187 qualitative). All records were retrieved and screened for relevance.

Inclusion and exclusion criteria

The following criteria were used to determine the eligibility of studies for this review:

- Quantitative study design—Studies that reported diagnostic values (PPVs, hazard ratios (HRs), odds ratios (ORs) and/or LRs) for the symptom, sign or test, or provide the necessary information needed to calculate these values (2×2 contingency tables could be reconstructed).
- Qualitative study design—Studies that explored the trajectory of symptoms from when symptoms were first experienced before diagnosis, or studies that described the onset of symptoms or first symptoms at presentation to clinician (primary care) were of interest for the purposes of the review.
- Participants—Only adult populations recruited from hospitals, outpatient clinic, specialist clinic, specific community or the general population.
- Outcomes—The group with the positive outcome (LC) must have had a confirmed clinical diagnosis of LC that met diagnostic standards set by the health service provider.
- Others—Studies written in a language other than English, German, Spanish, Malay and Chinese were excluded. Studies on multi-site cancers were included provided that LC was distinguished from other cancers in reporting of results.

Exclusion criteria for the review were:

- Study design—Studies that reported symptoms post-treatment were excluded. These included studies on the management of

symptoms in advanced LC, studies measuring the effect of toxicity and quality of life studies on symptom burden where baseline reported only post-treatment symptoms that will not provide diagnostic values. Single case studies, case reports, editorials, symposiums, reviews (literature) and practical guidelines were excluded.

- Participants—Studies that reported symptoms of metastatic cancer, where LC is the secondary cancer, were excluded.

Study selection and quality assessment

The initial screening of the titles and abstracts was carried out independently by the first reviewer (JS). A second and third reviewer (LB and MS) each checked a random sample of 75 (2.5%) of the abstracts. All papers shortlisted were retrieved in full. A second reviewer (LB) checked 100%, and a third reviewer (MS) 25%, of the full papers that were shortlisted to ensure that they met the eligibility criteria. Methodological quality was assessed using the Scottish Intercollegiate Guidelines Network Checklist for cohort studies and case-control studies, and the Consolidated Criteria for Reporting Qualitative Research checklist was used to assess qualitative studies. Disagreements or uncertainties about satisfaction of quality criteria were discussed with the second and third reviewers (LB and MS) and consensus achieved.

Data extraction and analysis

The reviewers requested raw data of potentially relevant data from main authors to ensure a comprehensive inclusion of existing literature. Data on the type of study, characteristics of the study population, duration of follow-up and the effect sizes were extracted systematically and tabulated for each study that met inclusion criteria. If diagnostic values were not reported, the positive and negative likelihood ratios and sensitivity and specificity were calculated using the 2×2 contingency tables (10). PPVs were also calculated where possible.

Owing to the large variance in populations and research methodologies between the studies in this review, a meta-analysis to pool the predictive values of the symptoms reported in all the studies was considered to be inappropriate. As a result, the study used narrative summaries (narrative review approach) and meta-synthesis of data to analyse the quantitative and qualitative evidence, respectively.

Results

A total of 6037 papers were retrieved using the search strategy (11). Result of the search strategy and selection process is shown in the flow diagram in Figure 1. Duplicates were removed using EndNote reference manager. The final update was performed

in July 2012. In total, 9054 articles (including the 3017 hand-searched journals) were assessed for relevance. Out of which, 11 studies (5 quantitative and 6 qualitative) were eligible for inclusion in the final review.

Results of quantitative studies

Table 1 and Appendix (II) (see supplementary data) summarizes the main characteristics of the quantitative studies.

Symptomatic prevalence

In general, cough, dyspnoea and haemoptysis were the most commonly 'measured' symptoms, based on patient-reported symptoms and those recorded by a general practitioner (GP)/primary care clinician (shown in Table 2). In the four studies that reported symptom frequencies, systemic symptoms such as appetite loss, weight loss, fatigue and fever/flu were less frequently reported (8,12,13). On the whole, the duration of symptoms onset to diagnosis reported in the studies ranged from 6 months to no more than 2 years before diagnosis.

Symptomatic diagnosis

Table 2 and Appendix (III) (see supplementary data) show symptoms that were recorded in the study and which of these predicted LC ($P \leq 0.05$). In two of the studies, P -values or confidence intervals (CIs) were not provided (12,14).

Table 3 present the PPVs and ORs of the individual symptoms reported in the case-control studies and cohort studies, respectively. Most symptoms reported in the study of Hamilton *et al.* (8) had reasonably high ORs (OR 4.40 to 16.24). However, even with the high ORs, the likelihood of a patient presenting with that symptom having LC (as indicated by the PPV) was low (<1.0 for most symptoms).

PPVs of 2.0 and higher ($2.4 \leq \text{PPV} \leq 7.5$) were consistently found for the symptom haemoptysis across the studies. Jones *et al.* (9) evaluated the diagnostic value of haemoptysis in LC, reported an increase in PPV from 5.8, at 6 months after symptom onset, to 7.5—3 years after the first symptom occurrence was recorded in the male group. Similar increase was observed in the female cohort: 3.3—4.3. A gender difference in PPV was also identified (9). However, no statistically significant differences in gender-specific PPVs were identified in Hippenley-Cox and Coupland's study (2011) (13); the remaining three studies did not carry out separate analysis for men and women (8,12,14).

Hippenley-Cox and Coupland (13) calculated the HRs of individual symptoms in the derivative cohort to develop a model strategy that predicts LC risk in a population (13). Only variables of HR <0.80 or HR >1.20 were included into the final model. The study also observed that haemoptysis had a higher HR in comparison to the other symptoms for the final model for LC for both genders (Male HR: 21.5 and Female HR: 23.9) followed by weight loss, appetite loss and cough (symptom variables included

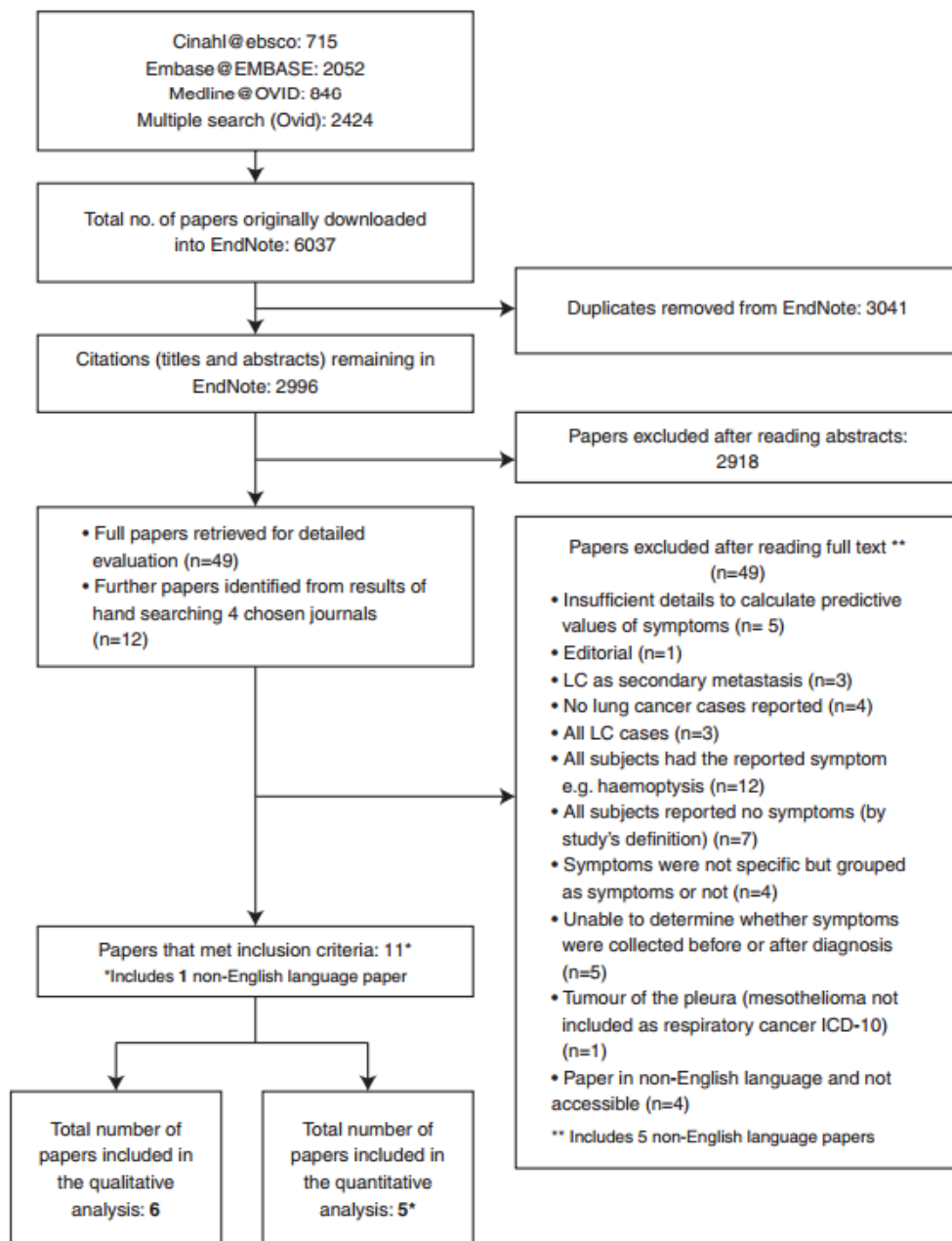


Figure 1. Flow diagram of results of search strategy adopted from the QUOROM statement flow diagram (11)

Table 1. Study methodology—study design and characteristics of exposure data (symptom)

Study (year)	Country (city)	Study design	Data source of LC	Sample characteristic	Period symptom data were collected before diagnosis	Method of recording symptom
Hamilton <i>et al.</i> (2005)	UK	Case-control study	GP records	Population of Exeter aged above 40 years with cases cohort diagnosed with primary LC	≤2 years	Coded using the ICD-2 coding system
Hoppe (1977)	Germany	Cohort study	Medical records	People attending a chest clinic in the state of Northrhine-Westphalia on suspicion of LC	<6 months	–
Hippesley-Cox and Coupland (2011)	UK	Cohort study	GP record (using ICD-9 or ICD-10 codes), patient's electronic record (EMIS), QResearch database (including 564 practices in England and Wales)	Population of primary care patients registered to a GP practice in England and Wales; 30–84 years	2 years	Established predictors were recorded from patient's electronic records using a symptom checklist
Jones <i>et al.</i> (2007)	UK	Cohort study	GP records, patient's record	Primary care population registered to a GP practice aged ≤100 and had reported the occurrence of haemoptysis before	≤3 years	Occurrences of alarm symptom recorded from patient's record
Kubík <i>et al.</i> (2001)	Czech Republic	Case-control study	Participant reported	Cases were female Czech patients with confirmed LC receiving treatment at the hospital in Prague. Controls were women spouses, relatives or friends of other patients of the same hospital (aged 25–84 years)	<2 years	MRC Questionnaire on respiratory symptoms (interviewer administered).

ICPC, International Classification of Primary Care; ICD, International Classification of Diseases; EMIS, Egton Medical Information Systems.

in the final model). In their final model, the top 0.5% of their risk score produced a PPV of 9.5 (8.8–10.3).

The remaining two studies did not identify any statistically significant associations of symptoms with LC (12,14).

Key methodological strengths and limitations of the studies

The included studies displayed several methodological weaknesses that included: the lack of standardized data collection, the potential for recall bias and recording bias, retrospective design, selection bias, confounding and limited generalizability.

Data collection issues

A variety of data collection methods were applied across the studies. Four of the studies extracted their symptom data retrospectively using medical and GP records or medical databases (8,9,12,13); Kubík *et al.* (14) used a standardized questionnaire,

the Medical Research Council (MRC) respiratory questionnaire (see Appendix (IV) in [supplementary data](#)), to prospectively record symptoms. A strength of most studies was that they included symptoms reported and recorded at the time of presentation to the clinician, as far as indicated, which reduced the likelihood of recall bias or retrospective reinterpretation of symptoms following diagnosis.

However, reports of symptoms based on these records can be restricted as they only reflect the occurrence of symptoms at the time of reporting and often tend to be diagnosis focused (15) rather than symptom focused, potentially only including symptoms thought to be relevant to a differential diagnosis, resulting in partial recording of symptoms. Therefore, medical notes or records can be subject to recording bias, as clinicians may have recorded symptoms more thoroughly if LC was suspected (8). The use of prospectively completed checklists or questionnaires might avoid recording bias if administered systematically to patients. However, this

Table 2. Symptoms reported that were independently associated with LC

Study (year)	Criteria for inclusion in the final model	Symptoms													
		Cough	Haemoptysis	Weight loss	COPD	Dyspnoea	Chest pain	Appetite loss	Wheezing	Cough and phlegm	Fatigue	Hoarseness	Heaviness in chest	Fever/flu	General unwell
Hamilton <i>et al.</i> (2005)	$P < 0.01$	✓	✓	✓	✓	✓	✓	✓	-	-	✓	-	-	-	-
Hippesley-Cox and Coupland (2011)	HR <0.9 or >1.2	✓	✓	✓	✓	-	-	✓	-	-	✓	-	-	-	-
Kubik <i>et al.</i> (2001)	-	✓	-	-	-	✓	-	-	✓	✓	-	-	-	-	-
Jones <i>et al.</i> (2007)	-	-	✓	-	-	-	-	-	-	-	-	-	-	-	-
Hoppe (1977)	-	None of the symptom diagnostic values were found to be statistically significant													

Check mark denotes symptom recorded in the study.

stands the risk of only recording common chest symptoms (see Appendix (IV) in [supplementary data](#) for MRC questionnaire) and excluding systemic symptoms (14).

Retrospective versus prospective study design

In most of the studies (8,9,12,13), the symptoms were recorded/measured before diagnosis but were obtained from previously recorded data (medical records) not systematically recorded for research purposes. The limitations of retrospective study design for these studies largely relates to recording bias.

Incomplete reporting of symptoms and recording bias

In general, recording bias could over inflate the LRs of symptoms, if GPs/clinicians are more likely to record symptoms if LC is suspected. Similarly, the ratios could also be under-estimated if patients with LC under-reported their symptoms which would then go unrecorded (13). Four of the five quantitative studies evaluated within this review used clinical records to ascertain symptoms. Therefore, the ORs and LRs obtained might lead to the under- or over-estimation of the true predictive values of symptoms.

Selection bias

Studies using clinical records that capture data on every patient in a region are less likely to be subject to selection bias than studies that rely on the recruitment or inclusion of individual patients (8,9,13). However, the eligibility criteria for LC diagnoses in some of the clinical records-based studies included in this review could have some selection effect. For example, Kubík *et al.* (14), using a prospective study design (recruitment rate 87%) with the potential for self-selection bias, only included histologically confirmed LC cases. Therefore, it is possible that those with advanced LC were more likely to be excluded (advanced LC patients will less likely be subjected to further invasive investigation to confirm diagnosis).

Risk factors

Two of the studies did not adjust for the most likely risk factors (COPD and smoking status) (9,12). Smoking status and COPD have independent associations with LC and might cause symptoms; from 40% to 90% of LC is preceded by COPD (16). Therefore, COPD and smoking might account for the observed effect of a symptom on the disease outcome (LC) in these studies.

Limited generalizability

All studies recruited adults of both genders above the age of 18 (see Table 1) except Kubík *et al.* (14); this study only represented women ≥ 18 , in the Czech Republic presenting at a secondary care centre at the time of recruitment. This inevitably reduces the generalizability of their findings to the general population.

Table 3. Diagnostic values (ORs, PPVs and HRs) for symptoms reported in case-control studies and cohort studies

	Case-control studies				Cohort studies							
	Kubik <i>et al.</i> (2001)		Hamilton <i>et al.</i> (2005)		Hoppe (1977)		Jones <i>et al.</i> (2007)		Hippesley-Cox and Coupland (2011)			
	Cases	Controls	Cases	Controls								
	140	280	247	1235								
Time period before diagnosis	2 Years		2 Years		6 Months	6 Months after symptom recorded	3 Years after symptom recorded	2 Years				
Symptoms	PPV* (95% CI)	OR (95% CI)	PPV* (95% CI)	OR (95% CI)	OR (95% CI)	PPV	PPV	PPV		Only HR <0.8 or >1.2 included HR (95% CI)		
Gender	-	-	-	-	-	M	F	M	F	M	F	
Cough	-	-	0.40 (0.3-0.5)	4.40	-	-	-	-	-	1.47 (1.23-1.75)	1.9 (1.56-2.32)	
Chronic cough with phlegm	0.44 (0.28-0.6)	1.92	-	-	-	-	-	-	-	-	-	
Chronic bronchitis	-	-	-	-	1.51 (0.85-2.67)	-	-	-	-	1.51 (1.34-1.69)	1.82 (1.57-2.11)	
Dyspnoea	-	-	0.66 (0.5-0.8)	6.99	-	-	-	-	-	-	-	
Haemoptysis	-	-	2.40* (1.4-4.1)	16.24	0.96 (0.32-2.82)	5.8 (5.0-6.7)	3.3 (2.3-4.3)	7.5 (6.6-8.5)	4.3 (3.4-5.3)	21.5 (19.3-23.9)	23.9 (20.6-27.6)	
Chest pain	-	-	0.82 (0.6-1.1)	4.92	0.79 (0.34-1.82)	-	-	-	-	-	-	
Weight loss	-	-	1.1 (0.8-1.6)	8.14	-	-	-	-	-	6.09 (5.33-6.95)	4.52 (3.8-5.38)	
Appetite loss	-	-	0.87 (0.6-1.3)	5.69	-	-	-	-	-	4.71 (3.69-6.1)	4.14 (3.15-5.45)	
Fatigue	-	-	0.43 (0.3-0.6)	3.07	-	-	-	-	-	-	-	
Abnormal spirometry	-	-	1.6 (0.9-2.9)	9.39	-	-	-	-	-	-	-	
Worsening cough	-	-	-	5.45 (3.81-7.79)	-	-	-	-	-	-	-	
Flu/fever	-	-	-	-	0.58 (0.21-1.6)	-	-	-	-	-	-	
General unwell	-	-	-	-	0.76 (0.31-1.85)	-	-	-	-	-	-	

OR, odds ratio.

*Positive predictive values (PPVs) in %.

†For a separate validation cohort of the same study.

Most of the studies were of unselected primary care populations (8,9,12,13). Therefore, the results might be transferable to consulting patients in comparable primary care populations. However, the spectrum of disease will differ between referred secondary care and primary care populations and ratios (LRs, ORs, HRs) and PPVs obtained in secondary care populations [e.g. Kubik *et al.* (14)], or already referred primary care populations [e.g. Hoppe (12)] are unlikely to be generalizable to patients presenting with symptoms in primary care; Hoppe (12) retrospectively extracted data from medical records in primary care but specifically selected only those who were referred to chest clinics on the basis of possible LC. These latter two studies involving populations referred to secondary care only identified one symptom with an independent association with LC (cough with phlegm) (14).

Results of qualitative studies

Table 4 summarizes the qualitative studies according to the characteristics of the sample population, data collection method and data analysis performed. The purpose of the review of the qualitative studies was to explore patient's interpretation of symptoms before diagnosis, and whether patients recalled any change in symptom with disease duration to provide a more complete overview of symptoms experienced before LC diagnosis. All studies were retrospective but interviews were conducted close to the time of diagnosis and therefore likely to represent the patient's symptom experience before diagnosis (17–22).

Time intervals between symptom onset and diagnosis

Using an interval event chart to demonstrate the time intervals between key events in the pathways to diagnosis, operable and inoperable LC patients recalled health changes that occurred 12 months (median) before the time of diagnosis (18). However, the study was underpowered and found no significant differences between the two groups ($P > 0.05$) across the time intervals between key events and time of diagnosis. Patients' recollections of key events were verified using GP and medical notes, and high levels of agreements were obtained (18).

Pre-diagnostic bodily experiences (symptoms and health changes)

Participants reported experiencing a broad spectrum of bodily experiences prior to diagnosis (18,20–22). Both systemic (lethargy, weakness, fatigue, weight loss and appetite change) and chest and respiratory symptoms were reported. Often, systemic symptoms have been related to advanced stages of the disease in previous literature but several studies have reported systemic symptoms in patients with operable LC or patients in the earlier stage of the disease (18,20,22).

Participants were unable to discern between 'normal' and 'symptomatic', particularly systemic symptoms resulting in the normalization of symptom (18,20); where symptoms were attributed to occurrences in daily living or 'everyday' bodily changes rather than a health problem (18,20,22). In light of this 'normalization' of symptoms, symptoms then had to become severe before they were presented to the GP/clinician (19).

Generally, most of the literature agreed that the existence of co-morbidity (especially those respiratory related) often complicated the process of becoming aware of a new and different disorder (masks new symptom) (19,21,22). There was a tendency to attribute symptoms to other acute or chronic conditions (21).

Discussion

This review has updated evidence obtained from previous reviews and, unlike previous reviews, was not limited to primary care studies. As there is a lack of evidence about early symptom epidemiology in LC, the inclusion of studies in secondary care and non-UK health systems could potentially enable the identification of diagnostics values of symptoms across a broader spectrum of the disease. However, diagnostic values of symptoms obtained in non-UK or non-primary care health services research are not generalizable to symptoms presented in UK primary care (23).

Based on the PPVs ($P > 0.05$), there is little evidence to suggest that symptoms other than haemoptysis consistently predicted LC. This is in keeping with previous studies and reviews (6). However, individual studies within this review have identified other symptoms that were independently associated with an LC diagnosis such as appetite loss, weight loss, fatigue and fever/flu presentations; some of which, for example, appetite loss despite increasing LC risk from 4- to 5-fold (13), are currently not included in the UK NICE guidelines as grounds for referral and may be worthy of further investigation. That being said, stronger claims beyond this cannot be made due to the methodological difficulty of non-comparability of findings across the different studies that extend beyond the spectrum of disease. For example, methods of recording symptoms were highly divergent between studies.

Four of the five quantitative studies reviewed used routine data sets; i.e. symptoms recorded in GP notes or electronic medical databases, with the potential for recording bias and incomplete presentation of symptoms (8,9,12,13). The one study that did prospectively and systematically record symptoms used a predefined symptom checklist, raising the possibility that symptoms with diagnostic value for LC were omitted (14). Studies have not yet to date, systematically recorded systemic symptoms as well as non-systemic respiratory symptoms prospectively for research purposes. The possibility of recording biases by GPs/clinicians, and the incomplete recording of symptoms in published

Table 4. Summary of qualitative studies

Study (year)	Aim	Study design	Method of data collection	Period of data collection	No. recruited/eligible (%)	Sample characteristics	Method of data analysis
O'Driscoll <i>et al.</i> (1999)	Recorded detailed notes of how patients described their breathlessness and its impact have been analysed and presented in order to offer descriptive material regarding the experience of breathlessness in LC. These data were from a previous study looking to develop and evaluate a breathlessness intervention	Retrospective	Assessment notes recorded by nurse research-practitioner regarding the reported experience of the patient, or how both nurse and patient agree to describe the symptom and the patient's coping strategies (a collaboration)	Notes were recorded at each subsequent visit to the nursing clinic	52	Thirty men (58%), and 22 women (42%) with a LC diagnosis who had completed chemotherapy or radiotherapy and experienced breathlessness aged ranging from 33 to 76 years (mean age: 60 years). Patients must have been healthy enough to be able to provide adequate material for data analysis	Content analysis—with frequency counts; descriptive analysis
Comer <i>et al.</i> (2005)	To develop a detailed picture of the pathway to diagnosis by mapping the pre-diagnosis symptom history and the events leading up to diagnosis of a group of patients diagnosed with LC	Retrospective	Directed interviews—semi-structured (entitled 'what happened to me?') and structured approach	Twenty patients—interviewed between 3 days and 4 weeks post-diagnosis, and two patients were interviewed 2–3 months after diagnosis	22/30 (73)	Twelve males and 10 females, aged 42–82 years recently diagnosed with LC to map. A third (<i>n</i> = 7) of the patients had operable LC and the remaining (<i>n</i> = 15) had inoperable LC	Thematic analysis
Comer <i>et al.</i> (2006)	To further analyse the data of previous study (Comer <i>et al.</i> 2005) to re-address unexpected findings/themes that emerged so as to better understand how individuals through the way they responded to their health changes, might have influenced the timing of their LC diagnosis	Retrospective	Results were from the same study conducted (Comer <i>et al.</i> 2005) therefore the study design and sample characteristic are the same (refer to the study above).				

Table 4. Continued

Study (year)	Aim	Study design	Method of data collection	Period of data collection	No. recruited/eligible (%)	Sample characteristics	Method of data analysis
Levealahri <i>et al.</i> (2007)	To explore how people with inoperable LC frame and conceptualize the onset of their sickness; testing a theoretical construct of viewing illness as a biographical continuity over existing theory of disruption	Retrospective narratives	Semi-structured explorative interviews with open questions	Within 1 year post diagnosis	37	Thirty-seven patients (21 women and 6 men) aged 48–86 years who 'survived' the first year post LC diagnosis. Twenty-four had their LC staged (19 were diagnosed at stage IIIb–IV, and 5 people at 'earlier' stages) but all inoperable LC. Participants survived between 2 months to over 3 years after the interview.	Narrative analysis—deductive approach
Tod <i>et al.</i> (2007)	To identify factors influencing delay in reporting symptoms of LC	Retrospective	Semi-structured interviews. Questions were informed by the study of Corner <i>et al.</i> (2003).	Within 1 year after LC diagnosis of 6–18 months	20	Purposive sample of people diagnosed with LC in the previous 6 months or longer. 18 participants were diagnosed 6 months ago and 2 were 18-month survivors. 8 females and 12 males with age ranging from 47 to 81 years.	Framework analysis (Ritchie & Spencer 1994)
Molassiotis <i>et al.</i> (2010)	To map the pathway from initial persistent change in health to diagnosis of cancer and explore the patient and system factors mediating this process	Retrospective accounts	In-depth interviews opened with broad questions, asking patients' to recall when they first became aware of a change in their health	In 2–3 weeks after initial cancer diagnosis and prior to or at initiation of treatment	75 cancer diagnosis; 14 (18.7%) LC diagnosis	In general, the sample characteristic represent patients from seven diagnostic groups, including LC ($n = 14$). Patients aged from 18 to 93 years (mean 58.5 years). Consecutive sample consisted of attendees of outpatient cancer clinics at a large centre hospital.	Content analysis—with frequency counts; descriptive analysis

prospective studies, limits the value of current evidence regarding the diagnostic value of symptoms for LC. Undoubtedly, a symptom profile specific to early stage of the disease would be useful as reduction in mortality is related to early stage LC diagnosis. However, beyond reporting the frequency of symptoms, very few studies explored and even fewer reported the severity or change in frequency of the symptoms that might occur over time. Hamilton *et al.* (8) reported episodes of cough with each subsequent clinic consultation suggesting persistence of the symptom but did not report severity or duration of symptoms as they were not recorded in most notes (personal communication). The level of reporting of the disease staging and histology also varied across the studies. Some studies provided more detail of the types of LC histology than others but none of the studies included in this review reported the stages of the disease (e.g. I-II, Ia-Ib or TNM). One study reported the percentage of LC cases with operable LC (12) but none of the quantitative studies distinguished between symptoms of operable and inoperable LC.

The qualitative studies report LC patients' recollections of their symptom experience before LC diagnosis, with the potential for reinterpretation of symptoms in light of a diagnosis. These studies indicated that patients experienced bodily changes including systemic and non-systemic symptoms months before diagnosis. Although they reported the occurrence and interpretation of symptoms, the studies did not report the characteristics of symptoms experienced before LC diagnosis in any detail. These retrospective reports indicated that some symptoms were not interpreted as being serious or alarming enough to prompt help seeking and, therefore, were not presented to primary care clinicians until late. Furthermore, less severe and normalized symptoms might never be presented to primary care clinicians by consulting patients (24). This could possibly have implications for the symptoms recorded in the GP clinical notes. Therefore, primary care records-based studies might not provide predictive values of early symptoms not fully elicited by clinicians, or not thought serious enough by the GP/clinician to record. The literature also suggests that symptoms that might precede LC further complicated the process of recognizing new symptoms for both patients and clinicians, when there is a pre-existing co-morbidity such as COPD.

This is the first review to have included qualitative studies of the available evidence in symptomatic detection of LC. Although qualitative studies do not inform the diagnostic value of symptoms in LC, they provide information about participants' interpretations of symptoms occurring before diagnosis that might inform the design of future prospective studies and methods of symptom data collection. For example, systematic methods of eliciting and recording symptoms are required in primary care studies. Furthermore, having recognized that patients could easily under-appraise the significance of symptoms, it is important that community-based prospective studies are designed to accurately record symptoms that patients do not interpret as serious or present to GPs/clinicians.

Recommendation for future research

At present, there is not enough evidence to suggest a signs-and-symptom profile for LC from the pool of study samples. Prospective studies that record symptoms systematically, and identify symptoms experienced at early stages of LC, are required in order to strengthen the evidence base for symptomatic diagnosis in LC. Future research might also explore the predictive values of changes in symptoms in those with pre-existing chronic respiratory disease, which is highly prevalent in LC patients.

Supplementary material

Supplementary material is available at *Family Practice online*.

Declaration

Funding: Faculty of Health Sciences, University of Southampton (PhD studentship to JS).

Ethical approval: The study does not require ethical approval.

Conflict of interest: All authors have no conflicts of interest to declare. We certify that there is no conflict of interest with any financial organization, individual authors' commitments, editors or reviewers regarding the material discussed in the manuscript.

References

1. GLOBOCON 2008. Lung Cancer Incidence and Mortality Worldwide in 2008. <http://globocan.iarc.fr/factsheets/cancers/lung.asp> (accessed on 1 August 2013).
2. Janssen-Heijnen ML, Coebergh JW. The changing epidemiology of lung cancer in Europe. *Lung Cancer* 2003; 41: 245–58.
3. Richards M. EURO-CARE-4 studies bring new data on cancer survival. *Lancet Oncol* 2007; 8: 752–3.
4. Walters S, Maringe C, Coleman MP *et al.*; ICBP Module 1 Working Group. Lung cancer survival and stage at diagnosis in Australia, Canada, Denmark, Norway, Sweden and the UK: a population-based study, 2004–2007. *Thorax* 2013; 68: 551–64.
5. National Institute for Health and Clinical Excellence. *Referral Guidelines for Suspected Cancer*. London: National Collaborating Centre for Primary Care, 2005.
6. Hamilton W, Sharp D. Diagnosis of lung cancer in primary care: a structured review. *Fam Pract* 2004; 21: 605–11.
7. Shapley M, Mansell G, Jordan JL, Jordan KP. Positive predictive values of ≥5% in primary care for cancer: systematic review. *Br J Gen Pract* 2010; 60: e366–77.
8. Hamilton W, Peters TJ, Round A, Sharp D. What are the clinical features of lung cancer before the diagnosis is made? A population based case-control study. *Thorax* 2005; 60: 1059–65.
9. Jones R, Latinovic R, Charlton J, Gulliford MC. Alarm symptoms in early diagnosis of cancer in primary care: cohort study using General Practice Research Database. *BMJ* 2007; 334: 1040.
10. Altman DG, Bland JM. Diagnostic tests. 1: Sensitivity and specificity. *BMJ* 1994; 308: 1552.

11. Moher D, Cook DJ, Eastwood S, Olkin I, Rennie D, Stroup DF. Improving the quality of reports of meta-analyses of randomised controlled trials: the QUOROM statement. Quality of Reporting of Meta-analyses. *Lancet* 1999; 354: 1896–900.
12. Hoppe R. [An analysis of 20,000 cases of suspected lung cancer (author's transl)]. *Prax Klin Pneumol* 1977; 31: 872–84.
13. Hippisley-Cox J, Coupland C. Identifying patients with suspected lung cancer in primary care: derivation and validation of an algorithm. *Br J Gen Pract* 2011; 61: e715–23.
14. Kubík A, Zatloukal P, Boyle P *et al.* A case-control study of lung cancer among Czech women. *Lung Cancer* 2001; 31: 111–22.
15. Kroenke K. Studying symptoms: sampling and measurement issues. *Ann Intern Med* 2001; 134: 844–53.
16. Young RP, Hopkins RJ, Christmas T, Black PN, Metcalf P, Gamble GD. COPD prevalence is increased in lung cancer, independent of age, sex and smoking history. *Eur Respir J* 2009; 34: 380–6.
17. O'Driscoll M, Corner J, Bailey C. The experience of breathlessness in lung cancer. *Eur J Cancer Care (Engl)* 1999; 8: 37–43.
18. Corner J, Hopkinson J, Fitzsimmons D, Barclay S, Muers M. Is late diagnosis of lung cancer inevitable? Interview study of patients' recollections of symptoms before diagnosis. *Thorax* 2005; 60: 314–9.
19. Corner J, Hopkinson J, Roffe L. Experience of health changes and reasons for delay in seeking care: a UK study of the months prior to the diagnosis of lung cancer. *Soc Sci Med* 2006; 62: 1381–91.
20. Leveälähti H, Tishelman C, Ohlén J. Framing the onset of lung cancer biographically: narratives of continuity and disruption. *Psychooncology* 2007; 16: 466–73.
21. Tod AM, Craven J, Allmark P. Diagnostic delay in lung cancer: a qualitative study. *J Adv Nurs* 2008; 61: 336–43.
22. Molassiotis A, Wilson B, Brunton L, Chandler C. Mapping patients' experiences from initial change in health to cancer diagnosis: a qualitative exploration of patient and system factors mediating this process. *Eur J Cancer Care (Engl)* 2010; 19: 98–109.
23. Summerton N. Cancer recognition and primary care. *Br J Gen Pract* 2002; 52: 5–6.
24. Brindle L, Pope C, Corner J, Leydon G, Banerjee A. Eliciting symptoms interpreted as normal by patients with early-stage lung cancer: could GP elicitation of normalised symptoms reduce delay in diagnosis? Cross-sectional interview study. *BMJ Open* 2012; 2: pii: e001977.

Appendix 3 Results of database search strategy (Supplementary Data)

Databases and years searched	#	Search files/ Search terms	Number retrieved	Number of hits reviewed
Ovid MEDLINE(R) 1946 to June Week 3 2012	1	exp *Lung Neoplasms/di, ep, et [Diagnosis, Epidemiology, Etiology]	17904	
	2	sympto*.ti,ab.	625199	
	3	exp *Diagnosis, Differential/	8001	
	4	#2 or #3	632998	
	5	#1 and #4		846
Embase @ EMBASE 1980 to 2012 Week 25	1	exp *lung tumor/di, ep, et [Diagnosis, Epidemiology, Etiology]	38924	
	2	sympto*.ti,ab.	880332	

Appendix 3

	3	*cancer diagnosis/	12078	
	4	#2 or #3	891622	
	5	#1 and #4		2052
CINAHL @ EBSCO	S1	MH Lung Neoplasms	13182	
	S2	(MH "Lung Neoplasms+/DI/EP/ET")	3258	
	S3	TI sympto*	23670	
	S4	AB sympto*	88534	
	S5	(MH "Diagnosis")	4423	
	S6	S3 or S4 or S5	105513	
	S7	S1 or S2	13278	

	S8	S6 and S7		715
Multi-database Search @ Ovid				
Embase 1980 to 2012 Week 25, Ovid MEDLINE(R) 1946 to June Week 3 2012, Ovid MEDLINE(R) Daily Update June 28, 2012, Ovid MEDLINE(R) In-Process & Other Non-Indexed Citations June 28, 2012	1	exp *Lung Neoplasms/di, ep, et Embase <1980 to 2012 Week 25> (38891) Ovid MEDLINE(R) <1946 to June Week 3 2012> (17904) Ovid MEDLINE(R) Daily Update <June 28, 2012> (12) Ovid MEDLINE(R) In-Process & Other Non-Indexed Citations <June 28, 2012> (0)	56807	
	2	exp *Lung tumor/di, ep, et Embase <1980 to 2012 Week 25> (38891)	38891	
	3	exp *Diagnosis, Differential/ Embase <1980 to 2012 Week 25> (9979) Ovid MEDLINE(R) <1946 to June Week 3 2012> (8001) Ovid MEDLINE(R) Daily Update <June 28, 2012> (1) Ovid MEDLINE(R) In-Process & Other Non-	17981	

Appendix 3

		Indexed Citations <June 28, 2012> (0)		
	4	*cancer diagnosis/ Embase <1980 to 2012 Week 25> (8659) Ovid MEDLINE(R) <1946 to June Week 3 2012> (0) Ovid MEDLINE(R) Daily Update <June 28, 2012> (0) Ovid MEDLINE(R) In-Process & Other Non-Indexed Citations <June 28, 2012> (0)	8659	
	5	sympto*.ti,ab. Embase <1980 to 2012 Week 25> (838068) Ovid MEDLINE(R) <1946 to June Week 3 2012> (625199) Ovid MEDLINE(R) Daily Update <June 28, 2012> (538) Ovid MEDLINE(R) In-Process & Other Non-Indexed Citations <June 28, 2012> (32530)	1496335	
	6	#3 or #4 or #5 Embase <1980 to 2012 Week 25> (855668) Ovid MEDLINE(R) <1946 to June Week 3 2012> (632995)	1521732	

		Ovid MEDLINE(R) Daily Update <June 28, 2012> (539) Ovid MEDLINE(R) In-Process & Other Non-Indexed Citations <June 28, 2012> (32530)		
	7	#1 or #2 Embase <1980 to 2012 Week 25> (38891) Ovid MEDLINE(R) <1946 to June Week 3 2012> (17904) Ovid MEDLINE(R) Daily Update <June 28, 2012> (12)	56807	
	8	#6 and #7 Embase <1980 to 2012 Week 25> (2076) Ovid MEDLINE(R) <1946 to June Week 3 2012> (839) Ovid MEDLINE(R) Daily Update <June 28, 2012> (0)	2915	
	9	remove duplicates from 8 Embase <1980 to 2012 Week 25> (1593) Ovid MEDLINE(R) <1946 to June Week 3 2012> (822) Ovid MEDLINE(R) Daily Update <June 28, 2012> (0)		2424

Appendix 3

Total RESULTS from databases				6037
------------------------------	--	--	--	-------------

Appendix 4 MRC Respiratory questionnaire

MRC RESPIRATORY QUESTIONNAIRE

Respiratory Questionnaire

Questionnaire based on the MRC (UK) Respiratory Questionnaire 1986, which has been extensively validated. This questionnaire is intended to be completed by an interviewer rather than by the patient. Additional questions have been added to cover clinical aspects of bronchial hyperresponsiveness validated by the Department of Occupational and Environmental Medicine, National Lung Institute¹ The British Occupational Health Research Foundation (BOHRF)² concluded that in the clinical setting questionnaires that identify symptoms of wheeze and/or shortness of breath which improve on days away from work or on holidays have a high sensitivity, but relatively low specificity for occupational asthma.

Preamble

I am going to ask some questions, mainly about your chest. I would like you to answer **Yes** or **No** whenever possible.

If the subject is disabled from walking from any condition other than heart and lung disease, please begin questionnaire at **Question 5** and mark the adjacent box ☐

Breathlessness and Wheezing

During the last month:

1. Are you troubled by shortness of breath when hurrying on level ground or walking up a slight hill? Yes ☐ No ☐
 If Yes to 1:
2. Do you get short of breath walking with other people of your age on level ground? Yes ☐ No ☐
 If Yes to 2:
3. Do you have to stop for breath when walking at your own pace on level ground? Yes ☐ No ☐
4. If you run, or climb stairs fast do you ever
 - a. cough? Yes ☐ No ☐
 - b. wheeze? Yes ☐ No ☐
 - c. get tight in the chest? Yes ☐ No ☐
5. Is your sleep ever broken
 - a. by wheeze? Yes ☐ No ☐

Appendix 4

6. Do you ever wake up in the morning (or from your sleep if a shift worker)

a. with wheeze?	Yes <input type="checkbox"/>	No <input type="checkbox"/>
b. difficulty with breathing?	Yes <input type="checkbox"/>	No <input type="checkbox"/>

7. Do you ever wheeze

a. if you are in a smoky room?	Yes <input type="checkbox"/>	No <input type="checkbox"/>
b. if you are in a very dusty place?	Yes <input type="checkbox"/>	No <input type="checkbox"/>

If Yes to either Q5, Q6, Q7:

8. Are your symptoms better

a. at weekends (or equivalent if shift worker)?	Yes <input type="checkbox"/>	No <input type="checkbox"/>
b. when you are on holidays?	Yes <input type="checkbox"/>	No <input type="checkbox"/>

If Yes to Question 8 please record details of any occupational exposure to respiratory hazards eg isocyanates, wood dust, aluminium pot room.

Cough

- | | | |
|--|------------------------------|-----------------------------|
| 9. Do you usually cough first thing in the morning in winter? | Yes <input type="checkbox"/> | No <input type="checkbox"/> |
| 10. Do you usually cough during the day – or at night – in the winter? | Yes <input type="checkbox"/> | No <input type="checkbox"/> |

If Yes to Q9. or Q10. :

- | | | |
|--|------------------------------|-----------------------------|
| 11. Do you cough like this on most days for as much as three months each year? | Yes <input type="checkbox"/> | No <input type="checkbox"/> |
|--|------------------------------|-----------------------------|

Phlegm

- | | | |
|--|------------------------------|-----------------------------|
| 12. Do you usually bring up phlegm from your chest first thing in the morning in winter? | Yes <input type="checkbox"/> | No <input type="checkbox"/> |
| 13. Do you usually bring up any phlegm from your chest during the day – or at night – in winter? | Yes <input type="checkbox"/> | No <input type="checkbox"/> |

If Yes to Q12. or Q13 :

- | | | |
|--|------------------------------|-----------------------------|
| 14. Do you bring up phlegm like this on most days for as much as three months each year? | Yes <input type="checkbox"/> | No <input type="checkbox"/> |
|--|------------------------------|-----------------------------|

Periods of cough and phlegm

15. In the past three years, have you had a period of (increased) cough and phlegm lasting for three weeks or more? Yes ☐ No ☐

If Yes to Q15 :

- Q16. Have you had more than one such episode? Yes ☐ No ☐

Chest Illnesses

17. During the past three years, have you had any chest illness that has kept you from your usual activities for as much as a week? Yes ☐ No ☐

If Yes to Q17. :

18. Did you bring up more phlegm than usual in any of these illnesses? Yes ☐ No ☐

If Yes to Q18 :

19. Have you had more than one illness like this in the past three years? Yes ☐ No ☐

Past Illnesses

20. Have you ever had, or been told that you have had:

- a. An injury, or operation affecting your chest? Yes ☐ No ☐
- b. Heart trouble? Yes ☐ No ☐
- c. Bronchitis? Yes ☐ No ☐
- d. Pneumonia Yes ☐ No ☐
- e. Pleurisy? Yes ☐ No ☐
- f. Asthma? Yes ☐ No ☐
- g. Other chest trouble? Yes ☐ No ☐
- h. Hay fever? Yes ☐ No ☐

Tobacco Smoking

21. Do you smoke? Yes ☐ No ☐

If No to Q21 :

- Q22. Have you ever smoked as much as one Cigarette a day for as long as one year? Yes ☐ No ☐

Appendix 4

If **No** to **Question 21 or 22**, omit remaining questions on smoking.

23. How old were you when you started smoking regularly? _____

- 24a. Do (did) you smoke manufactured cigarettes? Yes ☐ No ☐

If **Yes** to Q24a :

How many do you (did) you usually smoke per day? _____

Q24b. on weekdays? _____

Q24c. at weekends? _____

25. Do you smoke any other forms of tobacco? Yes ☐ No ☐

If **Yes** to Q25 :

Record details under **Additional Notes**

For ex-smokers

- Q26. When did you give up smoking altogether? Month _____ Year _____

Additional Notes

Appendix 5 IPCARD Questionnaire



Identifying Symptoms that Predict Chest and Respiratory Disease (IPCARD) Questionnaire

First Name _____ Surname _____ Date of Birth _____

Address _____

 _____ Postcode _____

Preferred Telephone No. _____ Email address _____
 (Home or mobile) _____

Date questionnaire issued: _____

IPCARD Study Consent Form

Thank you for helping with this research. Please indicate that you have read and understood the following statements by initialling in the box and signing the declaration below.

		Please initial each section.
1.	I confirm that I have read and understood the Patient Information Sheet version 2 dated 03/10/2012 for the above study and have had the opportunity to ask questions.	
2.	I consent to follow-up of my health status through my medical records and central registries. I understand that sections of any of my medical notes may be looked at by responsible individuals from the research team at the University, or by regulatory authorities, for the purposes of this research and its management.	
3.	We would like to talk to some people about their experience of completing this form. You will be provided with more information before you finally decide whether or not to be interviewed. Not all of those who offer to help will be invited to interview. I understand that Interviews will be audio-recorded. I agree to a researcher contacting me by phone or email to discuss being interviewed.	

Name of participant (please print)

Today's Date

Signature

We are interested in all aspects of your health, including:

- Aspects of your health unrelated to your appointment today
- Your everyday health
- Any changes in your health
- Changes in your health that have not resulted in you feeling ill


Please answer questions fully even if the question does not appear relevant to your current health complaints.

There may be questions that do not apply to you. If this is the case you will be asked to skip to the next section or question - you will not need to complete these sections of the form.

Please answer the questions by shading the circle (like this ●) for the relevant option.


Section 1 - Chest and upper body aches, pain or discomfort

Q1 Have you ever experienced any discomfort in your chest, upper body or shoulders?

No ☐ Please go to page 4, Section 2 

Yes and I still have the pain/discomfort ☐ Please go to question 3 (Q3)

Yes but I no longer have the pain/discomfort ☐ → Q2 Have you had pain/discomfort in the last three months?

Yes ☐ Please go to question 3 (Q3) No ☐  Please go to Page 5, Section 2

Q3 Please indicate whether the statements below accurately describe chest or upper body aches, pains or discomfort you have experienced currently or within the last 3 months by marking yes or no for each statement.

		Yes	No
a)	A niggle, pain or ache that feels like wind or indigestion but not associated with eating	<input type="radio"/>	<input type="radio"/>
b)	Discomfort or pain when laying/sitting in a particular position	<input type="radio"/>	<input type="radio"/>
c)	Discomfort or pain that feels like bruising	<input type="radio"/>	<input type="radio"/>
d)	Discomfort or pain that is not brought on by physical activity	<input type="radio"/>	<input type="radio"/>
e)	Discomfort or pain that comes and goes	<input type="radio"/>	<input type="radio"/>
f)	Discomfort or pain that feels like a muscle "pulled"	<input type="radio"/>	<input type="radio"/>
g)	Ache or pain in centre of chest or ribs	<input type="radio"/>	<input type="radio"/>
h)	Ache or pain in the side of chest or ribs	<input type="radio"/>	<input type="radio"/>

- i) Pain started in shoulder blade Yes ☐ No ☐
- j) Pain moved round from back to front of chest ☐ ☐



Q4 Please mark where the centre of your pain is (or pains are) with an 'X' on the images above.

Q5 Please indicate whether the aches, pain or discomfort described in questions 1 to 4 also occurred 4-12 months ago

Yes

☐

No

☐

Q6 Please indicate whether the aches, pain or discomfort described in questions 1 to 4 also occurred more than 12 months ago

Yes

☐

No

☐

Q7 In general are your aches, pain or discomfort worse than they were 3 months ago?

Yes

☐ Please go to question 9(Q9)

No

☐ Please go to question 8(Q8)

Q8 In general are your aches, pain or discomfort worse than they were 12 months ago? Yes ☐ No ☐

Q9 Please mark one number on the scale to indicate how much discomfort or distress the pain caused when at its worst

(0)

(1)

(2)

(3)

(4)

(5)

(6)

(7)

(8)

(9)

No discomfort/distress

Much discomfort/distress

Q10 Please mark one number on the scale to indicate how much the chest pain interfered with everyday life and activities when at its worst

(0)

(1)

(2)

(3)

(4)

(5)

(6)

(7)

(8)


(9)

Not at all

Very much indeed

Section 2 – Cough

Q11 Have you ever had a cough that lasted for more than 3 weeks?

No ☐ Please go to page 6, Section 3 

Yes and I still have the cough ☐ Please go to question 13 (Q13)

Yes but I no longer have the cough ☐ → Q12 Have you had a cough in the last three months Yes ☐ No ☐

Please go to question 13 (Q13)

Q13 Please indicate when you first had a cough that lasted for more than 3 weeks.

Within the last 3 months

☐

4-12 months ago

☐

More than 12 months ago

☐

Q14 Please indicate whether the statements below accurately describe your most recent cough/coughs (that lasted for more than 3 weeks) and how often you have had the type of cough described by that statement.

		Never	Once	Occasionally	Most of the time
a)	An irritating cough (feels like an irritation in the throat or chest)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
b)	A tickly cough	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
c)	A cough that starts in the throat	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
d)	A cough that feels like clearing the throat	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
e)	A wheezy cough	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
f)	Cough that feels as though it arises in one or other lung or side of the chest	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
g)	Cough that interrupts speaking	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
h)	A cough without phlegm	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
i)	A cough that usually produces phlegm in the morning	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
j)	A cough that produces phlegm at any time of the day	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

- | | | | | | |
|----|--|-----------------------|-----------------------|-----------------------|-----------------------|
| k) | A hard or harsh cough without phlegm | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> |
| l) | A hard or harsh cough that produces phlegm | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> |

Q15 Please indicate whether any of the descriptions below accurately describe a cough that you have had within the last six months (which has lasted for more than 3 weeks).

- | | Yes | No |
|----------------------------------|-----------------------|-----------------------|
| a) Cough comes and goes | <input type="radio"/> | <input type="radio"/> |
| b) Cough affected by the weather | <input type="radio"/> | <input type="radio"/> |
| c) A smoker's cough | <input type="radio"/> | <input type="radio"/> |

Q16 In general is your cough worse than it was 3 months ago?

- | | |
|--|--|
| Yes | No |
| <input type="radio"/> Please go to question 18 (Q18) | <input type="radio"/> Please go to question 17 (Q17) |

Q17 In general is your cough worse than it was 12 months ago? Yes ☐ No ☐

Q18 Please describe any changes in your cough over time here:

Q19 Please mark **one number on the scale** to indicate how much discomfort or distress the coughing caused when at its worst

☐ 0 ☐ 1 ☐ 2 ☐ 3 ☐ 4 ☐ 5 ☐ 6 ☐ 7 ☐ 8 ☐ 9

No discomfort/distress Much discomfort/distress

Q20 Please mark **one number on the scale** to indicate how much coughing interfered with everyday life and activities when at its worst


☐ 0 ☐ 1 ☐ 2 ☐ 3 ☐ 4 ☐ 5 ☐ 6 ☐ 7 ☐ 8 ☐ 9

Not at all Very much indeed

Section 3 - Breathing changes

Q21 Have you ever experienced any of the following?

- becoming short of breath more easily than you used to
- unexpected shortness of breath
- noise/unusual sensation when breathing
- any difficulty breathing

No ☐ Please go to page 8, Section 4 

Yes and I still have breathing difficulties/changes ☐ Please go to question 23 (Q23)

Yes but I no longer have breathing difficulties/changes ☐ → Q22 Have you had breathing difficulties/changes in the last three months? Yes ☐ No ☐

Please go to question 23 (Q23)

Q23 Please indicate when you first had breathing difficulties/changes

Within the last 3 months	4-12 months ago	More than 12 months ago
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Q24 Please indicate whether the statements below accurately describe your breathing difficulties and how often you have had the type of difficulties described by that statement within the last 12 months.

		Never	Once	Occasionally	Most of the time
a)	Breathlessness after walking a short distance on the flat	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
b)	Breathlessness that comes on unexpectedly	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
c)	Breathlessness when lying down	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
d)	Breathlessness on resting	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
e)	Breathing problems that require the use of an inhaler	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
f)	Breathlessness that feels like an anxiety attack or that is associated with a feeling of anxiousness	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Q25 Have you experienced breathing problems that are only present or get worse at certain times of the year? Yes ☐ No ☐

Q26 Please indicate whether any of the statements below accurately describe breathlessness that you have experienced within the last 6 months by marking yes or no for each statement:

	Yes	No
a) Feeling out of breath	<input type="radio"/>	<input type="radio"/>
b) Unable to get enough air	<input type="radio"/>	<input type="radio"/>
c) Tightness in chest	<input type="radio"/>	<input type="radio"/>
d) Breathing is shallow	<input type="radio"/>	<input type="radio"/>
e) Breathing is rapid	<input type="radio"/>	<input type="radio"/>
f) Feels like a weight on your chest	<input type="radio"/>	<input type="radio"/>

Q27 Please indicate whether any of the statements below accurately describe your breathing within the last six months by marking yes or no for each statement:

	Yes	No
a) Strange sensation felt in lung when breathing	<input type="radio"/>	<input type="radio"/>
b) Wheezing noise when breathing in (for more than 2 weeks)	<input type="radio"/>	<input type="radio"/>
c) Wheezing noise when breathing out (for more than 2 weeks)	<input type="radio"/>	<input type="radio"/>
d) Wheezing sensation when in a particular position	<input type="radio"/>	<input type="radio"/>

Q28 In general is your breathlessness worse than it was 3 months ago?

Yes	No
<input type="radio"/> Please go to question 30 (Q30)	<input type="radio"/> Please go to question 29 (Q29)

Q29 In general is your breathlessness worse than it was 12 months ago? Yes ☐ No ☐

Q30 Please mark **one number on the scale** to indicate how much discomfort or distress the change in breathing or breathlessness caused when at its worst.


(0)	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
No discomfort/distress					Much discomfort/distress				

Q31 Please mark **one number on the scale** to indicate how much the change in breathing or breathlessness interfered with everyday life and activities when at its worst.

(0)	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Not at all					Very much indeed				

Section 4 - Tiredness

Q32 Have you experienced any unexpected tiredness within the last 12 months?

No ☐ Please go to page 9, Section 5 

Yes and I still have unexpected tiredness ☐ Please go to question 34 (Q34)

Yes but I no longer have unexpected tiredness ☐ → Q33 Have you had unexpected tiredness in the last three months? Yes ☐ No ☐

Please go to question 34 (Q34)

Q34 Please indicate when you **first** experienced any unexpected tiredness

Within the last 3 months ☐ 4-12 months ago ☐ More than 12 months ago ☐

Please answer the questions below that are all about unexpected tiredness **within the last 12 months**

	Never	Once	Occasionally	Most of the time
Q35 Have you felt tired more easily than you used to?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Q36 Have you felt as though you needed to sleep during the day?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Q37 Have you felt like you wanted to sit down or stop activity?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Q38 In general is your tiredness worse than it was 3 months ago?

Yes ☐ Please go to question 40 (Q40) No ☐ Please go to question 39 (Q39)

Q39 In general is your tiredness worse than it was 12 months ago? Yes ☐ No ☐

Q40 Please mark **one number on the scale** to indicate how much discomfort or distress the tiredness caused when at its worst

(0) (1) (2) (3) (4) (5) (6) (7) (8) (9)

No discomfort/distress Much discomfort/distress


Q41 Please mark **one number on the scale** to indicate how much the tiredness interfered with everyday life and activities when at its worst.

(0) (1) (2) (3) (4) (5) (6) (7) (8) (9)

Not at all Very much indeed

Section 5 - Coughing up Blood

Q42 Have you ever coughed up any blood?

No ☐ Please go to Section 6 (below), on this page 

Yes and I am still coughing up blood ☐ Please go to question 44 (Q44)

Yes but I am no longer coughing up blood ☐ → Q43 Have you coughed up blood in the last three months? Yes ☐ No ☐

Please go to question 44 (Q44)

Q44 Please indicate when you first coughed up any blood

Within the last 3 months ☐ 4-12 months ago ☐ More than 12 months ago ☐

No ☐ Once ☐ Occasionally ☐ Most of the time ☐

Q45 Have you ever coughed up mostly blood (blood with little or no phlegm)? ☐ ☐ ☐ ☐

Q46 Have you ever coughed up phlegm with small amounts of blood? ☐ ☐ ☐ ☐

Section 6 - Chest and respiratory infections and colds

Yes No

Q47 Have you currently got a phlegmy chest or chest infection? ☐ ☐

Q48 Have you currently got a cold, flu or any other type of infection that has caused a cough or affected your breathing? ☐ ☐

Q49 How many times have you had a chest infection within the last 12 months?

0 ☐ 1 ☐ 2-3 ☐ More than 3 ☐

Q50 Have you had noticeably more chest infections within the last 12 months than in the year before (13-24 months ago)? Yes ☐ No ☐

Q51 How many times have you had an infection that has caused a cough or affected your breathing, a cold or flu within the last 12 months?

0 ☐ 1 ☐ 2-3 ☐ More than 3 ☐


Q52 Have you had noticeably more colds or flu within the last 12 months than in the year before (13-24 months ago)? Yes ☐ No ☐

Section 7 - Changes in Weight

	Yes	No
Q53 Do you have to eat more than you used to in order to maintain a steady weight?	<input type="radio"/>	<input type="radio"/>
Q54 Do you now weigh less than you have for most of your adult life?	<input type="radio"/>	<input type="radio"/>
Q55 Within the last 12 months have you unintentionally lost weight that you have not regained?	<input type="radio"/>	<input type="radio"/>
Q56 Have you gained weight within the last 12 months?	<input type="radio"/>	<input type="radio"/>

Section 8 - Hot or Cold Sweats

Q57 Have you experienced hot or cold sweats during the night or day within the last two years?

No ☐ Please go to Section 9 (below), on this page 

Yes and I am still experiencing hot or cold sweats ☐ Please go to question 59 (Q59)

Yes but I am no longer experiencing hot or cold sweats ☐ → Q58 Have you had hot or cold sweats in the last three months? Yes ☐ No ☐

Please go to question 59 (Q59)

Q59 Please indicate when you first experienced hot or cold sweats during the night or day

Within the last 3 months	4-12 months ago	More than 12 months ago
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	Never	Once
		Occasionally
		Most of the time

Q60 Have you experienced hot or cold sweats in the night? ☐ ☐ ☐ ☐

Q61 Have you experienced hot or cold sweats in the day? ☐ ☐ ☐ ☐

Q62 Do you think all of your hot or cold sweats are probably caused by the menopause? Yes ☐ No ☐ Not sure ☐

Section 9 - Eating Changes

	Yes	No
Q63 Has your appetite increased within the last 12 months?	<input type="radio"/>	<input type="radio"/>
Q64 Has your appetite decreased within the last 12 months?	<input type="radio"/>	<input type="radio"/>
Q65 Have you experienced any taste changes within the last 2 years?	<input type="radio"/>	<input type="radio"/>

Yes No

Q66 Have you currently gone off certain foods you used to eat?

☐ ☐**Section 10 - Arms, Legs and Joints**

No Within the last 3 months 4-12 months ago 1-2 years ago

Q67 Have you experienced a **new** aching sensation in any joints or any new joint pain in the last 2 years?☐ ☐ ☐ ☐Q68 Have you experienced any **new** unusual sensations or tingling in your arms or legs in the last 2 years?☐ ☐ ☐ ☐**Section 11 - Voice Changes**Q69 Have you experienced any **ongoing** changes in the sound of your voice when speaking?

Yes

No

☐ Please go to question 70 (Q70)☐ Please go to Section 12 (below), on this page

Q70 Please indicate when you first experienced changes in the sound of your voice when speaking

Within the last 3 months

4-12 months ago

More than 12 months ago

☐☐☐**Section 12 - Skin Changes**

Q71 Have you experienced any changes in the condition of your skin in the last two years?

Yes

No

☐ Please go to question 72 (Q72)☐ Please go to page 12, Section 13

Q72 Please describe any changes in the condition of your skin.

Section 13 - Any other illnesses

Q73 Have you been told by a doctor that you had the illnesses listed below? (Please mark one response for each condition):

	Never	Within the last 3 months	4-12 months ago	1-5 years ago	More than 5 years ago
a) Pneumonia	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
b) Bronchitis or chronic bronchitis	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
c) Asthma	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
d) Seasonal allergy (e.g. hay fever or seasonal breathing problems).	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
e) Chronic Obstructive Pulmonary Disease (COPD)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
f) Heart Disease or Angina	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
g) Anaemia	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
h) Cancer	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
i) Emphysema or pulmonary fibrosis	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
j) An asbestos related illness	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
k) Arthritis	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Q74 Please describe any other serious illnesses you have had within the last 2 years.

Q75 Have any of your blood relatives (brothers, sisters, parents or children) had the illnesses listed below? (Please mark all that apply):

	Yes	No
a) Asthma	<input type="radio"/>	<input type="radio"/>
b) Bronchitis	<input type="radio"/>	<input type="radio"/>
c) Heart disease or Angina	<input type="radio"/>	<input type="radio"/>
d) Lung cancer	<input type="radio"/>	<input type="radio"/>
e) Tuberculosis (TB)	<input type="radio"/>	<input type="radio"/>

Section 14 - Other changes in your health

Q76 Please describe any changes in your health or anything unusual or different about your body and health, you have noticed during the last 2 years.

Q77 Please describe any medication, or illnesses not mentioned above, which might have caused your symptoms:

Section 15 - Smoking History

Q78 Have you ever smoked?

(Smoking is defined as smoking one cigarette/pipe/cigar a day for as long as one year.)

Yes

☐

Please go to
question 79 (Q79)

No

☐

Please go to
The end of the questionnaire
(below)

Q79 When did you start smoking?

Year _____

Q80 What is the total number of years in your life that you have smoked?

Q81 Do you currently smoke?

Yes ☐

Q82 On average (over your lifetime), how much do you smoke on a week day?

_____ cigarettes

_____ cigars

_____ oz tobacco

No ☐

Q83 When did you give up smoking altogether?

Year _____

Q84 On average, how much did you used to smoke on a week day?

_____ cigarettes

_____ cigars

_____ oz tobacco

Thank you very much for your time.

Please return this questionnaire to: Faculty of Health Sciences, Building 67, The University of Southampton, Southampton, SO17 1BJ in the prepaid envelope provided.

If you have any concerns or questions about the study or issues raised by this study please contact:
Dr Lucy Brindle, Chief Investigator, University of Southampton, 02380 598526.



ID Number

Appendix 6 Introductory Letter



Study title: IPCARD (Identifying Symptoms that Predict Chest and Respiratory Disease) Chest Clinic Study

Date:

Dear Sir/Madam,

You are being invited to take part in a research study. Before you decide, it is important that you understand why the study is being carried out and what it will involve. This letter and the following information sheet will give you details that you might like to discuss with your family, friends or hospital staff before you make any decision. Take time to decide whether or not you wish to take part. Please ask if anything is unclear to you or if you would like to know more. You can contact a researcher directly or you can ask a member of staff to contact the research team. Our contact details are at the bottom of this page.

Purpose of the study:

The purpose of this study is to identify symptoms which might improve success in the early diagnosis of chest disease in the future. To do this we will use a questionnaire which records a lot of information about the health of those attending chest clinics. We want to find out which bits of the information are useful in identifying chest diseases. This study will take place over a period of eighteen months in the Faculty of Health Sciences, at the University of Southampton.

Thank you for taking the time to read this and to consider taking part in this research.

Yours Sincerely,

Joanna Shim
(PhD Researcher)
Address: Building 45, Highfield Campus

University of Southampton

SO17 1BJ

Lead Researcher: Dr Lucy Brindle
Address: Faculty of Health Sciences

Building 67

University of Southampton

SO17 1BJ

Appendix 6

Office: 02380 522360

Fax: +44(0)2380 594792

Email: js1g08@soton.ac.uk

Office: 02380 598526

Email: L.A.Brindle@soton.ac.uk

Appendix 7 Participant Information Sheet (A)



PARTICIPANT INFORMATION SHEET (A)

IPCARD (Identifying Symptoms that Predict Chest and Respiratory Disease) Chest Clinic Study

You are being invited to take part in a research study. Before you decide whether or not to take part, it is important for you to understand why the research is being done and what it will involve. Please take your time to read the information carefully.

What is the purpose of the study?

We want to gain a better understanding of the health of those who have been referred to chest clinics to help us to identify the symptoms of early chest disease. In order to do so, we will be using a questionnaire that records symptoms. This study aims to find out if it is possible to identify respiratory diseases from these symptoms. We are interested in all aspects of your health, including aspects of your health unrelated to your clinic visit today. Hence, we will also be asking questions about a range of conditions such as anaemia, angina (heart disease), skin and joint problems and cancer. The questionnaire records a lot of information about your health and we want to find out which information is useful in identifying, or ruling out, chest diseases such as asthma, chronic obstructive pulmonary disease and lung cancer. We hope to use the results to identify as early as possible patients who have or will develop serious chest disease, so that in the future, patients can be offered treatment at a stage when cure is more likely. To do this, we need the help of a large number of patients, some of whom have chest diseases, some of whom have cancers, and some of whom have no major illness at all.

Why have I been invited to participate?

We are inviting those aged 40 or over who are attending this clinic today, and have not yet had a CT scan result, to participate.

Do I have to take part?

No. It is entirely up to you to choose whether or not to take part and your decision will not affect your medical care or treatment in any way. If you decide not to take part, you do not have to give a reason. If you decide to take part you are still free to withdraw at any time and without giving a reason.

What will happen to me if I take part?

If you agree to take part, you will be asked to complete the questionnaire and either leave it at your chest clinic or return it by post in the pre-paid envelope. We expect that on average, the questionnaire will take twenty minutes to complete. If after returning the questionnaire you decide that you no longer wish to take part, please contact the Chief Investigator by letter, email or phone. Her contact details are on page 2 of this information sheet.

If you decide to take part and complete the questionnaire, we may contact you again to ask if you are willing to take part in an interview. The researcher would talk with you for about an hour, asking you about your health and about your answers to the questionnaire. However, most of those who fill in the questionnaire will not be contacted again.

Appendix 7

We will also ask those who complete the questionnaire whether we can then look at their medical records. This allows us to find out about your diagnosis. If you have given us permission, we will look at your medical records and check for any changes in your health.

What are the possible disadvantages and risks of taking part?

The study does not involve any treatment or tests, so there is no physical risk involved.

What are the possible benefits of taking part?

This research will not directly benefit you, but what you tell us may help future patients with chest disease.

Will the information about me be kept confidential?

We will follow ethical and legal practice and all information that is collected about you during the course of the research will be kept strictly confidential. In addition, any information about you which leaves the hospital will have your name and address removed so that you cannot be recognised. Things that you say in your interview may be quoted in reports of study findings, however, names and personal details will be removed to protect your identity. The audio-recordings and your data will be securely stored at Southampton University for 10 years in a password protected system.

What if I have any questions, concerns or want to complain?

If you have any questions about any aspect of this study, or would like more information before you make up your mind whether to participate, you can contact the researcher, Joanna Shim (PhD student) at the Faculty of Health Sciences, University of Southampton by calling 02380 522360

If you have any further questions or concerns please contact the study's Chief Investigator, Dr Lucy Brindle, (Faculty of Health Sciences, Building 67, University Road, Southampton, SO17 1BJ; Telephone 02380 598526; Email L.A.Brindle@soton.ac.uk). If you remain unhappy, please contact Martina Prude (Head of Research Governance) at the University of Southampton (Telephone 02380 598848).

The normal NHS complaints mechanisms are also available to you through the Patient Advisory and Liaison Service (PALS): Southampton University Hospital Trust on 023 80798498; or Salisbury/ Winchester (to add)

Who is organising and funding the research?

This research, which aims to identify early symptoms of chest and respiratory disease, is funded by the Faculty of Health Sciences, Southampton University and is being carried out by a group of researchers at Southampton University. The researcher in the clinic today might be a PhD student, and the data collected will also form part of a PhD.

Who has reviewed the study?

The research would have been looked at by an independent group of people called a Research Ethics Committee, to protect your safety, rights, well-being and dignity. This study has been reviewed and given a favourable opinion by Berkshire NHS Ethics Committee (REC No: 12/SC/0490).

What should I do if I want to take part?

If you are happy to take part in the study, kindly please fill in the questionnaire, seal it in the pre-paid envelope provided, and either leave it at the chest clinic or return it by post within 24-48 hours.

Thank you for taking time to consider participating in this research.

Appendix 8 Participant Information Sheet (B)



PARTICIPANT INFORMATION SHEET (B)

IPCARD (Identifying Symptoms that Predict Chest and Respiratory Disease) Chest Clinic Study

You are being invited to take part in a research study. Before you decide whether or not to take part, it is important for you to understand why the research is being done and what it will involve. Please take your time to read the information carefully.

What is the purpose of the study?

We want to gain a better understanding of the health of those who have been referred to chest clinics to help us identify the symptoms of early chest disease. In order to do so, we will be using a questionnaire that records symptoms. This study aims to find out if it is possible to identify respiratory diseases from these symptoms. We are interested in all aspects of your health, including aspects of your health unrelated to your clinic visit today. Hence, we will also be asking questions about a range of conditions such as anaemia, angina (heart disease), skin and joint problems and cancer. The questionnaire records a lot of information about your health and we want to find out which information is useful in identifying, or ruling out, chest diseases such as asthma, chronic obstructive pulmonary disease and lung cancer. We hope to use the results to identify as early as possible patients who have or will develop serious chest disease, so that in the future, patients can be offered treatment at a stage when cure is more likely. To do this, we need the help of a large number of patients, some of whom have chest diseases, some of whom have cancers, and some of whom have no major illness at all.

Why have I been invited to participate?

We are inviting those aged 40 or over who are attending this clinic today, and have not yet had a CT scan result, to participate.

Do I have to take part?

No. It is entirely up to you to choose whether or not to take part and your decision will not affect your medical care or treatment in any way. If you decide not to, you do not have to give a reason. If you decide to take part you are still free to withdraw at any time and without giving a reason.

What will happen to me if I take part?

If you agree to take part, you will be asked to complete the questionnaire and either leave it at your chest clinic or return it by post in the pre-paid envelope. We expect that on average, the questionnaire will take twenty minutes to complete. If after returning the questionnaire you decide that you no longer wish to take part, please contact the Chief Investigator by letter, email or phone. Her contact details are on page 2 of this information sheet.

We will also ask those who complete the questionnaire whether we can then look at their medical records. This allows us to find out about your diagnosis. If you have given us permission, we will look at your medical records and check for any changes in your health.

Appendix 8

What are the possible disadvantages and risks of taking part?

The study does not involve any treatment or tests, so there is no physical risk involved.

What are the possible benefits of taking part?

This research will not directly benefit you, but what you tell us may help future patients with chest disease.

Will the information about me be kept confidential?

We will follow ethical and legal practice and all information that is collected about you during the course of the research will be kept strictly confidential. In addition, any information about you which leaves the hospital will have your name and address removed so that you cannot be recognised. Things that you say in your interview may be quoted in reports of study findings, however, names and personal details will be removed to protect your identity. Your data will be securely stored at Southampton University for 10 years in a password protected system.

What if I have any questions, concerns or want to complain?

If you have any questions about any aspect of this study, or would like more information before you make up your mind whether to participate, you can contact the researcher, Joanna Shim (Researcher) at the Faculty of Health Sciences, University of Southampton by calling 02380 522360.

If you have any further questions or concerns please contact the study's Chief Investigator, Dr Lucy Brindle, (Faculty of Health Sciences, Building 67, University Road, Southampton, SO17 1BJ; Tel: 02380 598526; Email L.A.Brindle@soton.ac.uk) . If you remain unhappy, please contact Martina Prude (Head of Research Governance) at the University of Southampton (Tel: 02380 598848).

The normal NHS complaints mechanisms are also available to you through the Patient Advisory and Liaison Service (PALS): Southampton University Hospital Trust on 023 80798498
Salisbury (to add)
Winchester (to add)

Who is organising and funding the research?

This research, which aims to identify early symptoms of chest and respiratory disease, is funded by the Faculty of Health Sciences, Southampton University, and is being carried out by a group of researchers at Southampton University. The researcher in the clinic today might be a PhD student, and the data collected will also form part of a PhD.

Who has reviewed the study?

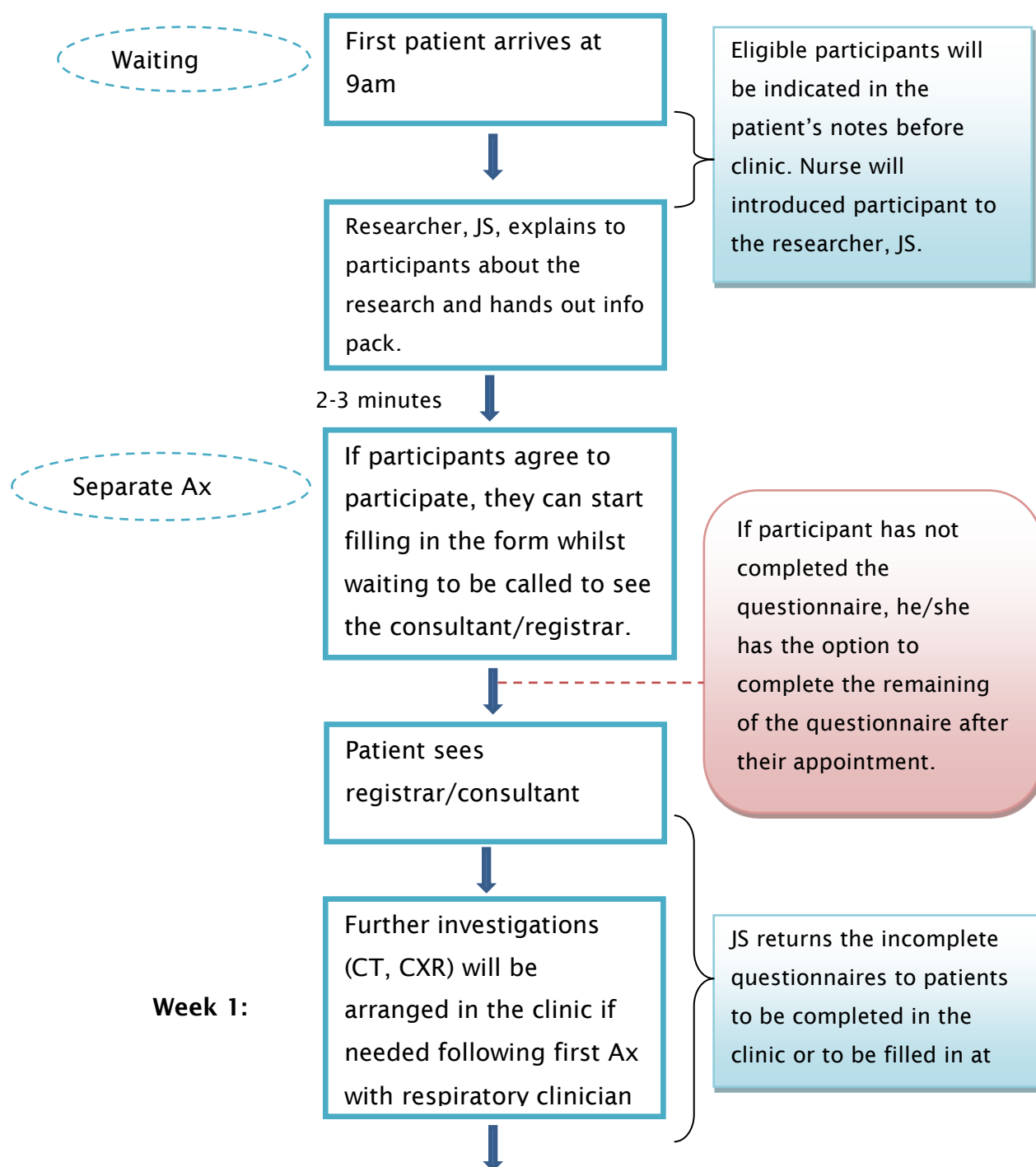
The research would have been looked at by an independent group of people called a Research Ethics Committee, to protect your safety, rights, well-being and dignity. This study has been reviewed and given a favourable opinion by Berkshire NHS Ethics Committee (REC No: **12/SC/0490**).

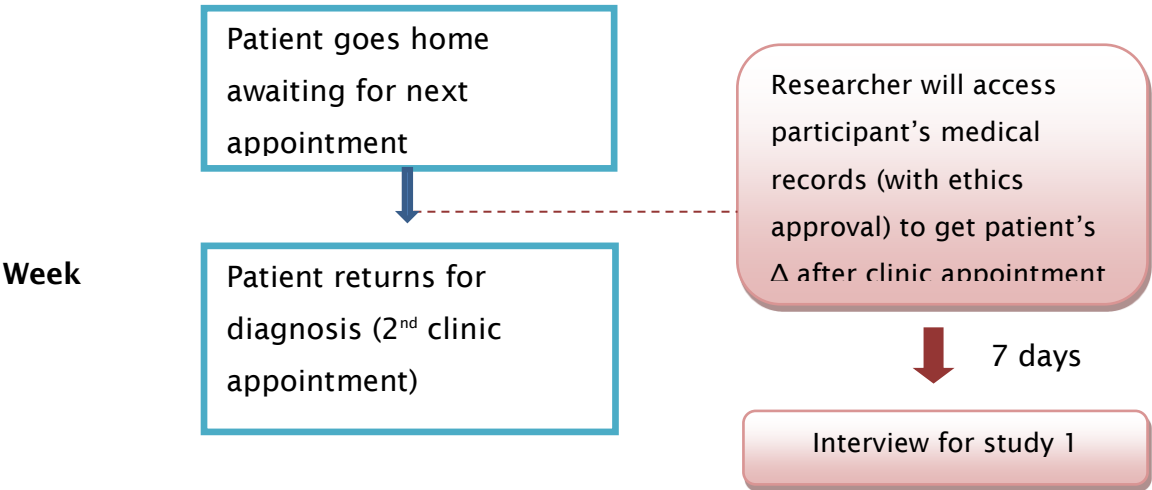
What should I do if I want to take part?

If you are happy to take part in the study, kindly please fill in the questionnaire, seal it in the pre-paid envelope provided, and either leave it at the chest clinic or return it by post within 24-48 hours.

Thank you for taking time to consider participating in this research.

Appendix 9 Schema for recruitment in Southampton clinic:





Appendix 10 Interview Schedule Outline

I. Semi-structured interview:

- a. To record patient's health and illness experiences over the last 2 years from their first symptom (health problem) leading up to their referral to lung-shadow clinic.
- b. Tips: The interview should be facilitated to focus upon the research questions. The person's own words for any health problem that they describe should be used throughout the interview.

II. Structured interview:

- a. To explore more specific symptom presentation and health changes.
- b. Dialogue: (Follow interview schedule for the main IPCARD study)

III. Cognitive interview:

To investigate the ease of interpretation of the questionnaire items, reasons for the missing data (non-completion) and participant's questionnaire responses. To explore discrepancies identified between the narrative of health and illness (Section A) and responses to section B to the responses to the questionnaire items.

Instructions to interviewer:

The interview should be sensitive to issues of subject burden. If the subject chooses not to complete the interview this should be documented, with brief reason (e.g. tiredness, pain), on the front sheet. The subject may request someone else to be present during the interview. This should be documented on the front sheet.

Section I: Mapping patient's experiences of health and illness in the last 2 years and their journey from the first problem/change in health to their referral to the lung-shadow clinic.

Instructions to Interviewer:

This section will focus on:

- a) an overview of the subject's experience of health including any symptoms/problems/changes in health that they have noticed and;
- b) the time of these events

Recording problem/symptoms/changes:

The nature of the problem/change in health (e.g. breathlessness) and the date it started should be explored. Also record with timeline (if possible) if the problem changed during the journey from noticing-now (e.g. breathlessness became worse).

Dialogue:

"Thank you for offering your time to participate in this interview. Your feedback will help us learn how to better improve our questionnaire that will help identify any chest problems. The interview should take about 60-90 minutes to complete. If at any point of the interview you wish to stop or have a break, please let me know. Also, if you want any questions repeated or clarified, please ask. I will be tape recording the interview. Do I have your permission to record this interview?"

"I will be asking you to talk about things that have happened to you from when you first noticed a change in your health up to the time when you were referred to the lung clinic and about all aspects of your health in the last 2 years."

"I would like to build up a detailed picture of your experiences of health over the last two years. I am interested in anything that you noticed about your health during this time even if you thought it was minor or not connected to your recent chest x-ray."

1. "For a start, could you tell me in your own words your health experiences over the last 2 years; how it changed (any changes that you noticed) up to now?"
2. "When did you first notice something was wrong, or a change in your health?"
3. "Could you tell me about what you noticed?"
4. "When did this happen?"

5. "Have you experienced any other changes in your health during the last 2 years?"
6. "Has there been anything else that you have visited your doctor about?"
7. "Has there been anything else at all relating to your health that you have noticed during the last two years even if minor?"

Section III: Cognitive Interview

Dialogue:

"Thank you again for taking part in the interview (Ask them if they need a break?). For this part of the interview, I would like to find out your experience filling out the questionnaire and how you feel about it. I am going to read out to you, questions from the questionnaire and I would like you to tell me what was going through your mind as you answer it (almost as if you are thinking out loud)."

Responses:	Prompts/ Probes:
If participants struggle	Could you tell me what is going through your mind?
	Tell me what you are thinking.
	You responded (xyz), how did you decide on an answer?
If participants show ease in answering	Respond with encouragement: That's great. Thinking out loud like this is just what I need.
	Good. Your comments are helping me understand what you are thinking about.

Appendix 10

PROBE Questions:

General

- Are there any parts of the questionnaire that you do not agree with or do not particularly like?
- What do you believe the question to be asking?
- What thoughts came to mind while reading (section/phrase/question)?
- What, to you, does the term/ word (x) mean?
- I noticed that you hesitated. Tell me what you were thinking?

Recall information

- How easy was it for you to recall these events in order to answer the question?
- How hard was this to answer?
- You commented on (their experience) when completing question (x). What thoughts came to mind when doing so?

Requesting clarification on timing

- How did you get to that answer of (x) period?
- How well do you remember this?
- How did you remember that you had (x) symptom for (y) months.

Appendix 11 Frequency and percentage of missing observations for all variables

Variable	No. of missing	% of missingness	No. of observed
Q1_Pain_	4	1.11	355
Q2_Pain	3	0.84	356
Q3a_Pain	18	5.01	341
Q3b_Pain	18	5.01	341
Q3c_Pain	23	6.41	336
Q3d_Pain	21	5.85	338
Q3e_Pain	18	5.01	341
Q3f_Pain	22	6.13	337
Q3g_Pain	17	4.74	342
Q3h_Pain	16	4.46	343
Q3i_Pain	25	6.96	334
Q3j_Pain	44	12.26	315
Q6_Pain	25	6.96	334
Q7_Pain	23	6.41	336
Q8_Pain_	25	6.96	334
Q9_Pain	28	7.80	331
Q10_Cgh_	12	3.34	347
Q11_Cgh	10	2.78	349
Q12_Cgh_	13	3.62	346
Q13a_Cgh	23	6.41	336
Q13b_Cgh	28	7.80	331
Q13c_Cgh	36	10.00	323
Q13d_Cgh_	31	8.64	328
Q13e_Cgh_	32	8.91	327
Q13f_Cgh_	48	13.37	311
Q13g_Cgh	32	8.91	327
Q13h_Cgh	31	8.63	328
Q13i_Cgh	24	6.69	335
Q13j_Cgh	24	6.69	335
Q13k_Cgh	42	11.70	317

Appendix 11

Q13l_Cgh	33	9.19	326
Q14a_Cgh	37	10.30	322
Q14b_Cgh_	42	11.70	317
Q14c_Cgh_	41	11.42	318
Q15_Cgh	23	6.41	336
Q16_Cgh_	23	6.41	336
Q18_Cgh	26	7.24	333
Q19_BrChnges_	19	5.29	340
Q20_BrChnges	16	4.46	343
Q21_BrChnges_	15	4.18	344
Q22a_BrChnges	19	5.29	340
Q22b_BrChnges	27	7.52	332
Q22c_BrChnges	25	6.96	334
Q22d_BrChnges	30	8.36	329
Q22e_BrChnges	27	7.52	332
Q22f_BrChnges	30	8.36	329
Q23_BrChnges	33	9.19	326
Q24a_BrChnges	19	5.29	340
Q24b_BrChnges	26	7.24	333
Q24c_BrChnges	21	5.85	338
Q24d_BrChnges	31	8.64	328
Q24e_BrChnges	29	8.08	330
Q24f_BrChnges	26	7.24	333
Q25a_BrChnges	24	6.69	335
Q25b_BrChnges	19	5.29	340
Q25c_BrChnges	25	6.96	334
Q25d_BrChnges	26	7.24	333
Q26_BrChnges	20	5.57	339
Q27_BrChnges	22	6.13	337
Q28_BrChnges	26	7.24	333
Q29_Tired_	17	4.74	342
Q30_Tired	11	3.06	348
Q31_Tired	14	3.90	345
Q32_Tired	14	3.90	345
Q33_Tired_	16	4.46	343

Q34_Tired_	19	5.29	340
Q35_Tired	15	4.18	344
Q36_Tired	18	5.01	341
Q37_Tired	22	6.13	337
Q38_CghBlood_	10	2.79	349
Q39_CghBlood	8	2.23	351
Q40_CghBlood_	12	3.34	347
Q41_CghBlood	16	4.46	343
Q42_CghBlood_	15	4.18	344
Q43_ChInfectn	17	4.74	342
Q44_ChInfectn	24	6.69	335
Q45_ChInfectn_	16	4.46	343
Q46_ChInfectn	18	5.01	341
Q47_ChInfectn_	25	6.96	334
Q48_ChInfe~n	36	10.03	323
Q49_Weight	22	6.13	337
Q50_Weight	15	4.18	344
Q51_Weight	17	4.74	342
Q52_Weight	19	5.29	340
Q53_HCswat_	14	3.90	345
Q54_HCswat	15	4.18	344
Q55_HCswat_	15	4.18	344
Q56_HCswat	12	2.34	347
Q57_HCswat	16	4.46	343
Q58_HCswat	16	4.46	343
Q59_EatChnges	16	4.46	343
Q60_EatChnges	20	5.57	339
Q61_EatChnges	21	5.85	338
Q62_EatChnges	48	13.37	311
Q63_New_JointPain_	21	5.85	338
Q64_New_JointPain	21	5.85	338
Q65_Voice	15	4.18	344
Q66_Voice	16	4.46	343
Q67_Skin	21	5.85	338
Q69a_Pneumo	22	6.13	337

Appendix 11

Q69c_Asthma_	24	6.69	335
Q69d_Allergy_	27	7.52	332
Q69e_COPD_	29	8.08	330
Q69f_HD_Angina_	25	6.96	334
Q69g_Anemia_	27	7.52	332
Q69h_Cancer_	28	7.80	331
Q69j_Asb	24	6.69	335
Q69k_Arthritis	27	7.52	332
Q71a_FamHx	34	9.47	325
Q71b_FamHx	46	12.80	313
Q71c_FamHx	38	10.58	321
Q71d_FamHx	40	11.14	319
Q71e_FamHx	44	12.26	315
Q73_Smoke	12	3.34	347

Appendix 12 Imputation Model

Generic symptom variables	Imputed variables
Pain	(Q10_Cgh_) (Q43_ChInfectn) (Q44_ChInfectn) (Q69k_Arthritis) (Q13k_Cgh) (Q63_New_JointachePain_) (Q69e_COPD_) (Q69f_HD_Angina_) (Q69h_Cancer_) (Q69j_Asbes) (Q15_Cgh)
Cough	(Q19_BrChnges) (Q43_ChInfectn) (Q44_ChInfectn) (Q69e_COPD_) (Q41_CghBlood) (Q42_CghBlood_) (Q65_Voice) (Q69a_Pneumo) (Q69c_Asthma_) (Q69d_Allergy_) (Q69j_Asbes)
Breathing change	(Q10_Cgh_) (Q29_Tired_) (Q43_ChInfectn) (Q44_ChInfectn) (Q52_Weight) (Q69a_Pneumo) (Q69c_Asthma_) (Q69d_Allergy_) (Q69e_COPD_) (Q69f_HD_Angina_) (Q69j_Asbes)
Tiredness	(Q50_Weight) (Q69f_HD_Angina_) (Q43_ChInfectn) (Q44_ChInfectn) (Q69e_COPD_)
Coughing up blood	(Q43_ChInfectn) (Q44_ChInfectn) (Q69e_COPD_)
Chest infection	(Q36_Tired) (Q10_Cgh_) (Q19_BrChnges) (Q69a_Pneumo) (Q69c_Asthma_) (Q69d_Allergy_) (Q69e_COPD_) (Q69j_Asbes)
Weight	(Q76_Smoke) (Q29_Tired_) (Q43_ChInfectn) (Q44_ChInfectn) (Q60_EatChnges) (Q61_EatChnges) (Q62_EatChnges)
Hot/Cold sweats	(Q43_ChInfectn) (Q44_ChInfectn) (Q69a_Pneumo)
Eating changes	(Q43_ChInfectn) (Q44_ChInfectn) (Q49_Weight) (Q50_Weight) (Q51_Weight) (Q52_Weight) (Q73_Smoke)

Appendix 12

New joint aches	(Q69k_Arthritis) (Q69f_HD_Angina_) (Q1_Pain_)
Voice	(Q69e_COPD_) (Q10_Cgh_) (Q43_ChInfectn) Q44_ChInfectn)
Skin	(AGE)

Appendix 13 Collinearity diagnostics

```
collin Q10_Cgh_ever Q43_ChInfectn Q44_ChInfectn Q69k_Arthritis_ever Q13k_Cgh_R Q63_New_JPain_12mths Q69e_COPD_ever
Q69f_HD_Angina_ever Q69h_Cancer_ever Q69j_Asbies_ever Q15_Cgh
(obs=264)
```

Collinearity Diagnostics

Variable	VIF	SQRT VIF	Tolerance	R- Squared
Q10_Cgh_ever	1.41	1.19	0.7100	0.2900
Q43_ChInfectn	1.57	1.25	0.6386	0.3614
Q44_ChInfectn	1.54	1.24	0.6499	0.3501
Q69k_Arthritis_ever	1.07	1.04	0.9317	0.0683
Q13k_Cgh_R	1.38	1.17	0.7262	0.2738
Q63_New_JPain_12mths	1.08	1.04	0.9265	0.0735
Q69e_COPD_ever	1.09	1.05	0.9150	0.0850
Q69f_HD_Angina_ever	1.06	1.03	0.9416	0.0584
Q69h_Cancer_ever	1.06	1.03	0.9450	0.0550
Q69j_Asbies_ever	1.07	1.03	0.9341	0.0659
Q15_Cgh	1.26	1.12	0.7942	0.2058

Mean VIF 1.24

	Eigenval	Cond Index
1	4.8830	1.0000
2	1.2620	1.9670
3	0.9627	2.2521
4	0.9410	2.2780
5	0.8603	2.3824
6	0.7145	2.6142
7	0.6849	2.6700
8	0.5358	3.0189
9	0.4710	3.2200
10	0.2972	4.0531
11	0.2708	4.2460
12	0.1168	6.4656

```
Condition Number 6.4656
Eigenvalues & Cond Index computed from scaled raw sscp (w/ intercept)
Det(correlation matrix) 0.3167
```


Appendix 14 Missing data pattern

. mi misstable patterns

		Missing-value patterns (1 means complete)															
Percent	Pattern	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
30%		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
1		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
<1		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
<1		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
<1		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
<1		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1

Appendix 14

[illegible]

[illegible]

Appendix 14

<1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	0	0
	0	1	1	1	1	1	1	0	1	1	1	1	1	1	1
	0	1	1	1	1	1	1	1	1	0	1	1	0	0	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1
<1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1
<1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
<1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1
	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0
	1	0	0	0	0	0	0	0	0	0	1	1	1	0	0
	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
<1	1	1	0	0	0	0	0	0	0	1	0	0	0	0	0
	0	0	0	1	1	0	0	0	0	1	0	0	0	1	1
	1	0	1	0	0	0	1	0	1	1	1	1	1	0	0
	1	0	0	0	1	1	1	0	0	0	1	1	1	0	0
	0	1	1	1	0	0	0	0	0	1	0	0	1	1	1
	1	1	0	0	0	0	0	0	0	0	1	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	1	1	0
<1	1	1	0	0	0	1	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1
	1	0	0	0	0	0	1	1	0	0	0	0	0	0	0
	1	0	0	1	0	0	0	0	1	0	1	1	1	1	1
	0	0	0	0	1	1	0	1	1	1	1	0	0	0	1
	0	0	1	0	0	0	0	0	0	1	0	0	0	0	1
	1	1	1	1	0	1	1	0	1	1	1	0	1	1	1
<1	1	1	0	0	1	1	0	0	1	1	1	1	1	1	0
	0	1	1	1	1	0	0	0	1	1	1	1	1	1	1
	1	1	1	0	1	1	1	1	0	1	1	1	1	1	1
	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1

<1	1 1 1 0	1 0 0 1	1 1 1 1	0 0 0 0
	0 1 1 1	1 1 1 1	1 1 1 1	0 1 1 1
	1 0 1 1	0 0 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 0 0 1	1 1 1 1
	1 1 1 1	1 1 0 0	1 1 0 0	1 1 1 1
	1 1 1 0	1 0 1 0	0 1 1 0	0 0 0 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 0
<1	1 1 1 0	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 0	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 0 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 1 1	0 0 0 1	1 1
<1	1 1 1 1	0 1 1 0	0 0 0 0	1 1 1 1
	1 0 0 0	1 1 1 0	0 1 1 1	1 1 0 1
	0 1 1 0	1 1 1 0	1 0 0 1	1 1 1 1
	1 1 1 0	0 0 1 0	1 1 0 1	0 1 0 0
	0 0 1 0	0 0 0 1	1 1 1 0	0 0 0 0
	0 0 0 1	1 1 1 0	1 1 0 0	1 1 0 0
	0 0 0 0	1 0 0 0	0 0 0 1	0 1
<1	1 1 1 1	0 1 1 1	0 0 1 1	1 1 1 0
	1 1 1 1	0 1 1 1	1 1 1 1	1 1 1 1
	1 1 0 1	1 1 1 1	1 1 1 0	1 1 1 1
	1 1 1 1	0 0 0 1	1 1 1 1	1 1 0 0
	0 0 0 0	0 1 1 1	1 1 1 1	0 0 0 0
	0 0 0 1	1 1 1 1	1 1 0 1	1 1 1 0
	0 0 0 1	1 0 0 1	0 0 0 1	0 0
<1	1 1 1 1	0 1 1 1	0 0 1 1	1 1 1 1
	1 1 1 1	0 1 1 1	1 1 0 0	1 1 1 1
	1 1 0 1	1 1 1 0	0 1 1 1	1 0 0 0
	1 0 0 0	0 1 0 0	0 1 1 1	1 1 0 0
	1 1 1 1	0 0 1 0	1 1 0 0	1 1 1 1
	1 0 0 0	0 0 0 0	0 1 0 0	0 0 0 0
	0 0 0 0	0 0 0 1	0 0 0 1	0 0
<1	1 1 1 1	0 1 1 1	0 0 1 1	1 1 1 1
	1 1 1 1	0 1 1 1	1 1 1 1	1 0 0 0
	0 1 1 1	1 1 0 1	1 1 1 1	1 1 1 1
	0 1 1 0	1 1 1 0	1 1 0 0	0 0 0 0
	1 1 1 1	0 0 1 1	0 0 1 1	1 1 1 1
	1 1 0 1	1 1 0 1	1 0 1 1	1 1 1 0
	0 0 0 0	1 0 0 1	0 0 0 0	0 1
<1	1 1 1 1	0 1 1 1	0 0 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 0 1	1 1 1 1
	1 1 1 0	1 1 1 0	1 1 1 1	1 1 0 0
	0 1 1 1	0 0 1 1	1 1 1 1	0 1 1 1
	1 1 0 1	1 1 1 0	1 1 1 1	1 1 1 0
	0 0 0 0	1 0 0 1	0 0 0 1	0 1
<1	1 1 1 1	1 0 1 1	1 1 1 1	0 0 0 1
	1 1 1 1	1 1 1 1	1 1 1 1	0 1 1 1
	1 0 1 1	0 0 1 1	1 0 1 1	0 1 1 1
	1 1 1 1	1 1 1 1	1 0 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1

Appendix 14

<1	1	1	1	1	1	0	1	1	1	1	1	1	0	0	0	1
	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1
	1	0	1	1	0	0	1	1	1	1	1	1	1	0	0	1
	1	0	0	1	1	1	1	1	0	0	1	1	1	1	1	1
	1	1	1	1	0	1	1	0	1	1	0	0	1	1	1	1
	1	1	0	0	0	0	0	0	0	1	1	0	0	0	0	0
	0	0	0	0	1	1	0	1	1	1	0	1	0	1		
<1	1	1	1	1	1	0	1	1	1	1	1	1	0	0	0	1
	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1
	1	0	1	1	0	0	1	1	1	1	1	1	1	0	1	0
	1	0	0	1	1	1	1	1	0	0	1	1	1	1	1	1
	1	1	1	1	1	1	1	0	0	1	0	0	1	1	1	1
	1	1	1	1	1	1	0	0	0	1	0	0	0	1	0	1
	1	1	1	1	0	1	1	1	1	1	1	1	1	1		
<1	1	1	1	1	1	0	1	1	1	1	1	1	0	0	0	1
	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1
	1	0	1	1	0	0	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	0	1	1	1	1	1	1		
<1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	0
	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
<1	1	1	1	1	1	1	1	0	1	1	0	0	1	1	1	1
	1	0	0	0	0	1	1	1	0	0	0	1	1	0	1	1
	1	1	0	1	1	1	1	1	0	0	0	0	0	1	1	1
	1	1	1	1	0	0	0	1	1	1	1	1	1	1	1	1
	1	0	0	0	1	1	0	1	1	1	1	1	0	0	0	0
	0	0	1	1	1	1	1	1	1	0	1	1	1	1	1	1
	1	1	1	1	1	1	1	0	1	1	1	0	1	1		
<1	1	1	1	1	1	1	1	0	1	1	1	0	1	1	1	1
	1	0	0	1	1	1	1	0	0	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	0		
<1	1	1	1	1	1	1	1	0	1	1	1	0	1	1	1	1
	1	0	0	1	1	1	1	0	0	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1
	1	1	0	0	1	1	1	1	1	1	1	1	0	1	0	0
	1	0	1	1	1	1	1	1	1	0	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1		

[illegible]

Appendix 14

<1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1
	1	1	1	1	1	1	1	1	1	0	0	1	1	1	1
	1	1	1	1	1	1	1	0	1	1	1	0	0	0	0
	1	0	0	1	1	1	1	1	0	1	1	1	1	1	1
	0	1	1	1	1	1	1	0	1	1	0	0	1	1	1
	1	1	1	0	0	0	0	0	0	1	1	0	0	0	1
	1	1	1	1	0	1	1	1	1	1	1	1	1	1	
<1	1	1	1	1	1	1	1	1	1	1	1	0	0	1	1
	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1
	1	0	1	1	0	0	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
<1	1	1	1	1	1	1	1	1	1	1	1	0	1	0	1
	1	0	1	1	1	1	1	1	1	1	1	0	1	1	1
	1	1	1	1	0	0	1	0	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	0	1	1	1	
<1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1
	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1
	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1
	1	0	1	1	1	1	1	1	1	1	1	1	0	1	
<1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1
	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
<1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	0	1	1	1	1	1	1	1	1	0	0	1	1
	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1
	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	0	1	1	1	0	1	0
	1	1	0	1	1	1	1	1	0	0	1	0	1	1	
<1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	0	1	0	1	1	1	1	1	1	0	1	1
	1	0	1	1	1	0	0	1	1	0	1	0	1		
<1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0
	1	1	1	1	0	1	1	1	1	1	1	0	0		

<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 0 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 0
	1 1 0 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 0 1 1	1 1 1 1	1 1 1 1
	1 1 0 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 0 1 1	1 1 1 1	1 1 1 1	0 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 1 1	1 1 1 1	0 1 1 1
	1 1 1 1	0 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 0	1 1 1 0	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	0 0 1 1	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 0
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 0 1 1	0 0 1 1	1 0 0 0	0 1 1 1
	1 1 0 1	1 1 1 0	0 1 1 0	1 1 1 0
	1 1 1 1	0 1 1 1	1 1 1 0	1 1 0 1
	1 1 1 1	1 1 1 0	1 1 1 0	1 1 1 1
	1 1 0 1	1 1 1 0	0 1 1 0	1 1 0 0
	0 0 0 1	0 0 0 1	0 0 0 0	0 0
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 0 1	1 1 1 1	1 1 1 1	1 1 1 0
	0 1 1 1	1 1 0 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 1 1	1 1 1 1	0 0 0 1
	1 1 1 1	1 1 1 0	1 0 0 0	1 1 1 1
	1 1 0 1	1 1 1 1	0 1 1 0	0 0 0 1
	1 0 0 1	1 0 0 1	0 1 0 1	0 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 0 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 0
	1 1 1 0	1 1 0 0	1 1 1 1	1 1 0 0
	1 1 1 1	1 0 1 1	1 1 1 0	1 1 1 1
	1 1 0 1	1 1 1 0	1 1 1 0	1 1 0 0
	0 0 0 0	1 0 0 1	0 0 0 1	0 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 0	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 0 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 0 0 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 0 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	0 0 0 1	1 1 1 1	1 1 1 1
	1 1 1 0	1 1 1 1	0 1 1 1	1 1 1 1
	1 1 1 0	1 1 0 0	1 1 1 1	1 1 0 0
	0 1 1 1	1 0 0 1	1 1 1 1	1 1 1 1
	1 1 0 1	1 1 1 1	1 1 1 1	1 1 1 0
	0 0 0 0	1 0 0 0	0 0 0 1	0 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	0 0 1 1	1 1 1 1	1 1 1 1
	1 1 0 0	1 1 1 1	0 1 0 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 0 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 0	1 1 1 1	1 1

Appendix 14

[illegible]

<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 0 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 0 1
	1 1 1 1	1 0 1 1	0 1 1 1	1 1 1 1
	0 1 1 1	1 1 1 0	1 1 0 1	1 1 1 1
	1 1 1 0	0 0 1 1	1 1 0 1	0 0 1 1
	0 1 1 1	1 1 1 1	1 1 1 1	0 0
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 0 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 0	1 1 1 1	0 0 0 0
	1 1 1 1	0 1 0 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 1 1	1 1 0 0	0 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 0	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 0	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 0	1 1 1 1
	1 1 1 1	1 1 1 0	0 1 1 0	1 1 0 1
	1 1 1 0	1 1 1 1	0 0 0 1	0 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 0	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	0 1 1 1	1 1 1 1
	1 0 1 1	1 1 1 1	1 1 0 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 0
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 0 0 0	1 0 1 0
	1 1 1 1	1 0 1 0	1 1 1 0	0 0 0 0
	0 0 0 1	0 1 1 1	0 0 1 1	0 1 1 0
	1 1 1 1	0 1 1 0	1 1 0 0	1 1 1 1
	1 1 1 0	0 0 0 0	0 0 1 0	0 0 0 0
	0 1 0 0	0 1 0 1	0 0 0 1	0 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 0 1 1	1 0 0 1
	1 1 1 1	1 1 0 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 1 1	1 1 0 0	0 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 0	1 1 1 1	1 1 1 0	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 0 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	0 1 1 0	1 1 1 1
	1 1 1 1	0 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	0 0 0 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 0 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 0 0 0	1 0 0 1
	1 1 1 1	0 1 1 1	1 1 1 1	1 1 1 0
	1 1 1 1	0 1 0 0	1 1 0 1	1 1 1 1
	1 1 0 0	0 0 1 0	0 1 1 0	0 0 0 0
	0 0 1 0	1 0 0 1	0 0 0 1	0 0

[illegible]

<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 0	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 0	1 1 1 1	0 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 0 1	1 1 1 1
	1 1 1 1	1 1 1 0	1 1 1 1	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 0 0 0
	0 1 1 1	1 1 0 1	1 1 1 0	1 1 1 1
	0 0 1 1	1 1 1 1	1 1 0 0	0 0 1 1
	1 1 1 1	1 1 1 1	0 0 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 0 0 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 0	1 0
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 0 0 0
	1 1 1 1	1 1 0 1	1 1 1 1	1 1 0 1
	1 1 1 1	1 1 1 1	0 1 1 0	1 0 1 1
	1 1 1 1	1 1 1 1	0 0 0 0	1 1 1 1
	1 1 1 0	0 0 0 1	0 0 1 0	1 0 1 1
	1 1 1 1	0 1 1 1	1 1 1 0	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 0 0 0
	1 1 1 1	1 1 0 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 1 1	1 1 1 0	0 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 0	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 0 0 0
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	0 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 0	1 0
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 0 0
	0 1 1 1	1 1 0 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 1 1	1 1 0 0	0 1 0 1
	1 1 1 1	0 1 1 1	1 1 1 1	1 1 1 1
	1 1 0 1	1 1 1 1	1 1 1 1	1 1 1 0
	0 0 0 0	1 1 0 1	1 1 0 0	0 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 0 0
	0 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 1 1	0 1 0 0	0 1 1 1
	1 1 1 1	1 1 1 0	1 1 0 1	1 1 1 1
	0 1 0 0	0 0 1 1	0 1 0 0	0 0 1 1
	0 0 0 0	1 0 1 1	0 0 0 0	0 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 0 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 0 1 1	1 1 1 1	1 1 1 1
	1 1 0 0	1 1 1 1	1 1 1 1	0 0 0 0
	0 0 1 1	1 1 1 1	1 1 0 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 0 1	1 1 1 1	1 1 1 1
	0 1 1 1	1 1 1 1	1 1 0 0	0 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 0	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 0 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 0 1	1 1 1 1	1 1 1 1
	1 1 1 1	0 0 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 0
	0 0 1 1	1 1 0 1	1 0 0 1	0 1

Appendix 14

<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 0 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 0	1 0 1 1	1 0 1 0	1 0
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 0 1	1 1 1 1	0 1 1 1	1 1 1 1
	1 1 1 0	1 1 1 0	1 1 1 1	1 1 1 1
	1 1 1 1	1 0 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 0 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	0 1 1 1	1 1 1 1	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 0	1 1 1 1	1 1 1 1	1 0 0 1
	1 0 0 1	1 1 1 1	0 1 1 1	1 1 1 1
	1 1 1 0	1 1 0 0	1 1 0 1	1 1 1 1
	1 1 1 0	0 0 0 1	1 1 1 1	0 0 1 1
	1 1 1 1	0 1 1 0	1 1 1 1	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 0	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 0 0 1	1 1 1 1
	1 1 1 1	1 1 0 1	0 1 1 1	1 1 1 0
	1 1 1 1	1 1 0 1	1 0 1 1	1 1 1 1
	1 1 1 1	1 1 1 0	1 1 1 1	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 0	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 0
	1 1 1 1	0 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 0
	1 1 0 0	1 0 0 1	0 1 1 1	0 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 0 1 1	1 1 1 1	1 0 1 1
	1 1 1 1	1 1 1 1	1 1 1 0	0 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 0	0 1 1 1	0 1 1 1	1 1 1 0
	1 1 1 1	1 1 1 1	1 1 1 0	1 1
<1	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 0 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1 1 1
	1 1 1 1	1 1 1 1	1 1 1 1	1 1

[illegible]

[illegible]

<1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	0	1	1	1	1	0	1	1	1
	1	1	1	1	1	0	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	0	1	1	1
	1	1	1	1	0	1	0	1	1	0	1		
<1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	0	1	1	1	0	1	1	1	1
	1	1	1	1	1	0	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	0	1	0	1	1	0	1		
<1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	0	1	0	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	0	1	1	1		
<1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	0	1	1	1	1	1	1	1
	1	0	0	0	1	1	1	1	1	0	0	0	0
	0	0	1	1	1	1	1	1	1	0	0	1	1
	1	1	1	1	1	1	1	1	1	0	1		
<1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	0	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	0	0	0	0
	0	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
<1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	0	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	0	1	1	1	1	1	1	1	1	1	1	1	1
<1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	0	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	0	1	1	1	1	1
<1	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	1	1	1	1	1	1
	1	1	1	1	1	1	1	0	1	1	1	1	1

[illegible]

Appendix 14

[illegible]

[illegible]

[illegible]

Appendix 14

[illegible]

[illegible]

Appendix 14

Variables are







```

Row 1:  (1) Q2_Pain (2) Q1_Pain_ (3) Q39_CghBlood (4) Q38_CghBlood_ (5) Q11_Cgh (6) Q30_Tired
        (7) Q40_CghBlood_ (8) Q56_HCsweat (9) Q10_Cgh_ (10) Q12_Cgh_ (11) Q65_Voice (12) Q53_HCsweat_
        (13) Q32_Tired (14) Q31_Tired (15) Q35_Tired (16) Q41_CghBlood
Row 2:  (1) Q42_CghBlood_ (2) Q54_HCsweat (3) Q55_HCsweat_ (4) Q66_Voice (5) Q50_Weight
        (6) Q45_ChInfectn_ (7) Q43_ChInfectn (8) Q57_HCsweat (9) Q58_HCsweat (10) Q59_EatChnges
        (11) Q20_BrChnges (12) Q21_BrChnges_ (13) Q33_Tired_ (14) Q3h_Pain (15) Q3b_Pain
        (16) Q3g_Pain
Row 3:  (1) Q3e_Pain (2) Q29_Tired_ (3) Q51_Weight (4) Q46_ChInfectn (5) Q34_Tired_ (6) Q36_Tired
        (7) Q3a_Pain (8) Q19_BrChnges_ (9) Q52_Weight (10) Q63_New_JointachePain_
        (11) Q64_New_JointachePain (12) Q60_EatChnges (13) Q67_Skin (14) Q24a_BrChnges
        (15) Q25b_BrChnges (16) Q22a_BrChnges
Row 4:  (1) Q3d_Pain (2) Q27_BrChnges (3) Q26_BrChnges (4) Q15_Cgh (5) Q61_EatChnges (6) Q69a_Pneumo
        (7) Q49_Weight (8) Q16_Cgh_ (9) Q24c_BrChnges (10) Q37_Tired (11) Q3f_Pain (12) Q3i_Pain
        (13) Q3c_Pain (14) Q7_Pain (15) Q13a_Cgh (16) Q13i_Cgh
Row 5:  (1) Q44_ChInfectn (2) Q69c_Asthma_ (3) Q69j_Asbes (4) Q69f_HD_Angina_ (5) Q13j_Cgh
        (6) Q18_Cgh (7) Q47_ChInfectn_ (8) Q25a_BrChnges (9) Q8_Pain_ (10) Q6_Pain (11) Q25c_BrChnges
        (12) Q22c_BrChnges (13) Q69i_Emphys (14) Q69b_BronchitisCB_ (15) Q69d_Allergy_
        (16) Q69g_Anaemia_
Row 6:  (1) Q69k_Arthritis (2) Q69h_Cancer_ (3) Q13b_Cgh (4) Q24b_BrChnges (5) Q24f_BrChnges
        (6) Q25d_BrChnges (7) Q28_BrChnges (8) Q22b_BrChnges (9) Q22e_BrChnges (10) Q9_Pain
        (11) Q69e_COPD_ (12) Q22d_BrChnges (13) Q24e_BrChnges (14) Q24d_BrChnges (15) Q22f_BrChnges
        (16) Q13d_Cgh_
Row 7:  (1) Q13e_Cgh_ (2) Q13h_Cgh (3) Q13g_Cgh (4) Q13l_Cgh (5) Q23_BrChnges (6) Q14a_Cgh
        (7) Q13c_Cgh (8) Q48_ChInfectn (9) Q14c_Cgh_ (10) Q14b_Cgh_ (11) Q13k_Cgh (12) Q3j_Pain
        (13) Q13f_Cgh_ (14) Q62_EatChnges

```

Appendix 15 Tetrachoric correlations to determine response cut-offs

Examples of codings for individual response cut-off

<pre>gen Q12_Cgh_1=. replace Q12_Cgh_1=0 if (Q12_Cgh_==0) replace Q12_Cgh_1=1 if (Q12_Cgh_==1) replace Q12_Cgh_1=0 if (Q12_Cgh_==2) replace Q12_Cgh_1=0 if (Q12_Cgh_==3)</pre>		First experience 3 months ago
<pre>gen Q12_Cgh_2=. replace Q12_Cgh_2=0 if (Q12_Cgh_==0) replace Q12_Cgh_2=1 if (Q12_Cgh_==1) replace Q12_Cgh_2=1 if (Q12_Cgh_==2) replace Q12_Cgh_2=0 if (Q12_Cgh_==3)</pre>		First experience <12 months ago
<pre>gen Q12_Cgh_3=. replace Q12_Cgh_3=0 if (Q12_Cgh_==0) replace Q12_Cgh_3=1 if (Q12_Cgh_==1) replace Q12_Cgh_3=1 if (Q12_Cgh_==2) replace Q12_Cgh_3=1 if (Q12_Cgh_==3)</pre>		First experience >12 months ago
<pre>gen Q13a_Cgh_1=. replace Q13a_Cgh_1=0 if (Q13a_Cgh_==0) replace Q13a_Cgh_1=1 if (Q13a_Cgh_==1) replace Q13a_Cgh_1=1 if (Q13a_Cgh_==2) replace Q13a_Cgh_1=1 if (Q13a_Cgh_==3)</pre>		Never/ever
<pre>gen Q13a_Cgh_2=. replace Q13a_Cgh_2=0 if (Q13a_Cgh_==0) replace Q13a_Cgh_2=0 if (Q13a_Cgh_==1) replace Q13a_Cgh_2=1 if (Q13a_Cgh_==2) replace Q13a_Cgh_2=1 if (Q13a_Cgh_==3)</pre>		Occasionally to most of the time
<pre>gen Q13a_Cgh_3=. replace Q13a_Cgh_3=0 if (Q13a_Cgh_==0) replace Q13a_Cgh_3=0 if (Q13a_Cgh_==1) replace Q13a_Cgh_3=0 if (Q13a_Cgh_==2) replace Q13a_Cgh_3=1 if (Q13a_Cgh_==3)</pre>		Only most of the time
<pre>gen Q13b_Cgh_1=. replace Q13b_Cgh_1=0 if (Q13b_Cgh_==0) replace Q13b_Cgh_1=1 if (Q13b_Cgh_==1) replace Q13b_Cgh_1=1 if (Q13b_Cgh_==2) replace Q13b_Cgh_1=1 if (Q13b_Cgh_==3)</pre>		
<pre>gen Q13b_Cgh_2=. replace Q13b_Cgh_2=0 if (Q13b_Cgh_==0) replace Q13b_Cgh_2=0 if (Q13b_Cgh_==1) replace Q13b_Cgh_2=1 if (Q13b_Cgh_==2) replace Q13b_Cgh_2=1 if (Q13b_Cgh_==3)</pre>		
<pre>gen Q13b_Cgh_3=. replace Q13b_Cgh_3=0 if (Q13b_Cgh_==0) replace Q13b_Cgh_3=0 if (Q13b_Cgh_==1) replace Q13b_Cgh_3=0 if (Q13b_Cgh_==2) replace Q13b_Cgh_3=1 if (Q13b_Cgh_==3)</pre>		

Appendix 15

Stata Output: Tetrachoric correlations

```
. tetrachoric Q1_Pain_ever Q1_Pain_current Q1_Pain_3mths LCdiagnosis
(obs=355)
```

```
matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0187
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q1_Pai~r	Q1_Pai~t	Q1_Pai~s	LCdiag~s
Q1_Pain_ever	1.0000			
Q1_Pain_cu~t	1.0000	1.0000		
Q1_Pain_3m~s	1.0000	1.0000	1.0000	
LCdiagnosis	-0.1496	0.0484	-0.0763	1.0000

```
.
. tetrachoric Q10_Cgh_ever Q10_Cgh_current Q10_Cgh_3mths LCdiagnosis
(obs=347)
```

```
matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0008
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q10_Cg~r	Q10_Cg~t	Q10_Cg~s	LCdiag~s
Q10_Cgh_ever	1.0000			
Q10_Cgh_cu~t	1.0000	1.0000		
Q10_Cgh_3m~s	0.9970	1.0000	1.0000	
LCdiagnosis	0.0987	0.1003	0.1156	1.0000

```
. tetrachoric Q12_Cgh_1 Q12_Cgh_2 Q12_Cgh_3 LCdiagnosis
(obs=346)
```

```
matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0056
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q12_Cg~1	Q12_Cg~2	Q12_Cg~3	LCdiag~s
Q12_Cgh_1	1.0000			
Q12_Cgh_2	1.0000	1.0000		
Q12_Cgh_3	1.0000	1.0000	1.0000	
LCdiagnosis	0.2170	0.1820	0.1124	1.0000

```
. tetrachoric Q13b_Cgh_1 Q13b_Cgh_2 Q13b_Cgh_3 LCdiagnosis
(obs=331)
```

```
matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0012
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q13b_C~1	Q13b_C~2	Q13b_C~3	LCdiag~s
Q13b_Cgh_1	1.0000			
Q13b_Cgh_2	1.0000	1.0000		
Q13b_Cgh_3	1.0000	1.0000	1.0000	
LCdiagnosis	0.1399	0.0918	0.1041	1.0000

```
. tetrachoric Q13c_Cgh_1 Q13c_Cgh_2 Q13c_Cgh_3 LCdiagnosis
(obs=323)
```

```
matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0002
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q13c_C~1	Q13c_C~2	Q13c_C~3	LCdiag~s
Q13c_Cgh_1	1.0000			
Q13c_Cgh_2	1.0000	1.0000		
Q13c_Cgh_3	1.0000	1.0000	1.0000	
LCdiagnosis	0.0878	0.0878	0.0674	1.0000

```
. tetrachoric Q13d_Cgh_1 Q13d_Cgh_2 Q13d_Cgh_3 LCdiagnosis
(obs=328)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0127
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q13d_C-1	Q13d_C-2	Q13d_C-3	LCdiag-s
Q13d_Cgh_1	1.0000			
Q13d_Cgh_2	1.0000	1.0000		
Q13d_Cgh_3	1.0000	1.0000	1.0000	
LCdiagnosis	0.1262	0.1319	-0.0305	1.0000

```
. tetrachoric Q13e_Cgh_1 Q13e_Cgh_2 Q13e_Cgh_3 LCdiagnosis
(obs=327)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0046
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q13e_C-1	Q13e_C-2	Q13e_C-3	LCdiag-s
Q13e_Cgh_1	1.0000			
Q13e_Cgh_2	1.0000	1.0000		
Q13e_Cgh_3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0320	0.0136	-0.0830	1.0000

```
. tetrachoric Q13f_Cgh_1 Q13f_Cgh_2 Q13f_Cgh_3 LCdiagnosis
(obs=311)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0010
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q13f_C-1	Q13f_C-2	Q13f_C-3	LCdiag-s
Q13f_Cgh_1	1.0000			
Q13f_Cgh_2	1.0000	1.0000		
Q13f_Cgh_3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0555	-0.0363	-0.0807	1.0000

```
. tetrachoric Q13g_Cgh_1 Q13g_Cgh_2 Q13g_Cgh_3 LCdiagnosis
(obs=327)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0040
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q13g_C-1	Q13g_C-2	Q13g_C-3	LCdiag-s
Q13g_Cgh_1	1.0000			
Q13g_Cgh_2	1.0000	1.0000		
Q13g_Cgh_3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0757	-0.0643	-0.1542	1.0000

```
. tetrachoric Q13h_Cgh_1 Q13h_Cgh_2 Q13h_Cgh_3 LCdiagnosis
(obs=328)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0007
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q13h_C-1	Q13h_C-2	Q13h_C-3	LCdiag-s
Q13h_Cgh_1	1.0000			
Q13h_Cgh_2	1.0000	1.0000		
Q13h_Cgh_3	1.0000	1.0000	1.0000	
LCdiagnosis	0.0412	0.0174	0.0039	1.0000

```
. tetrachoric Q13i_Cgh_1 Q13i_Cgh_2 Q13i_Cgh_3 LCdiagnosis
(obs=335)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0021
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q13i_C-1	Q13i_C-2	Q13i_C-3	LCdiag-s
Q13i_Cgh_1	1.0000			
Q13i_Cgh_2	1.0000	1.0000		
Q13i_Cgh_3	1.0000	1.0000	1.0000	
LCdiagnosis	0.0044	0.0100	0.0700	1.0000

```
. tetrachoric Q13j_Cgh_1 Q13j_Cgh_2 Q13j_Cgh_3 LCdiagnosis
(obs=335)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0010
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q13j_C-1	Q13j_C-2	Q13j_C-3	LCdiag-s
Q13j_Cgh_1	1.0000			
Q13j_Cgh_2	1.0000	1.0000		
Q13j_Cgh_3	1.0000	1.0000	1.0000	
LCdiagnosis	0.0857	0.0420	0.0559	1.0000

Appendix 15

```
. tetrachoric Q13k_Cgh_1 Q13k_Cgh_2 Q13k_Cgh_3 LCdiagnosis
(obs=317)
```

```
matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0069
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q13k_C~1	Q13k_C~2	Q13k_C~3	LCdiag-s
Q13k_Cgh_1	1.0000			
Q13k_Cgh_2	1.0000	1.0000		
Q13k_Cgh_3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.1618	-0.1001	-0.2169	1.0000

```
. tetrachoric Q13l_Cgh_1 Q13l_Cgh_2 Q13l_Cgh_3 LCdiagnosis
(obs=326)
```

```
matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0065
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q13l_C~1	Q13l_C~2	Q13l_C~3	LCdiag-s
Q13l_Cgh_1	1.0000			
Q13l_Cgh_2	1.0000	1.0000		
Q13l_Cgh_3	1.0000	1.0000	1.0000	
LCdiagnosis	0.0689	0.0932	0.1831	1.0000

```
. tetrachoric Q19_BrChnges_ever Q19_BrChnges_current Q19_BrChnges_3mths LCdiagnosis
(obs=338)
```

```
matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0007
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q19_Br~r	Q19_Br~t	Q19_Br~s	LCdiag-s
Q19_BrChng-r	1.0000			
Q19_BrChng-t	1.0000	1.0000		
Q19_BrChng-s	1.0000	1.0000	1.0000	
LCdiagnosis	0.0274	0.0525	0.0646	1.0000

```
. tetrachoric Q21_BrChnges_1 Q21_BrChnges_2 Q21_BrChnges_3 LCdiagnosis
(obs=344)
```

```
matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0244
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q21_Br~1	Q21_Br~2	Q21_Br~3	LCdiag-s
Q21_BrChng-1	1.0000			
Q21_BrChng-2	1.0000	1.0000		
Q21_BrChng-3	1.0000	1.0000	1.0000	
LCdiagnosis	0.2312	0.2648	0.0393	1.0000

```
. tetrachoric Q22a_BrChnges_1 Q22a_BrChnges_2 Q22a_BrChnges_3 LCdiagnosis
(obs=340)
```

```
matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0006
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q22a_B~1	Q22a_B~2	Q22a_B~3	LCdiag-s
Q22a_BrChn-1	1.0000			
Q22a_BrChn-2	1.0000	1.0000		
Q22a_BrChn-3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0811	-0.0855	-0.1161	1.0000

```
. tetrachoric Q22b_BrChnges_1 Q22b_BrChnges_2 Q22b_BrChnges_3 LCdiagnosis
(obs=332)
```

```
matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0001
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q22b_B~1	Q22b_B~2	Q22b_B~3	LCdiag-s
Q22b_BrChn-1	1.0000			
Q22b_BrChn-2	1.0000	1.0000		
Q22b_BrChn-3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0400	-0.0345	-0.0513	1.0000

```
. tetrachoric Q22c_BrChnges_1 Q22c_BrChnges_2 Q22c_BrChnges_3 LCdiagnosis
(obs=334)
```

```
matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0052
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q22c_B~1	Q22c_B~2	Q22c_B~3	LCdiag-s
Q22c_BrChn-1	1.0000			
Q22c_BrChn-2	1.0000	1.0000		
Q22c_BrChn-3	1.0000	1.0000	1.0000	
LCdiagnosis	0.0005	-0.0414	-0.1020	1.0000

```
. tetrachoric Q22d_BrChnges_1 Q22d_BrChnges_2 Q22d_BrChnges_3 LCdiagnosis
(obs=329)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0002
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q22d_B-1	Q22d_B-2	Q22d_B-3	LCdiag-s
Q22d_BrChn-1	1.0000			
Q22d_BrChn-2	1.0000	1.0000		
Q22d_BrChn-3	1.0000	1.0000	1.0000	
LCdiagnosis	0.0243	0.0271	0.0074	1.0000

```
. tetrachoric Q22e_BrChnges_1 Q22e_BrChnges_2 Q22e_BrChnges_3 LCdiagnosis
(obs=332)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0028
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q22e_B-1	Q22e_B-2	Q22e_B-3	LCdiag-s
Q22e_BrChn-1	1.0000			
Q22e_BrChn-2	1.0000	1.0000		
Q22e_BrChn-3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0112	-0.0231	0.0513	1.0000

```
. tetrachoric Q22f_BrChnges_1 Q22f_BrChnges_2 Q22f_BrChnges_3 LCdiagnosis
(obs=329)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0028
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q22f_B-1	Q22f_B-2	Q22f_B-3	LCdiag-s
Q22f_BrChn-1	1.0000			
Q22f_BrChn-2	1.0000	1.0000		
Q22f_BrChn-3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0400	-0.0616	0.0138	1.0000

```
. tetrachoric Q29_Tired_ever Q29_Tired_current Q29_Tired_3mths LCdiagnosis
(obs=342)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0015
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q29_Ti-r	Q29_Ti-t	Q29_Ti-s	LCdiag-s
Q29_Tired_r	1.0000			
Q29_Tired_t	1.0000	1.0000		
Q29_Tired_s	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0075	0.0476	0.0075	1.0000

```
. tetrachoric Q31_Tired_1 Q31_Tired_2 Q31_Tired_3 LCdiagnosis
(obs=345)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0099
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q31_Ti-1	Q31_Ti-2	Q31_Ti-3	LCdiag-s
Q31_Tired_1	1.0000			
Q31_Tired_2	1.0000	1.0000		
Q31_Tired_3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0299	0.1131	-0.0258	1.0000

```
. tetrachoric Q32_Tired_1 Q32_Tired_2 Q32_Tired_3 LCdiagnosis
(obs=345)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0018
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q32_Ti-1	Q32_Ti-2	Q32_Ti-3	LCdiag-s
Q32_Tired_1	1.0000			
Q32_Tired_2	1.0000	1.0000		
Q32_Tired_3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0222	-0.0116	-0.0722	1.0000

```
. tetrachoric Q33_Tired_1 Q33_Tired_2 Q33_Tired_3 LCdiagnosis
(obs=343)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0015
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q33_Ti-1	Q33_Ti-2	Q33_Ti-3	LCdiag-s
Q33_Tired_1	1.0000			
Q33_Tired_2	1.0000	1.0000		
Q33_Tired_3	1.0000	1.0000	1.0000	
LCdiagnosis	0.0132	0.0041	-0.0416	1.0000

Appendix 15

```
. tetrachoric Q34_Tired_1 Q34_Tired_2 Q34_Tired_3 LCdiagnosis
(obs=340)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0009
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q34_Ti~1	Q34_Ti~2	Q34_Ti~3	LCdiag~s
Q34_Tired_1	1.0000			
Q34_Tired_2	1.0000	1.0000		
Q34_Tired_3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0274	-0.0308	0.0113	1.0000

```
.
. tetrachoric Q38_CghBlood_ever Q38_CghBlood_current Q38_CghBlood_3mths LCdiagnosis
(obs=349)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0014
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q38_Cg~r	Q38_Cg~t	Q38_Cg~s	LCdiag~s
Q38_CghBlo~r	1.0000			
Q38_CghBlo~t	1.0000	1.0000		
Q38_CghBlo~s	0.9973	1.0000	1.0000	
LCdiagnosis	-0.0161	0.0331	0.0286	1.0000

```
. tetrachoric Q40_CghBlood_1 Q40_CghBlood_2 Q40_CghBlood_3 LCdiagnosis
(obs=347)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0044
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q40_Cg~1	Q40_Cg~2	Q40_Cg~3	LCdiag~s
Q40_CghBlo~1	1.0000			
Q40_CghBlo~2	1.0000	1.0000		
Q40_CghBlo~3	1.0000	1.0000	1.0000	
LCdiagnosis	0.0972	0.0434	0.0035	1.0000

```
. tetrachoric Q41_CghBlood_1 Q41_CghBlood_2 Q41_CghBlood_3 LCdiagnosis
(obs=343)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0256
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q41_Cg~1	Q41_Cg~2	Q41_Cg~3	LCdiag~s
Q41_CghBlo~1	1.0000			
Q41_CghBlo~2	1.0000	1.0000		
Q41_CghBlo~3	1.0000	1.0000	1.0000	
LCdiagnosis	0.0791	0.0794	-0.1569	1.0000

```
. tetrachoric Q42_CghBlood_1 Q42_CghBlood_2 Q42_CghBlood_3 LCdiagnosis
(obs=344)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0222
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q42_Cg~1	Q42_Cg~2	Q42_Cg~3	LCdiag~s
Q42_CghBlo~1	1.0000			
Q42_CghBlo~2	1.0000	1.0000		
Q42_CghBlo~3	1.0000	1.0000	1.0000	
LCdiagnosis	0.0248	0.0030	-0.1928	1.0000

```
. tetrachoric Q53_HCsweat_ever Q53_HCsweat_current Q53_HCsweat_3mths LCdiagnosis
(obs=343)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0003
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q53_HC~r	Q53_HC~t	Q53_HC~s	LCdiag~s
Q53_HCswea~r	1.0000			
Q53_HCswea~t	1.0000	1.0000		
Q53_HCswea~s	1.0000	1.0000	1.0000	
LCdiagnosis	0.0315	0.0533	0.0306	1.0000

```
. tetrachoric Q55_HCsweat_1 Q55_HCsweat_2 Q55_HCsweat_3 LCdiagnosis
(obs=344)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0010
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q55_HC~1	Q55_HC~2	Q55_HC~3	LCdiag~s
Q55_HCswea~1	1.0000			
Q55_HCswea~2	1.0000	1.0000		
Q55_HCswea~3	1.0000	1.0000	1.0000	
LCdiagnosis	0.0514	0.0142	0.0057	1.0000

```
. tetrachoric Q56_HCsweat_1 Q56_HCsweat_2 Q56_HCsweat_3 LCdiagnosis
(obs=347)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0036
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q56_HC-1	Q56_HC-2	Q56_HC-3	LCdiag-s
Q56_HCswea-1	1.0000			
Q56_HCswea-2	1.0000	1.0000		
Q56_HCswea-3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.0156	-0.0341	-0.1012	1.0000

```
. tetrachoric Q57_HCsweat_1 Q57_HCsweat_2 Q57_HCsweat_3 LCdiagnosis
(obs=343)

matrix with tetrachoric correlations is not positive semidefinite;
it has 1 negative eigenvalue
maxdiff(corr,adj-corr) = 0.0020
(adj-corr: tetrachoric correlations adjusted to be positive semidefinite)
```

	Q57_HC-1	Q57_HC-2	Q57_HC-3	LCdiag-s
Q57_HCswea-1	1.0000			
Q57_HCswea-2	1.0000	1.0000		
Q57_HCswea-3	1.0000	1.0000	1.0000	
LCdiagnosis	-0.1576	-0.1715	-0.2194	1.0000

List of References

- Abayomi K, Gelman A and Levy M (2008) Diagnostics for multivariate imputations. *Applied Statistics* 57(Series C, Part 3):273-291
- Ades AE, Biswas M, Welton NJ, Hamilton W (2014) Symptom lead time distribution in lung cancer: natural history and prospects for early diagnosis. *International Journal of Epidemiology* 43(6):1865-73
- Agresti A (1996) *An Introduction to Categorical Data Analysis*. New York: John Wiley and Sons, Inc
- AJCC, American Joint Committee on Cancer (2002) *AJCC staging manual* (6th Edition). New York: Springer-Verlag
- Akaike H, Parzen E, Tanabe K and Kitagawa G (1998) *Selected papers of hirotugu akaike*. New York: Springer
- Alavanja MCR, Brownson RC, Boice JD and Hock E (1992) Preexisting lung disease and lung cancer among nonsmoking women. *American Journal of Epidemiology* 136(6):623-632
- Alberg AJ and Samet JM (2003) Epidemiology of lung cancer. *Chest* 123(Supplement 1):21S-49S
- Allison PD (2002) *Missing Data*. Thousand Oaks, California: Sage Publication
- Altman DG and Bland JM (1994a) Diagnostic tests 2: Predictive values. *British Medical Journal* 309(6947):102
- Altman DG and Bland JM (1994b) Diagnostic tests 1: Sensitivity and specificity. *British Medical Journal* 308(6843):1552
- Altman DG and Bland JM (1994c) Diagnostic tests 3: receiver operating characteristic plots. *British Medical Journal* 309(6948):188
- American Cancer Society (2006) *Cancer facts and figures 2006*. Available from: <http://www.cancer.org/acs/groups/content/@nho/documents/document/cff2006pwsecuredpdfpdf> [Accessed 10th January 2012]
- American College of Asthma, Allergy and Immunology (2010) *Asthma Symptoms: Overview*. Available from: <http://acaai.org/asthma/symptoms> [Accessed 16th January 2015]
- Andersen BL and Cacioppo JT (1995) Delay in seeking a cancer diagnosis: delay stages and psychophysiological comparison processes. *British Journal of Social Psychology* 34(Pt 1):33-52

References

- Andersen HA and Prakash UBS (1982) Diagnosis of symptomatic lung cancer. *Seminar of Respiratory Medicine* 3:165-175
- Anthonisen NR, Dik N, Manfreda J and Roos LL (2001) Spirometry and obstructive lung disease in Manitoba. *Canadian Respiratory Journal: journal of Canadian Thoracic Society* 8(6):421-426
- Attia J (2003) DIAGNOSTIC TESTS: Moving beyond sensitivity and specificity: using likelihood ratios to help interpret diagnostic tests. *Australian Prescriber* 26(5): 111-113
- Bankhead C (2005) *Identifying Potentially Significant Diagnostic Factors For Ovarian Cancer in Primary Care: A Qualitative and Quantitative Study*. Unpublished PhD thesis Oxford University
- Bankhead CR, Collins C, Stokes-Lampard H, Rose P, Wilson S, Clements A, Mant D, Kehoe ST and Austoker J (2008) Identifying symptoms of ovarian cancer: a qualitative and quantitative study. *BJOG: An International Journal of Obstetrics & Gynaecology* 115(8):1008-1014
- Barnes PJ, Shapiro SD and Pauwels RA (2003) Chronic obstructive pulmonary disease: Molecular and cellular mechanisms. *The European Respiratory Journal* 22(4):672-688
- Barton MB, Elmore JG and Fletcher SW (1999) Breast symptoms among women enrolled in a Health Maintenance Organisation: frequency, evaluation, and outcome. *Annals of Internal medicine* 130(8):651-657
- Beale E M L (1970) Note on procedures for variable selection in multiple regression. *Technometrics* 12: 909-914
- Bechtel JJ, Kelley WA, Coons TA, Mohler P, Mohler A, James S and Petty TL (2009) Five-year Outcome of Lung Cancer Detection in Patients With and Without Airflow Obstruction in a Primary Care Outpatient Practice. *Journal of Thoracic Oncology* 4(11):1347-1351
- Beckles MA, Spiro SG, Colice GL and Rudd RM (2003) Initial evaluation of the patient with lung cancer: symptoms, signs, laboratory tests, and paraneoplastic syndromes. *Chest* 123(1):97s-104s
- Belsley DA, Kuh E and Welsch RE (1980) *Regression Diagnostics: identifying influential data and sources of collinearity*. New York: Wiley
- Berry WD and Feldman S (1985) *Multiple regression in practice*. Beverley Hills: SAGE University Press
- Billing JS and Wells FC (1996) Delays in the diagnosis and surgical treatment of lung cancer. *Thorax* 51(9):903-906

- Bjerager M, Palshof T, Dahl R, Vedsted P and Olesen F (2006) Delay in diagnosis of lung cancer in general practice. *British Journal of General Practice* 56(532):863-868
- Bland JM and Altman GA (2000) The odds ratio. *British Medical Journal* 320: 1468
- Bofetta P and Kogevinas M (1999) Epidemiological research and prevention of occupational cancer in Europe. *Environmental Health Perspectives* 107(Suppl 2):229-231
- Boucot KR, Seidman H and Weiss W (1977) The Philadelphia Pulmonary Neoplasm Research Project: The risk of lung cancer in relation to symptoms and roentgenographic abnormalities. *Environmental Research* 13(3):451-469
- Bowen EF and Rayner CJF (2002) Patient and GP led delays in the recognition of symptoms suggestive of lung cancer. *Lung Cancer* 37(2):227-228
- Boyatzis RE (1998) *Transforming Qualitative Information: Thematic Analysis and Code Development*. Thousand Oaks, CA: Sage Publications
- Boyer B, Valles AM and Edme N (2000) Induction and regulation of epithelial-mesenchymal transitions. *Biochemical Pharmacology* 60(8):1091-1099
- Brancato G, Macchia S, Murgia M, Signore M, Simeoni G, Blanke K, Körner T, Nimmergut A, Lima P, Paulino R and Hoffmeyer-Zlotnik JHP (2006) Handbook of recommended practices for questionnaire development and testing in the European statistical system. Available from: http://www.istat.it/en/files/2013/12/Handbook_questionnaire_development_2006.pdf [Accessed 10th October 2015]
- Braun V and Clarke V (2006) Using thematic analysis in psychology. *Qualitative Research in Psychology* 3(2):77-101
- Brett GZ (1968) The value of lung cancer detection by six-monthly chest radiographs. *Thorax* 23(4):414-420
- Brindle L, Pope C, Corner J, Leydon G and Banerjee A (2012) Eliciting symptoms interpreted as normal by patients with early-stage lung cancer: could GP elicitation of normalised symptoms reduce delay in diagnosis? Cross-sectional interview study. *British Medical Journal Open* 2(6):e001977
- Brindle LA, Dowswell G, James EP, Clifford S, Ocansey L, Hamilton W, Banerjee A, George S, Djearaman M, Aitchinson F, Grove A, Chee S, Rudran B, Miller B, Indrajeet D and Wilson S, on behalf of the IPCARD Feasibility Study team (2015) Using a participant-completed questionnaire to identify Symptoms that Predict Chest and Respiratory Disease (IPCARD): A Feasibility Study. Report to NSPCR

References

- Brindle LA, Hamilton W., Banerjee A and Dowswell G (2014) Symptoms that predict chest X-ray results suspicious for lung cancer in UK primary care: results from a prospective study. *European Journal of Cancer Care* 23(Supplement 1): 3-4
- British Lung Foundation (1998) The health costs of transport and air pollution. *Nursing standard (Royal College of Nursing)* 12(26):32-33
- British Lung Foundation (2003) Lung Report III- casting a shadow over the nation's health. Available from:
<http://www.blf.org.uk/Search?query=lung+report&searchButtonx=0&searchButtony=0> [Accessed May 2012]
- British Lung Foundation (2012) Lung Cancer. Available from:
<http://www.blf.org.uk/Conditions/Detail/lung-cancer#overview> [Accessed January 2012]
- Brody JS and Spira A (2006) State of the art Chronic obstructive pulmonary disease, inflammation, and lung cancer. *Proceedings of the American Thoracic Society* 3(6):535-537
- Brownson RC and Alavanja MC (2000) Previous lung disease and lung cancer risk among women (United States). *Cancer Causes Control* 11(9):853-858
- Brownson RC, Alavanja MC, Caporaso N, Simoes EJ and Chang JC (1998) Epidemiology and prevention of lung cancer in nonsmokers. *Epidemiologic Reviews* 20(2):218-236
- Bruinsma SM, Rietjens JA, Seymour JE, Anquinet L and van der Heide A. (2012) The experiences of relatives with the practice of palliative sedation: a systematic review. *Journal of pain and symptom management* 44(3): 431-5
- Buccheri G and Ferrigno D (2004) Lung cancer: clinical presentation and specialist referral time. *The European Respiratory Journal* 24(6):898-904
- Buist AS (1988) Smoking and other risk factors IN: Murray JF and Nadel JA (eds) *Textbook of respiratory medicine*. Philadelphia: WB Saunders Company
- Burgess CC, Ramirez AJ, Richards MA and Love SB (1998) Who and what influences delayed presentation in breast cancer?. *British Journal of Cancer* 77(8):1343-1348
- Calabrò E, Randi G, Vecchia CL, Sverzellati N, Marchiano A, Villani M, Zompatori M, Cassandro R, Harari S and Pastorino U (2010) Lung function predicts lung cancer risk in smokers: a tool for targeting screening programmes. *European Respiratory Journal* 35(1):146-151
- Cancer Research UK (2004) CancerStats Monograph London: Cancer Research UK

- Cancer Research UK (2009) Lung cancer incidence statistics. Available from: <http://www.cancerresearchuk.org/cancer-info/cancerstats/types/lung/incidence/> [Accessed April 2012]
- Cancer Research UK (2014) *Be Clear on Cancer: Evaluation Summary*. Available from: http://www.cancerresearchuk.org/prod_consump/groups/cr_common/@nre/@hea/documents/generalcontent/cr_119405.pdf [Accessed 2nd April 2015]
- Cancer Research UK (2014) Cancer Statistics Report: Cancer Incidence and Mortality in the UK, January 2014 (Incidence 2011, Mortality 2011). Available from: http://publications.cancerresearchuk.org/downloads/Product/CS_REPORT_TOP10INCMORT.pdf [Accessed 05th February 2014]
- Cancer Research UK cancer statistics report (2009) CancerStats: Cancer Statistics for the UK. Available from: <http://www.cancerresearchuk.org/cancer-info/cancerstats/> [Accessed January 2012]
- Carpenter JR and Goldstein H (2004) Multiple imputation in MLwiN. *Multilevel modelling newsletter* 16(2)
- Carpenter JR and Kenward MG (2008) *Missing data in clinical trials - a practical guide*. Birmingham: National Health Service Co-ordinating Centre for Research Methodology. Available from: www.missingdata.org.uk [Accessed 28th September 2014]
- Carpenter JR, Kenward MG and White IR (2007) *Sensitivity analysis after multiple imputation under missing at random - a weighting approach*. *Statistical Methods in Medical Research* 16(3):259-275
- Cassidy A, Myles JP, van Tongeren M, Page RD, Liloglou T, Duffy SW and Field JK (2008) The LLP risk model: an individual risk prediction model for lung cancer. *British Journal of Cancer* 98(2):270-276
- CEBM Centre for evidence-based medicine (2014) Likelihood Ratios. Available from: <http://www.cebm.net/likelihood-ratios/> [Accessed 10th October 2015]
- Coleman MP, Babb P and Damiecki P (1999) *Cancer survival trends in England and Wales, 1971-1995: Deprivation and NHS Region*. London: The Stationery Office
- Coleman MP, Forman D, Bryant H, Butler J, Rachet B, Maringe C, Nur U, Tracey E, Coory M, Hatcher J, McGahan CE, Turner D, Marrett L, Gjerstorff ML, Johannesen TB, Adolfsson J, Lambe M, Lawrence G, Meechan D, Morris EJ, Middleton R, Steward J, Richards MA and ICBP Module 1 Working Group (2011) Cancer survival in Australia, Canada, Denmark, Norway, Sweden,

References

- and the UK, 1995-2007 (the International Cancer Benchmarking Partnership): an analysis of population-based cancer registry data. *Lancet* 377(9760):127-38
- Corner J, Hopkinson J and Roffe L (2006) Experience of health changes and reasons for delay in seeking care: A UK study of the months prior to the diagnosis of lung cancer. *Social Science & Medicine* 62(6):1381-1391
- Corner J, Hopkinson J, Fitzsimmons D, Barclay S and Muers M (2005) Is late diagnosis of lung cancer inevitable? Interview study of patients' recollections of symptoms before diagnosis. *Thorax* 60(4):314-319
- Creswell JW (2003) *RESEARCH DESIGN: Qualitative, Quantitative and Mixed Methods Approaches* (2nd Edition). Thousand Oaks, California: Sage Publications
- Dasari V, Gallup M, Lemjabbar H, Maltseva I and Mc Namara N (2006) Epithelial-mesenchymal transition in lung cancer: is tobacco the "smoking gun"? *American Journal of Respiratory Cell and Molecular Biology* 35(1):3-9
- de Torres JP, Bastarrika G, Wisnivesky JP, Alcaide AB, Campo A, Seijo LM, Pueyo JC, Villanueva A, Lozano MD, Montes U, Montyenga L and Zulueta JJ (2007) Assessing the relationship between lung cancer risk and emphysema detected on low-dose CT of the chest. *Chest* 132(6):1932-1938
- Deeks JJ (2001) Systematic reviews in health care: Systematic reviews of evaluations of diagnostic and screening tests. *British Medical Journal* 323(7305):157-162
- Department of Health (2007) *Cancer Reform Strategy*. London: Department of Health
- DeVon HA, Block ME, Moyle-Wright P, Ernst DM, Hayden SJ, Lazzara DJ, Savoy SM and Kostas-Polston E (2007) A psychometric Toolbox for testing Validity and Reliability. *Journal of Nursing scholarship* 39 (2): 155-164.
- Diez-Herranz A (2001) COPD and lung cancer: practical implications. *Archivos de Bronconeumologia* 37(5):4673-4678
- Doll R and Hill AB (1950) Smoking and carcinoma of the lung Preliminary report. *British Medical Journal* 2(4682):739-748
- dos Santos Silva (1999) *Cancer Epidemiology: Principles and Methods*. Lyon, France: International Agency for Research on Cancer (IARC)
- Dragow F (1988) Polychoric and polyserial correlations IN: Kotz L and Johnson L (eds) *Encyclopedia of statistical sciences*. New York: Wiley & sons

- Eagan TM, Gulsvik A, Eide GE and Bakke PS (2004) Remission of respiratory symptoms by smoking and occupational exposure in a cohort study. *European Respiratory Journal* 23(4):589-594
- Edelman NH, Kaplan RM, Buist S, Cohen AB, Hoffman LA and Kleinhenz ME, Snider GL and Speizer FE (1992) Chronic obstructive lung disease. *Chest* 102(3):243S-256S
- Edwards SL, Roberts C, McKean ME, Cockburn JS, Jeffrey RR and Kerr KM (2000) Pre-operative histological classification of primary lung cancer: accuracy of diagnosis and use of the non-small cell category. *Journal of Clinical Pathology* 53(7):537-540
- Egger M, Davey Smith G and Phillips AN (1997) Meta-analysis: principles and procedures. *British Medical Journal* 315(7121):1533-1537
- Egger M, Schneider M and Davey Smith G (1998) Spurious precision? Meta-analysis of observational studies. *British Medical Journal* 316(7125):140-144
- Eisen T, Matakidou A, Houlston R and GELCAPS Consortium (2008) Identification of low penetrance alleles for lung cancer: the GEnetic Lung CAncer Predisposition Study (GELCAPS). *BMC Cancer* 8: 244-254
- Ellis BG and Thompson MR (2005) Factors identifying higher risk rectal bleeding in general practice. *British Journal of General Practice* 55(521):949-955
- Fan J, Upadhye S and Worster A (2006) Understanding receiver operating characteristic (ROC) curves. *Canadian Journal of Emergency Medicine* 8(1):19-20
- Faraway JJ (2015) *Linear models with R* (2nd Edition). London: CRC Press
- Feld R, Ginsberg RJ, Payne DG and Shepherd FA (1995) Lung IN: Abeloff MD, Armitage JO, Lichter AS and Niederhuber JE (eds) *Clinical Oncology*. New York: Churchill Livingstone
- Ferguson GT, Enright PL, Buist AS and Higgins MW (2000) Office spirometry for lung health assessment in adults: a consensus statement from the National Lung Health Education Program. *Chest* 117(4):1146-1161
- Fergusson R, Gregor A, Dodds R and Kerr G (1996) Management of lung cancer in South East Scotland. *Thorax* 51(6):569-574
- Field A (2009) *Discovering Statistics Using SPSS*. London: SAGE Publication
- Field JK and Youngson JH (2002) The Liverpool Lung Project: a molecular epidemiology study of early lung cancer detection. *The European respiratory journal: official journal of the European Society for Clinical Respiratory Physiology* 20(2):464-479

References

- Field JK, Chen Y, Marcus MW, Mcronald FE, Raji OY and Duffy SW (2013) The contribution of risk prediction models to early detection of lung cancer. *Journal of Surgical Oncology* 108(5):304-311
- Finch E, Brooks D, Stratford PW and Mayo NE (2002) *Physical Rehabilitation Outcome Measures: a guide to enhanced clinical decision making* (2nd Edition). Ontario: Lippincott, Williams & Wilkins
- Fontana RS, Sanderson DR, Taylor WF, Woolner LB, Miller WE, Muhm JR and Uhlenhopp MA (1984) Early lung cancer detection: results of the initial (prevalence) radiologic and cytologic screening in the Mayo Clinic Study. *The American review of respiratory disease* 130(4):561-565
- Fontana RS, Sanderson DR, Woolner LB, Taylor WF, Miller WE and Muhm JR (1986) Lung cancer screening: the Mayo program. *Journal of occupational medicine* 28(8):746-750
- Garg AX, Hackam D and Tonelli M (2008) Systematic review and meta-analysis: when one study is just not enough. *Clinical journal of the American Society of Nephrology: CJASN* 3(1): 253-260
- GLOBOCAN (2008) *Lung Cancer Incidence and Mortality Worldwide in 2008*. Available from: <http://globocan.iarc.fr/factsheets/cancers/lung.asp> [Accessed 1st August 2013]
- Goff BA, Mandel LS, Melancen CH and Munz HG (2004) Frequency of symptoms of ovarian cancer in women presenting to primary care clinics. *Journal of the American Medical Association* 291(22):2705-2712
- GOLD, Global Initiative for chronic Obstructive Lung Disease (2004) *Global Strategy for the Diagnosis, Management, and Prevention of Chronic Obstructive Pulmonary Disease*. Available from: <http://www.goldcopd.org/guidelines-global-strategy-for-diagnosis-managementhtml> [Accessed April 2012]
- GOLD, Global Initiative for chronic Obstructive Lung Disease (2001) *Global Strategy for the Diagnosis, Management, and Prevention of Chronic Obstructive Pulmonary Disease National Institutes of Health, National heart, Lung, and Blood Institute*. Available from: www.goldcopd.com [Accessed April 2012]
- Greenland S (1989) Modeling and variable selection in epidemiologic analysis. *American Journal of Public Health* 79(3):340-349
- Grippi MA (1990) Clinical aspects of lung cancer. *Chest* 25(1):12-24
- Groves RM (1989) *Measurement Errors Associated with the Questionnaire IN: Groves RM (eds) Survey Errors and Survey Costs*. Hoboken, New Jersey: John Wiley & Sons, Inc.

- Guolo A and Varin C (2015) Random-effects meta-analysis: the number of studies matters. *Statistical Methods in Medical Research* 0(0): 1-19
- Hackshaw AK, Law MR and Wald NJ (1997) The accumulated evidence on lung cancer and environmental tobacco smoke. *British Medical Journal* 315(7114):980-988
- Hamilton W (2009) The CAPER studies: five case-control studies aimed at identifying and quantifying the risk of cancer in symptomatic primary care patients. *British Journal of Cancer* 101(Suppl2):S80-S86
- Hamilton W and Peters TJ (2007) Cancer Diagnosis in Primary individual risk prediction model for lung cancer. *British Journal of Cancer* 98(2):270-276
- Hamilton W and Sharp D (2004) Diagnosis of lung cancer in primary care: a structured review. *Family Practice* 21(6):605-611
- Hamilton W, Peters TJ, Round A and Sharp D (2005) What are the clinical features of lung cancer before the diagnosis is made? A population based case-control study. *Thorax* 60(12):1059-1065
- Hamilton W, Sharp DJ, Peters TJ and Round AP (2006) Clinical features of prostate cancer before diagnosis: a population-based, case-control study. *British Journal of General Practice* 56(531) 756-762
- Harden A (2010) Mixed-Methods Systematic Reviews: Integrating Quantitative and Qualitative Findings. Available from: http://ktdrr.org/ktlibrary/articles_pubs/ncddrwork/focus/focus25/Focus25.pdf [Accessed 12th October 2015]
- Hauptmann M, Pohlabeln H, Lubin JH, Jockel KH, Ahrens W, Briske-Hohlfeld I and Wichmann HE (2002) The exposure-time-response relationship between occupational asbestos exposure and lung cancer in two German case-control studies. *American journal of industrial medicine* 41(2):89-97
- Health and Social Care Information Centre (2006) *National Lung cancer Audit Report 2006*. Available from: <http://www.hscic.gov.uk/catalogue/PUB02715/clin-audi-supp-prog-lung-canc-nlca-2005-rep2.pdf> [Accessed April 2012]
- Health and Social Care Information Centre (2011) *National Lung Cancer Audit Report*. Available from: <http://www.hscic.gov.uk/catalogue/PUB16019/clin-audi-supp-prog-lung-nlca-2014-rep.pdf> [Accessed 18th April 2013]
- Henschke CI, McCauley DI, Yankelevitz D, Naidich DP, McGuinness G, Miettinen OS, Libby DM, Pasmantier MW, Koizumi J, Altorki NK and Smith JP (1999) Early lung cancer action project: overall design and findings from baseline screening. *Lancet* 354(9173):99-105

References

- Hippesley-Cox and J Coupland C (2011) Identifying patients with suspected lung cancer in primary care: derivation and validation of an algorithm. *British Journal of General Practice* 61(592):e715-723
- Hoppe R (1977) An analysis of 20,000 cases of suspected lung cancer. *Praxis und Klinik der Pneumologie* 31(10):872-884
- Huisman M (2008) Missing Data Analysis. Presentation Workshop Bath, UK: 9th September 2008
- Hurt CN, Roberts K, Rogers TK, Griffiths GO, Hood K, Prout H, Nelson A, Fitzgibbon J, Barham A, Thomas-Jones E, Edwards RT, Yeo ST, Hamilton W, Tod A and Neal RD (2013) A feasibility study examining the effect on lung cancer diagnosis of offering a chest X-ray to higher-risk patients with chest symptoms: protocol for a randomized controlled trial. *Trials* 14:405
- Hyde L and Hyde CI (1974) Clinical manifestations of lung cancer. *Chest* 65(3):299-306
- Iyen-Omofoman B, Tata LJ, Baldwin DR, Smith CJ and Hubbard RB (2013) Using socio-demographic and early clinical features in general practice to identify people with lung cancer earlier. *Thorax* 68(5):451-459
- Janssen-Heijnen ML and Coebergh JW (2003) The changing epidemiology of lung cancer in Europe. *Lung Cancer* 41(3):245-58
- Janssen-Heijnen MLG, Gatta G, Forman R, Capocaccia R and Coebergh JW (1998) Variation in survival of patients with lung cancer in Europe, 1985-1989. *European Journal of Cancer* 34(14):2191-2196
- Jemal A, Thomas A, Murray T, Samuels A, Ghafoor A, Ward E and Thun M (2003) Cancer statistics, 2003. *CA Cancer Journal for Clinicians* 53(1):5-26
- Jensen AR, Mainz J and Overgaard J (2002) Impact of delay on diagnosis and treatment of primary lung cancer. *Acta Oncology* 41(2):147-152
- Jones R, Latinovic R, Charlton J and Gulliford MC (2007) Alarm symptoms in early diagnosis of cancer in primary care: cohort study using General Practice Research Database. *British Medical Journal* 334(7602):1040
- Kardos P and Gebhardt T (1996) Chronic persistent cough in general practice: diagnosis and therapy in 329 patients over the course of 2 years. *Pneumologie* 50(6):437-441
- Kennedy TC, Proudfoot SP, Franklin WA, Merrick TA, Saccomanno G, Corkill ME, Mumma DL, Sirgi KE, Miller YE, Archer PG and Prochazka A (1996) Cytopathological Analysis of Sputum in Patients with Airflow Obstruction and Significant Smoking Histories. *Cancer Research* 56(20):4673-4678

- Kiri VA, Soriano JB, Visick G and Fabbri LM (2010) Recent trends in lung cancer and its association with COPD: an analysis using the UK GP Research Database. *Primary Care Respiratory Journal* 19(1):57-61
- Kornmann O, Beeh KM, Beier J, Geis UP, Ksoll M and Buhl R (2003) Global Initiative for Obstructive Lung Disease Newly diagnosed chronic obstructive pulmonary disease Clinical features and distribution of the novel stages of the Global Initiative for Obstructive Lung Disease. *Respiration* 70(1):67-75
- Koshiol J, Rotunno M, Consonni D, Pesatori AC, Matteis SD, Goldstein AM, Chaturvedi AK, Wacholder S, Landi MT, Lubin JH and Caporaso NE (2009) Chronic Obstructive Pulmonary Disease and Altered Risk of Lung Cancer in a Population-Based Case-Control Study *PLoS. ONE* 4(10):e7380
- Koyi H, Hillerdal G and Brandén E (2002) Patients' and doctors' delays in the diagnosis of chest tumours. *Lung Cancer* 35(1):53-57
- Kroenke K (2001) Studying Symptoms: Sampling and Measurement Issues. *Annals of Internal Medicine* 134(9):844-53
- Kubik A, Parkin DM, Khat M, Erban J, Polak J and Adamec M (1990) Lack of benefit from semi-annual screening for cancer of the lung: follow-up report of a randomized controlled trial of population of high-risk males in Czechoslovakia. *International journal of cancer* 45(1):26-33
- Kubik A, Zatloukal P, Boyle P, Robertson C, Gandini S, Tomasek L, Gray N and Havel L (2001) A case-control study of lung cancer among Czech women. *Lung Cancer* 31(2-3):111-122
- Kumar P and Clark M (2005) *Clinical Medicine* (6th Edition). London: Elsevier Saunders
- Lalkhen AG and McCluskey A (2008) Clinical tests: sensitivity and specificity. *Continuing Education in Anaesthesia, Critical Care & Pain* 8 (6): 221-223
- Landis SH, Murray T, Bolden S and Wingo PA (1998) Cancer statistics, 1998. *CA Cancer Journal for clinicians* 48(1):6-29
- Lee KJ and Carlin JB (2010) Multiple imputation for missing data: fully conditional specification versus multivariate normal imputation. *American Journal of Epidemiology* 171(5):624-632
- Levealahti H, Tishelman C and Ohlen J (2007) Framing the onset of lung cancer biographically: narratives of continuity and disruption. *Psycho-Oncology* 16(5):466-473
- Levy M L, Fletcher M, Price DB, Hausen T, Halbert RJ and Yawn BP (2006) International Primary Care Respiratory Group (IPCRG) Guidelines:

References

- diagnosis of respiratory diseases in primary care. *Primary Care Respiratory Journal* 15(1):20-34
- Liedekerken BM, Hoogendam A, Buntinx F, van der Weyden T and de Vet HC (1997) Prolonged cough and lung cancer: the need for more general practice research to inform clinical decision-making. *British Journal of General Practice* 47(421):505
- Little R and Rubin DB (1987) *Statistical analysis with missing data*. New York: Wiley
- Little RJA and Rubin DB (2002) *Statistical Analysis with Missing Data* (2nd Edition). Hoboken, New Jersey: Wiley
- Magnani C, Agudo A, Gonzalez CA, Andrion A, Calleja A, Chellini E, Dalmaso P, Escolar A, Hernandez S, Ivadi C, Mirabelli D, Ramirez J, Turuguet D, Usel M and Terracini B (2000) Multicentric study on malignant pleural mesothelioma and non-occupational exposure to asbestos. *British journal of cancer* 83(1):104-111
- Mallinckrodt CH (2013) *Preventing and Treating Missing Data in Longitudinal Clinical Trials*. New York: Cambridge University Press
- Mannino DM, Aguayo SM, Petty TL and Redd SC (2003) Low lung function and incident lung cancer in the United States: data from the First National Health and Nutrition Examination Survey follow-up. *Archives of International Medicine* 163:1475-1480
- Mannino DM, Homa DM, Akinbami LJ, Ford ES and Redd SC (2002) Chronic obstructive pulmonary disease surveillance- United States, 1971- 2000 Practice. *Journal of Thoracic Oncology* 4(11):1347-1351
- Mantel N and Haenszel W (1959) Statistical aspects of the analysis of data from retrospective studies of disease. *Journal of the National Cancer Institute* 22(4): 719-748.
- Marchenko YV and Eddings WD (2011) *A note on how to perform multiple-imputation diagnostics in Stata*. College Station, Texas: StataCorp
- Mason GA (1941) Cancer of the Lung. *Post-graduate Medical Journal* 17(191):153-156
- Mastrangelo G, Ballarin MN, Bellini E, Bizzotto R, Zannol F, Gioffre F, Gobbi M, Tessadri G, Marchiori L, Marangi G, Bozzolan S, Lange JH, Valentini F and Spolaore P (2008) Feasibility of a screening programme for lung cancer in former asbestos workers. *Occupational Medicine (Oxford, London)* 58(3):175-180
- Mathers CD and Loncar D (2006) Projections of global mortality ad burden of disease from 2002 to 2030. *PLoS medicine* 3(11):e442

- Mayne ST, Buenconsejo J and Janerich DT (1999) Previous lung disease and risk of lung cancer among men and women nonsmokers. *American Journal of Epidemiology* 149:13-20
- McDonald JH (2014) *Handbook of Biological Statistics* (3rd Edition). Baltimore, Maryland: Sparky House Publishing
- Melamed MR, Flehinger BJ, Zaman MB, Heelan RT, Perchick WA and Martini N (1984) Screening for early lung cancer Results of the Memorial Sloan-Kettering study in New York. *Chest* 86(1):44-53
- Menard S (2010) *Logistic Regression: From Introductory to Advanced Concepts and Applications*. Thousand Oaks, California: Sage Publications, Inc
- Mitchell ED, Rubin G and Macleod U (2013) Understanding diagnosis of lung cancer in primary care: qualitative synthesis of significant event audit reports. *British Journal of General Practice* 63(606):e37-46
- Moher D, Cook DJ, Eastwood S, Olkin I, Rennie D and Stroup DF (1999) Improving the quality reports of meta-analyses of randomised controlled trials: the QUOROM statement Quality of Reporting of Meta-analyses. *Lancet* 354(9193):1896-1900
- Molassiotis A, Wilson B, Brunton L and Chandler C (2010) Mapping patients' experiences from initial change in health to cancer diagnosis: a qualitative exploration of patient and system factors mediating this process. *European Journal of Cancer Care* 19(1):98-109
- Molenberghs G and Kenward MG (2007) *Missing data in clinical studies*. Chichester: John Wiley
- Molenberghs G, Thijs H, Jansen I, Beunkens C, Kenward MG, Mallinkrodt C and Carroll RJ (2004) Analyzing incomplete longitudinal clinical trial data. *Biostatistics* 5:445-464
- Moody A, Muers M and Forman D (2004) Delays in managing lung cancer. *Thorax* 59(1):1-3
- Mulka O (2005) NICE suspected cancer guidelines. *British Journal of General Practice* 55(517):580-581
- National Lung Screening Trial Research Team, Aberle DR, Adams AM, Berg CD, Black WC, Clapp JD, Fagerstrom RM, Gareen IF, Gatsonis C, Marcus PM and Sicks JD (2011) Reduced lung-cancer mortality with low-dose computed tomographic screening. *The New England journal of medicine* 365(5):395-409
- Neal RD, Robbé IJ, Lewis M, Williamson I and Hanson J (2014) The complexity and difficulty of diagnosing lung cancer: findings from a national primary-care study in Wales. *Primary health care research & development* 8:1-14

References

- NHS Choices (2010) NHS Choices Annual report 2010. Available from: <http://www.nhs.uk/aboutNHSChoices/professionals/developments/Documents/annual-report/annual-report-2010.pdf> [Accessed May 2012]
- NHS Choices (2013) *Angina*. Available from: <http://www.nhs.uk/conditions/angina/Pages/Symptoms.aspx> [Accessed 16th January 2015]
- NICE, National Institute for Health and Clinical Excellence (2005) *The diagnosis and treatment of Lung Cancer Methods, Evidence & Guidance*. London: National Collaborating Centre for Acute Care
- NICE, National Institute for Health and Clinical Excellence (2005) *Referral guidelines for suspected cancer*. London: National Collaborating Centre for Primary Care
- NICE, National Institute for Health and Clinical Excellence (2011) *Lung cancer. The diagnosis and treatment of lung cancer*. London: National Collaborating Centre for Cancer
- Norland EVT (1990) Controlling error in evaluation instruments. Available from: <http://www.joe.org/joe/1990summer/tt2.html> [Accessed 12th October 2015]
- O'Brien RM (2007) A Caution Regarding Rules of Thumb for Variance Inflation Factors. *Quality & Quantity* 41(5):673-690
- O'Rourke N and Edwards R (2000) Lung cancer treatment waiting times and tumour growth. *Clinical Oncology* 12(3):141-144
- O'Driscoll M, Corner J and Bailey C (1999) The experience of breathlessness in lung cancer. *European journal of cancer care* 8(1):37-43
- Office for National Statistics (1997) *Health Inequalities*. London: The Stationery Office
- Office for National Statistics (2009) *Survival Rates in England, patients diagnosed 2001-2006 followed up to 2007*. London: The Stationery Office
- Office for National Statistics (2012) *General Lifestyle Survey 2010*. London: ONS.
- Ott JJ, Ullrich A and Miller AB (2009) The importance of early symptoms recognition in the context of early detection and cancer survival. *European journal of cancer* 45(16):2743-2748
- Papi A, Casoni G, Caramori G, Guzzinati I, Boschetto P, Ravenna F, Calia N, Petruzzelli S, Corbetta L, Cavallesco G, Forini E, Saetta M, Ciaccia A and Fabbri LM (2004) COPD increases the risk of squamous histological subtype in smokers who develop non-small cell lung carcinoma. *Thorax* 59(8):679-681

- Parkin DM (1998) The global burden of cancer. *Seminars in cancer biology* 8(4):219-235
- Parkin DM, Bray F, Ferlay J and Pisani P (2005) Global cancer statistics, 2002. *CA: a cancer journal for clinicians* 55(2):74-108
- Parkin DM, Pisani P, Lopez AD and Masuyer E (1995) At least one in seven cases of cancer is caused by smoking: global estimates for 1985. *International Journal of Cancer* 59(4):494-504
- Patton MQ (2002) *Qualitative Evaluation and Research Methods* (3rd Edition). Thousand Oaks, California: Sage Publications, Inc
- Pauwels RA and Rabe KF (2004) Burden and clinical features of chronic obstructive pulmonary disease (COPD) .*Lancet* 364(9434):613-620
- Peake MD (2015) Should we be pursuing the earlier diagnosis of lung cancer in symptomatic patients?. *Thorax* 67:379-380
- Pearce N and Bethwaite P (1997) Social class and male cancer mortality in New Zealand 1984-7. *New Zealand medical journal* 110(1045):200-202
- Pepe M S (2004) *The statistical evaluation of medical tests for classification and prediction*. Oxford: Oxford University Press
- Petty TL (1996) The worldwide epidemiology of chronic obstructive pulmonary disease. *Current opinion in Pulmonary Medicine* 2(2):84-89
- Petty TL (1997) The *predictive* value of spirometry Identifying patients at risk for lung cancer in the primary care setting. *Postgraduate Medicine* 101(3):128-140
- Petty TL (1998) Definitions, causes, course, and prognosis of chronic obstructive pulmonary disease. *Respiratory Care Clinics of North America* 4(3):345-358
- Pilot D and Hunger B (1999) *Nursing research: principals and methods*. Philadelphia: Lippincott Williams & Wilkins
- Porter JC and Spiro SG (2000) Detection of early lung cancer. *Thorax* 55(Supp 1):56-62
- Potton E, Mc Caughan F and Jane S (2009) Chronic obstructive pulmonary disease and lung cancer. *Respiratory Medicine* 5(2):34-37
- Powell HA, Iyen-Omofoman B, Baldwin DR, Hubbard RB and Tata LJ (2013) Chronic obstructive pulmonary disease and risk of lung cancer: the importance of smoking and timing of diagnosis. *Journal of thoracic Oncology* 8(1):6-11

References

- Punturieri A, Szabo E, Croxton TL, Shapiro SD and Dubinett SM (2009) Lung cancer and chronic obstructive pulmonary disease: needs and opportunities for integrated research. *Journal of National Cancer Institute* 101(8):554-9
- Rabe KF, Hurd S, Anzueto A, Barnes PJ, Buist SA, Calverley P, Fukuchi Y, Jenkins C, Rodriguez-Roisin R, van Weel C and Zielinski J; Global Initiative for Chronic Obstructive Lung Disease (2007) Global strategy for the diagnosis, management, and prevention of chronic obstructive pulmonary disease: GOLD executive summary. *American Journal of Respiratory and Critical Care Medicine* 176(6):532-555
- Radhakrishna RB (2007) Tips for Developing and Testing Questionnaires/Instruments. *Journal of Extension* 45(1)
- Raghunathan TE and Bondarenko I (2007) *Diagnostics for multiple imputations*. Available from: http://papers.ssrn.com/sol3/papers.cfm?abstract_id=1031750 [Accessed 27th July 2014]
- Raghunathan TE, Lepkowski JM, Van Hoewyk J and Solenberger P (2001) A multivariate technique for multiply imputing missing values using a sequence of regression models. *Survey Methodology* 27:85-95.
- Ransohoff DF and Feinstein AR (1978) Problems of spectrum and bias in evaluating the efficacy of diagnostic tests. *The New England journal Medicine* 299(17):926-30
- Richards M (2007) EURO CARE-4 studies bring new data on cancer survival. *Lancet Oncology* 8(9):752-753
- Robins JM and Greenland S (1992) Identifiability and Exchangeability for Direct and Indirect Effects. *Epidemiology* 3(2):143-155
- Rogers TK (2006) Lung cancer diagnosis: time for a new approach. *Lung Cancer in Practice* 3(1):8-9
- Roggli VL, Vollmer RT, Greenberg SD, McGavran MH, Spjut HJ and Yesner R (1985) Lung cancer heterogeneity: a blinded and randomised study of 100 consecutive cases. *Human Pathology* 16(6):569-579
- Roth K, Nilsen TI, Hatlen E, Sørensen KS, Hole T and Haaverstad R (2008) Predictors of long time survival after lung cancer surgery: a retrospective cohort study. *BMC Pulmonary Medicine* 8:22
- Rothman KJ (1986) *Modern Epidemiology*. Toronto: Little, Brown and Company
- Rothman KJ (2002) *Epidemiology: an introduction*. Oxford: Oxford University Press

- Rothman KJ, Greenland S and Lash TL (2008) *Modern Epidemiology* (3rd Edition). Philadelphia: Wolter Kluwer Health/Lippincott Williams & Wilkins
- Royston R (2005) Multiple imputation of missing values: update. *Stata Journal* 5(4):527-536
- Rubin DB (1976) Inference and missing data. *Biometrika* 63(3):581-592
- Rubin DB (1987) *Multiple imputation for nonresponse in surveys*. New York: Wiley
- Rubin DB (1996) Multiple imputation after 18+ years. *Journal of the American Statistical Association* 91:473-489
- Rubin DB, Stern H and Vehovor V (1995) Handling “Don’t Know” survey responses: The Case of the Slovenian Plebiscite. *Journal of the American Statistical Association* 90(431):822-828
- Salomaa ER, Liippo K, Taylor P, Palmgren J, Haapakoski J, Virtamo J and Heionen OP (1998) Prognosis of patients with lung cancer found in a single chest x-ray screening. *Chest* 114(6):1514-1518
- Salomaa ER, Sallinen, Hiekkanen H and Liippo K (2005) Delays in the diagnosis and treatment of lung cancer. *Chest* 128(4):2282-2288
- Scagliotti G (2001) Symptoms, signs and staging of lung cancer. *European Respiratory Monograph* 17:86-119
- Schafer J (1999) Multiple imputation: A primer. *Statistical Methods in Medical Research* 8(1):3-15
- Schols AM, Soeters PB, Dingemans AM, Mostert R, Frantzen PJ and Wouters EF (1993) Prevalence and characteristics of nutritional depletion in patients with stable COPD eligible for pulmonary rehabilitation. *American Review of Respiratory Disease* 147(5):1151-1156
- Schwartz AG, Cote ML, Wenzlaff MPH, van Dyke A, Chen W, Ruckdeschel JC, Gadgeel S and Soubani AO (2009) Chronic obstructive lung diseases and risk of non-small cell lung cancer in women. *Journal of Thoracic Oncology* 4(3):291-299
- Scottish Cancer Registry, Information Services Division (ISD) (2013) *ISD Scotland, Practices and their Populations*. Available from: <http://www.isdscotland.org/Health-Topics/General-Practice/Practices-and-Their-Populations/> [Accessed 2nd April 2015]
- Scottish Executive (2002) *Scottish referral guidelines for suspected cancer*. Edinburgh: Scottish Executive, NHS Scotland

References

- Sekine Y, Katsura H, Koh E, Hiroshima K and Fujisawa T (2012) Early detection of COPD is important for lung cancer surveillance. *European Respiratory Journal* 39(5):1230-1240
- Sellers TA, Bailey-Wilson JE, Elston RC, Wilson AF, Elston GZ, Ooi WL and Rothschild H (1990) Evidence for Mendelian inheritance on the pathogenesis of lung cancer. *Journal of National Cancer Institute* 82(15):1272-1279
- Shahab L, Jarvis MJ, Britton J and West R (2006) Prevalence, diagnosis and relation to tobacco dependence of chronic obstructive pulmonary disease in a nationally representative population sample. *Thorax* 61(12):1043-1047
- Shapley M, Mansell G, Jordan JL and Jordan KP (2010) Positive predictive values of $\geq 5\%$ in primary care for cancer: systematic review. *British Journal of General Practice* 60(578):e366-77
- Shields PG (2000) Epidemiology of tobacco carcinogenesis. *Current oncology reports* 2(3):257-262
- Siafakas NM (2006) Definition and differential diagnosis of chronic obstructive pulmonary disease IN: Siafakas NM (ed) *European Respiratory Monograph 38: Management of Chronic Obstructive Pulmonary Disease*. UK: European Respiratory Society Journals Limited 1-6
- SIGN, Scottish Intercollegiate Guidelines Network (2005) Management of patients with lung cancer. Available from: <http://www.sign.ac.uk/pdf/sign80.pdf> [Accessed April 2012]
- Skillrud DM, Offord KP and Miller RD (1986) Higher risk of lung cancer in chronic obstructive pulmonary disease A prospective, matched, controlled study. *Annals of Internal Medicine* 105(4):503-507
- Sone S, Li F, Yang Z, Takashima S, Maruyama Y, Hasegawa M, Wang J, Kawakami S and Honda T (2000) Characteristics of small lung cancers invisible on conventional chest radiography and detected by population based screening using spiral CT. *The British Journal of Radiology* 73(866):137-145
- Sone S, Takashima S, Li F, Yang Z, Honda T, Maruyama Y, Hasegawa M, Yamanda T, Kubo K, Hanamura K and Asakura K (1998) Mass screening for lung cancer with mobile spiral computed tomography scanner. *Lancet* 351(9111):1242-1245
- Spiro GS, Gould MK, Colice GL and American College of Chest Physicians (2007) Initial Evaluation of the Patient with Lung Cancer: symptoms, Signs, Laboratory Tests, and Paraneoplastic Syndromes: ACCP Evidence-Based Clinical Practice Guidelines (2nd Edition). *Chest* 132(3 Suppl):149S-160S

- Spiro SG and Silvestri GA (2005) One Hundred Years of Lung Cancer. *American Journal of Respiratory and Critical Care Medicine* 172(5):523-529
- Stang P, Lydick E, Silberman C, Kempel A and Keating ET (2000) The prevalence of COPD: using smoking rates to estimate disease frequency in the general population. *Chest* 117(5 Suppl 2):354S-359S
- Stata (2013) *STATA USER'S GUIDE*. Texas: Stata Press
- Steenland K, Loomis D, Shy C and Simonsen N (2001) Review of occupational lung carcinogenesis. *American Journal of Industrial Medicine* 29(5):474-490
- Sterne JAC and Kirkwood B (2003) *Essential medical statistics* (2nd Edition). Malden, Mass: Blackwell Science
- Strauss GM (1997) Measuring effectiveness of lung cancer screening: from consensus to controversy and back. *Chest* 112(4 Suppl):216S-228S
- Strauss GM, Gleason RE and Sugarbaker DJ (1997) Screening for lung cancer Another look; a different view. *Chest* 111(3):754-768
- Summerton N (1999) *Diagnosing cancer in primary care*. Abingdon: Radcliffe Medical
- Summerton N (2002) Cancer recognition and primary care. *British Journal of General Practice* 62(474):5-6
- Sun S, Schiller JH and Gazdar AF (2007) Lung cancer in never smokers-a different disease. *Nature Reviews Cancer* 7:778-790
- Taylor JMG, Cooper KL, Wei JT, Sarma RV, Raghunathan TE and Heeringa SG (2002) Use of multiple imputation to correct for non-response bias in a survey of urologic symptoms among African-American men. *American Journal of Epidemiology* 156(8):774-782.
- The Joanna Briggs Institute (2014) *Joanna Briggs Institute Reviewers' Manual: Methodology for JBI Mixed Methods Systematic Reviews*. South Australia: The Joanna Briggs Institute
- Thomas JS, Lamb D, Ashcroft T, Corrin B, Edwards CW, Gibbs AR, Kenyon WE and Stephens RJ, Whimster WF (1993) How reliable is the diagnosis of lung cancer using small biopsy specimens?. *Thorax* 48(11):1135-1139
- Thompson SG (1994) Why sources of heterogeneity in meta-analysis should be investigated. *British Medical Journal* 309(6965): 1351-1355
- Tockman MS, Anthonisen NR, Wright EC and Donithan MG (1987) Airway obstruction and the risk of lung cancer. *Annals of Internal Medicine* 106(4):512-518

References

- Tockman MS, Erozan YS, Gupta P, Piantadosi S, Mulshine JL and Ruckdeschel JC (1994) The early detection of second primary lung cancers by sputum immunostaining LCEWDG Investigators Lung Cancer Early Detection Group. *Chest* 106(6 Suppl):385S- 390S
- Tockman MS, Gupta PK, Myres JD, Frost JK, Baylin SB, Gold EB, Chase AM, Wilkinson PH and Mulshie JL (1988) Sensitive and specific monoclonal antibody recognition of human lung cancer antigen on preserved sputum cells: a new approach to early lung cancer detection. *Journal of clinical oncology: official journal of the American Society of Clinical Oncology* 6(11):1685-1693
- Tod AM, Craven J and Allmark P (2007) Diagnostic delay in lung cancer: a qualitative study. *Journal of Advanced Nursing* 61(3):336-343
- Tokuhashi GK and Lilienfeld AM (1963) Familial aggregation of lung cancer in humans. *Journal of the National Cancer Institute* 30:289-232
- Tomatis L, Aitio A, Day NE, Heseltine E, Kaldor JM, Miller AB, Parkin DM and Riboli E (1990) *Cancer: causes, occurrence and control*. Lyon: International Agency for Research on Cancer IARC Scientific Publication 100
- Tomatis L, Kogevinas M, Pearce N, Susser M and Boffetta P (1997) *Poverty and lung cancer in social inequalities and cancer*. Lyon: International Agency for Research on Cancer IARC Scientific Publication 138:25-39
- Turner MC, Chen Y, Krewski D, Calle EE and Thun MJ (2007) Chronic obstructive pulmonary disease is associated with lung cancer mortality in a prospective study of never smokers. *American Journal of Respiratory Critical Care Medicine* 176(3):285-290
- van Buuren S, Boshuizen HC and Knook DL (1999) Multiple Imputation of Missing Blood Pressure Covariates in Survival Analysis. *Statistics in Medicine* 18(6):681-694
- Van den Eeden SK and Friedman GD (1992) Forced expiratory volume (1 second) and lung cancer incidence and mortality. *Epidemiology* 3:253-257
- Vestbo J and Lange P (2002) Can GOLD stage 0 provide information of prognostic value in chronic obstructive pulmonary disease?. *American Journal of Respiratory and Critical Care Medicine* 166(3):329-332
- Viegi G, Pistelli F, Sherrill DL, Maio S, Baldacci S and Carrozzi L (2007) Definition, epidemiology and natural history of COPD. *European Respiratory Journal* 30(5):993-1013
- Vittinghoff E, Glidden DV, Shiboski SC and McCulloch CE (2005) Regression Methods in Biostatistics: Linear, Logistic, Survival, and Repeated Measures Models. *The Stata Journal* 5(2):272-278

- Walter FM, Rubin G, Bankhead C, Morris HC, Hall N, Mills K, Dobson C, Rintoul RC, Hamilton W and Emery J (2015) Symptoms and other factors associated with time to diagnosis and stage of lung cancer: a prospective cohort study. *British Journal of Cancer* 112 Suppl:S6-S13
- Walters S, Maringe C, Coleman MP, Peake MD, Butler J, Young N, Bergström S, Hanna L, Jakobsen E, Kölbeck K, Sundstrøm S, Engholm G, Gavin A, Gjerstorff ML, Hatcher J, Johannesen TB, Linklater KM, McGahan CE, Steward J, Tracey E, Turner D, Richards MA and Rachet B; ICBP Module 1 Working Group (2013) Lung cancer survival and stage at diagnosis in Australia, Canada, Denmark, Norway, Sweden and the UK: a population-based study, 2004-2007. *Thorax* 68(6):551-564
- Wasswa-Kintu S, Gan WQ, Man SF, Pare PD and Sin DD (2005) Relationship between reduced forced expiratory volume in one second and the risk lung cancer: a systematic review and meta-analysis. *Thorax* 60(7):570-575
- Weiss ST, DeMeo DL and Postma DS (2003) COPD: problems in diagnosis and measurement. *European Respiratory Journal Supplement* 41:4s-12s
- Weller D, Vedsted P, Rubin G, Walter FM, Emery J, Scott S, Campbell C, Andersen RS, Hamilton W, Olesen F, Rose P, Nafees S, van Rijswijk E, Hiom S, Muth C, Beyer M and Neal RD (2012) The Aarhus statement: improving design and reporting of studies on early cancer diagnosis. *British Journal of Cancer* 106(7):1262-1267
- White IR, Daniel R and Royston P (2010) Avoiding bias due to perfect prediction in multiple imputation of incomplete categorical variables. *Computational Statistics and Data Analysis* 54(10):2267-2275.
- White IR, Royston P and Wood AM (2011) Multiple imputation using chained equations: Issues and guidance for practice. *Statistics in Medicine* 30(4):377-399.
- Whiting P, Rutjes AW, Reitsma JB, Glas AS, Bossuyt PM and Kleijnen J (2004) Sources of variation and bias in studies of diagnostic accuracy: a systematic review. *Annals of Internal Medicine* 140(3):189-202
- WHO (2011) World Health Organisation. Available from: <http://www.who.int/respiratory/copd/en/> [Accessed 27th January 2011]
- Wilcock A, Crosby V, Hughes A, Fielding K, Corcoran R and Tattersfield AE (2002) Descriptors of breathlessness in patients with cancer and other cardiorespiratory diseases. *Journal of Pain Symptom Manage* 23(3):182-9
- Willis GB (2005) *Cognitive interviewing: a tool for improving questionnaire design*. London: Sage Publications
- Wilson DO, Weissfeld JL, Balkan A, Schragin JG, Fuhrman CR, Fisher SN, Wilson J, Leader JK, Siegfried JM, Shapiro SD and Sciurba FC (2008) Association of

References

- radiographic emphysema and airflow obstruction with lung cancer. *American Journal of Respiratory and Critical Care Medicine* 178(7):738-744
- Wingo PA, Ries LA, Giovani GA, Miller DS, Rosenberg HM, Shopland DR, Thun MJ and Edwards BK (1999) Annual report to the nation on the status of cancer, 1973-1996, with a special section on lung cancer and tobacco smoking. *Journal of the National Cancer Institute* 91(8):675-690
- Wood A, White IR and Thompson SG (2004) Are missing outcome data adequately handled? A review of published randomised controlled trials. *Clinical Trials* 1(4):368-376
- Wu AH, Fontham ETH, Reynolds P, Greenberg RS, Buffler P, Liff J, Boyd P, Henderson BE and Correa P (1995) Previous lung disease and risk of lung cancer among lifetime nonsmoking women in the United States. *American Journal of Epidemiology* 141(11):1023-1032
- Yao X, Gomes MM, Tsoa MS, Allen CJ, Geddie W and Sekhon H (2012) Fine-needle aspiration biopsy versus core-needle biopsy in diagnosing lung cancer: a systematic review. *Current Oncology* 19(1):e16-e27
- Yoder LH (2006) An overview of lung cancer symptoms, pathophysiology, and treatment. *MEDSURG Nursing Journal: Cancer Caring and Conquering* 15(4):231-234
- Yoder LH (2006) Lung Cancer Epidemiology. *MEDSURG Nursing Journal: Cancer Caring and Conquering* 15(3):171-174
- Youden, WJ (1950) Index for rating diagnostic tests. *Cancer* 3(1):32-35
- Young RP, Hopkins RJ and Eaton TE (2007) Forced expiratory volume in one second: not just a lung function test but a marker of premature death from all causes. *European Respiratory Journal* 30(4):616-622
- Young RP, Hopkins RJ, Christma T, Black PN, Metcalf P and Gamble GD (2009) COPD prevalence is increased in lung cancer independent of age, gender and smoking history. *The European Respiratory Journal* 34(2):380-386
- Zheng T, Boffetta P and Boyle P (2011) *Epidemiology and Biostatistics*. Lyon, France: iPRI Scientific Publication. Available from: http://www.i-pri.org/wp-content/uploads/2013/12/Pages_IPRI_epidemiology-book.pdf [Accessed 16th June 2011]