

1 **Choosing where to set the threshold between low- and high-risk**
2 **patients: Evaluating a classification tool within a simulation**

3 Christina E Saville,^{a*} Honora K Smith,^b Katarzyna Bijak^c and Pauline
4 Leonard^d

5 *^aSchool of Health Sciences, University of Southampton, Southampton, UK;*

6 *^bSchool of Mathematical Sciences, University of Southampton, Southampton, UK;*

7 *^cSouthampton Business School, University of Southampton, Southampton, UK;*

8 *^dWhittington Health NHS Trust, UK*

9 **corresponding author C.E.Saville@soton.ac.uk*

1 **Evaluating a patient risk-classification tool within a simulation:**
2 **Poisson regression for generating patient characteristic combinations**

3 Health service providers must balance the needs of high-risk patients who require urgent
4 medical attention against those of lower-risk patients who require reassurance or less
5 urgent medical care. Based on their characteristics, we developed a tool to classify
6 patients as low- or high-risk, with correspondingly different patient pathways through a
7 service. Rather than choosing the threshold between low- and high-risk patients solely
8 considering classification accuracy, we demonstrate the use of discrete-event simulation
9 to find the best threshold from an operational perspective as well. Moreover, the
10 predictors in classification tools are often categorical, and may be inter-dependent.
11 Defining joint distributions of these variables from empirical data assumes that missing
12 combinations are impossible. Our new approach involves using Poisson regression to
13 estimate the joint distributions in the underlying population. We demonstrate our methods
14 on a practical example: setting the threshold between low- and high-risk patients with
15 proposed different pathways through a breast diagnostic clinic.

16 Keywords: simulation; health services; risk; credit scoring; regression

17 **1. Introduction**

18 Healthcare budgets constantly have to address the increased demand for the array of
19 diagnostic tests and treatments for patients with cancer because the successes of
20 research and development in care and diagnostics have led to more people living longer.
21 Health service providers must balance the needs of high-risk patients who require
22 urgent medical attention against those of lower-risk patients who require reassurance or
23 less urgent medical care. Classification methods, for example logistic regression or
24 decision trees, can be used to predict whether a patient is at low or high risk of a
25 disease. Then these patients can be routed along low- or high-risk pathways through
26 health services, recognising the different needs of these two patient groups.

27 When using these classification methods, managers must choose where to set the
28 threshold between low- and high-risk patients. Solely basing this decision on
29 classification accuracy (for example sensitivity or specificity measures), neglects the
30 operational impact of implementing risk-based pathways (for example waiting times or
31 resource use). In this paper, we propose using discrete-event simulation, in addition to

1 estimating classification accuracy, to help choose the threshold that provides the best
2 results from an operational perspective.

3 When evaluating risk-based patient management strategies in discrete-event simulation
4 models, the way in which patient characteristics are modelled requires careful
5 consideration. Patient characteristics affect not only the patient's risk group, but also the
6 patient's route through the simulation, as well as potentially their priority in queues and
7 service times. Patient characteristics are often categorical and inter-dependent. Some
8 possible combinations of characteristics may not appear in a data sample, but may exist
9 in the wider population. We propose an approach using Poisson loglinear (regression)
10 models for generating combinations of dependent categorical characteristics, allowing
11 generation of missing combinations.

12 We demonstrate our approach with a case study of a breast diagnostic clinic. The
13 classification tools developed, in this case logistic regression scorecards, are built from
14 data extracted from forms completed by general practitioners (GPs) referring patients to
15 the clinic. Currently all patients follow the same pathway through the clinic; we
16 investigate a proposal for patients classified as being at high risk of an abnormal result,
17 and thus needing diagnostic tests, to take a different pathway from low-risk patients.
18 Our method evaluates the classification tools in a simulation which routes patients along
19 low- and high-risk pathways through the clinic. Our aim is to find appropriate threshold
20 risk scores (cut-off scores) above which patients should be considered high-risk. As this
21 is a preliminary study, with data from a limited number ($n=179$) of patients, not all the
22 possible combinations of patient characteristics are present in the data. We therefore use
23 our method of Poisson regression to generate combinations of characteristics within the
24 simulation. We compare results in terms of clinic efficiency (proportion of time spent in
25 consultations or tests) and patients' total time at the clinic.

26 Although many researchers have applied operational research techniques to cancer care,
27 operational research studies addressing cancer diagnostic services are rare (Saville,
28 Smith, & Bijak, 2019), aside from tools for predicting cancer risk. Also, limited
29 research spans both primary and secondary cancer care services (Saville et al., *ibid*).

30 The main contributions of this paper from a theoretical and practical perspective are

- 1 • Combining classification and discrete-event simulation to find the best result
2 considering both predictive accuracy and operational performance
- 3 • Poisson regression for modelling combinations of characteristics that do not
4 necessarily appear in the sample
- 5 • Showing how GP referral information can be used to triage patients while still
6 giving all patients the chance to have tests if a clinician decides they are needed

7 The paper proceeds as follows. In Section 2, we present a brief overview of relevant
8 literature. In Section 3 we describe the healthcare background and the case study
9 setting. Section 4 details the approach used, including the development of the logistic
10 regression scorecards, a description of the simulation model and our method for
11 generating patient labels for simulation. In Section 5, we present the classification
12 accuracy and operational performance of different cut-off scores. In Section 6, we
13 discuss the contributions and limitations of the study, with future research directions.
14 The conclusion is given in Section 7.

15 **2. Literature review**

16 Over the last two decades, several authors have combined patient classification
17 techniques with discrete-event simulation (Bhattacharjee & Ray, 2016; Cannon et al.,
18 2013; Harper et al., 2003; Harper, 2002; Huang & Hanauer, 2016). Some of these
19 researchers used classification to model patient pathways accurately, by generating
20 groups of similar patients as simulation inputs (Harper, 2002) or predicting the
21 occurrence of health-related events during the simulation (Cannon et al., 2013; Harper
22 et al., 2003). On the other hand, Bhattacharjee & Ray (2016) classified patients using a
23 Classification and Regression Tree and then used discrete-event simulation to evaluate
24 the potential impact of sequencing appointments based on the patient classes. Huang &
25 Hanauer (2016) present a series of logistic regression models to predict no-shows. Each
26 model contains information about one more prior attendance, and can be used to decide
27 to what extent to overbook appointments. Here, discrete-event simulation was used to
28 evaluate the cost (waiting time plus overtime plus idle time) per patient for these
29 different models, and so to decide how many prior attendance variables should be
30 included. Unlike these papers, we describe how discrete-event simulation can help

1 choose the threshold score above which patients should be classified as high-risk.

2 Another relevant body of literature relates to setting thresholds in classification
3 algorithms when there are asymmetric costs (Pazzani et al., 1994, Zhao, 2008). This is
4 the case when misclassifying is worse in one direction than the other, so sensitivity may
5 be more important than specificity or vice versa. Many papers have built classification
6 models for predicting breast cancer, for example ((Ayer, Chhatwal, Alagoz, & Al, 2010;
7 Mangasarian, Street, & Wolberg, 1995; Pendharkar, Rodger, Yaverbaum, Herman, &
8 Benner, 1999). Unlike these, we are using classification models for predicting any kind
9 of abnormal result, including but not limited to cancer, from imperfect information (GP
10 referrals). This is to help with identifying which patients are likely to need imaging
11 tests. In our context, the cost of a missed abnormal result is also not as extreme as
12 missing a cancer case – these misclassified patients are sent to see a clinician who is
13 still able to send the patient for imaging tests (as today).

14 In simulation models of pathways through healthcare services, researchers have
15 modelled patient characteristics in three main ways (sometimes in combination). One
16 way is grouping patients with similar characteristics, either by specifying the probability
17 of belonging to each group (Bayer, Petsoulas, Cox, Honeyman, & Barlow, 2010;
18 Chemweno, Thijs, Pintelon, & van Horenbeek, 2014; Crawford, Parikh, Kong, &
19 Thakar, 2014), or by using group-specific arrival rates (Cooper, Davies, Roderick,
20 Chase, & Raftery, 2002; Gillespie et al., 2016; Monks et al., 2016). The relative
21 numbers of patients in each group are sometimes based on expert opinion (Chemweno
22 et al., 2014) or assumed to be the same as in data samples (Cooper et al., 2002;
23 Crawford et al., 2014). When the choice of groups is not obvious, patient data can be
24 analysed to find appropriate groups, for example Gillespie et al. (2016) group patients
25 with similar lengths of stay using Kaplan-Meier and log-rank tests. Elsewhere, different
26 clustering (Ceglowski, Churilov, & Wasserthiel, 2006; Isken & Rajagopalan, 2002) and
27 classification (Harper, 2002) techniques have been used to group similar patients. These
28 approaches provide insufficiently detailed characteristics for our situation.

29 A second way of modelling patient characteristics is inputting empirical data, either by
30 putting each real patient's information directly into the discrete-event simulation
31 (Eatock, Clarke, Picton, & Young, 2011; Khanna, Sier, Boyle, & Zeitz, 2016),
32 bootstrapping (Lord et al., 2013) or generating copies of each patient's set of

1 characteristics to compare different treatment strategies on the same cohort (Revankar et
2 al., 2014). The drawback of directly using empirical data is that it only includes those
3 patients seen in reality. This approach is therefore most suitable when the data sample is
4 deemed large enough to closely resemble the underlying population.

5 A third approach to this problem is using statistical distributions, either assuming
6 independence between characteristics (Burr et al., 2012; Crane, Kymes, Hiller, Casson,
7 & Karnon, 2013; Tran-Duy et al., 2014), or relationships between some characteristics
8 (Cooper et al., 2002; Lord et al., 2013; Pilgrim et al., 2008; Vataire et al., 2014; Wang et
9 al., 2017). These papers either use the empirical conditional distributions present in their
10 data or make assumptions about the relationships when data are unavailable. Using the
11 empirical conditional distribution relies on the relationships present in the sample being
12 representative of the wider population; combinations of characteristics not present in the
13 sample will not be simulated. In contrast, we propose an approach for generating
14 combinations of dependent, categorical characteristics, where combinations not present
15 in the data sample may be simulated.

16 **3. Healthcare background and case study**

17 Breast cancer is the most common cancer in the UK, making up 15% of new cancer
18 cases, with 99% of cases affecting women (Cancer Research UK, 2016b). Survival is
19 improving, with 85% of women diagnosed in England and Wales surviving the disease
20 for at least five years. However, the stage at which a cancer is detected greatly impacts
21 chances of survival, with only 26% of women with final stage disease surviving beyond
22 5 years (Cancer Research UK, 2018).

23 The most common route to breast cancer diagnosis is via referral by a General
24 Practitioner (GP) to a specialist diagnostic clinic, accounting for 60% of diagnoses
25 (Cancer Research UK, 2016a). These clinics are under strain; the covid-19 pandemic
26 has caused a backlog for cancer diagnostic services (Hanna, Aggarwal, Booth, &
27 Sullivan, 2020).

28 Currently, diagnostic clinics are organised as follows. A two-week wait target between
29 when a patient is referred and their attendance in clinic (Keogh, 2009) recognises the
30 urgency of confirming or eliminating a cancer diagnosis for both physical and mental
31 reasons. One-stop clinics are recommended, i.e. they should offer all necessary

1 diagnostic tests on a single day (Willett, Michell, & Lee, 2010). Two main options exist
2 for organising the sequence of services within the day. In some clinics, staff use
3 information provided by GPs on referral to identify those patients who should be sent
4 straight for imaging tests. The remaining patients are sent to see a clinician who decides
5 whether to request imaging. In other clinics, all patients see a clinician first. In this case,
6 the information provided by GPs may not be used at all.

7 All patients visiting breast diagnostic clinics will be worrying about the possibility of
8 cancer, although only a small proportion will receive a breast cancer diagnosis, for
9 example 4% of patients included in this study. Cancer is the most feared serious illness,
10 with women fearing breast cancer second most after brain cancer, according to a survey
11 commissioned by Cancer Research UK (2011). Clinic visits involve multiple stages,
12 meaning patients may wait multiple times, with little distraction from contemplating
13 their potential diagnosis. Thus, it is important that the proportion of time patients spend
14 in consultations and tests (where appropriate), as opposed to waiting between stages and
15 queuing, should be as high as possible. Moreover, it is critical that patients receive a
16 diagnosis confirming or excluding cancer as quickly as possible on the day of their
17 clinic visit.

18 The case study is based at the Whittington Health NHS Trust in North London, which
19 provides hospital and community services to a population of 500,000 in Islington,
20 Haringey, Barnet and Camden (Whittington Health NHS, 2019a). The Whittington
21 Hospital runs a one-stop clinic for diagnosis of patients with breast symptoms
22 (Whittington Health NHS, 2019b). In this clinic, all patients see a clinician first. We
23 model the potential operational impact of implementing risk-based pathways at this
24 clinic.

25 **4. Materials and Methods**

26 Our methods and data sources are outlined here; further details are available in the
27 supplemental material.

28 ***4.1 Data sources***

29 Between November 2015 and December 2016, patients were asked to fill in
30 questionnaires about the time they spent in different stages of their appointment (n=99).

1 This was complemented with a time and motion study where service and turnaround
2 times were measured by the PhD student.

3
4 Separately, between January and March 2016, we asked for patients' consent to use
5 their anonymised records to create a unique dataset linking GP referral information to
6 clinic tests and results. This dataset (n=179) was used both for developing the
7 scorecards and for generating patient labels to determine each patient's route through
8 the simulation.

9 ***4.2 Logistic regression scorecards to predict patient risk***

10 As classification tools, we develop two alternative logistic regression models that use
11 GP referral information to separate patients into groups at low and high risk of having
12 an abnormal result. For ease of interpretability, we transform the logistic regression
13 models into scorecards, using the "weights of evidence" approach common in credit
14 scoring applications (Thomas, 2009). Scorecards can be represented on an arbitrary
15 linear scale, are particularly suitable when the concept of risk is involved and make the
16 same predictions as the original logistic regression models. In this last respect
17 scorecards differ from the points system method proposed by Sullivan, Massaro, &
18 D'Agostino (2004), used for example as a breast cancer prediction rule (McCowan,
19 Donnan, Dewar, Thompson, & Fahey, 2011), where the points system provides an
20 approximation to the original model.

21
22 We considered the following seven commonly-reported characteristics for inclusion as
23 predictors in our model: "family history of cancer", "lump", "unilateral pain" (one-sided
24 pain), "other symptom", "urgency", "duration of symptoms" and "age". The
25 characteristic "other symptom" refers to rarer symptoms (those present in 15 or fewer
26 cases).

27
28 The outcome, a "normal" ($Y=1$) or "abnormal" ($Y=0$) diagnostic result, was derived
29 from patients' test results. A normal result means the patient has healthy breasts or does
30 not require imaging. Abnormal results cover both cancer and benign breast diseases,
31 including benign breast lumps such as cysts and fibroadenomas, infections (for example
32 mastitis and abscesses), and congenital problems, which cause the breast to have an

1 abnormal external appearance (Harvey, Down, Bright-Thomas, Winstanley, & Bishop,
 2 2014). Distinguishing between cancer and non-cancerous diseases without imaging is
 3 difficult (Harvey et al., 2014), which is why we propose routing patients at high risk of
 4 having an abnormal result to imaging first.

5
 6 A binary logistic regression model predicts the probability, p , of a normal result from
 7 patient-specific variables, X_1, X_2, \dots, X_n , obtained from GP referral information, as given
 8 in equation 1. The parameters $\beta_1, \beta_2, \dots, \beta_n$ show the relative importance of each
 9 characteristic in the prediction and β_0 is the intercept. These parameters are obtained
 10 from maximum likelihood estimates. In reality, there is some error, ε , that is not
 11 captured by the model.

$$12 \quad p := Prob(Y = 1|X_1, X_2, \dots, X_n) = \frac{1}{1 + e^{-(\beta_0 + \sum_{i=1}^n \beta_i X_i)}} \quad (1)$$

13 The weight of evidence, W , of a particular grouped attribute j , for example "young", of
 14 a characteristic i , for example "age", is the strength of evidence that patients who have
 15 the attribute will have an abnormal result. Let the number of normal (abnormal) results
 16 with attribute j be $n_j(a_j)$ and the total number of normal (abnormal) results be n_{total}
 17 (a_{total}). Then the weight of evidence variable, W_{ij} , is given by the following.

$$18 \quad W_{ij} = \ln\left(\frac{a_j}{a_{total}} \div \frac{n_j}{n_{total}}\right) \quad (2)$$

19 The scaled score, S , which is calculated from a scorecard, is related to the unscaled
 20 score, $\beta_0 + \sum_{i=1}^n \beta_i X_i$, which appears in the logistic regression, as follows.

$$21 \quad S = \left(\beta_0 + \sum_{i=1}^n \beta_i X_i\right) \cdot \text{factor} + \text{offset} \quad (3)$$

22 The factor and offset are the solution to the following system of linear equations. This
 23 follows since the unscaled score is equal to the log odds. The pair ($Score, Odds$) is the
 24 alignment point, e.g. it is assumed that the score 300 corresponds to odds of 12(:1) of

1 having an abnormal result. The *PDO* is the specified number of points to double the
 2 odds, e.g. if *PDO* is 20, then the odds double for every increase in 20 points.

$$3 \quad \text{Score} = \ln(\text{Odds}) \cdot \text{factor} + \text{offset} \quad (4)$$

$$4 \quad \text{Score} + \text{PDO} = \ln(2 \cdot \text{Odds}) \cdot \text{factor} + \text{offset} \quad (5)$$

5 Using the weights of evidence codings as the *X* variables, the scaled score for a
 6 particular patient becomes

$$7 \quad S = \left(\beta_0 + \sum_{i=1}^n \beta_i W_{ij} \right) \cdot \text{factor} + \text{offset} \quad (6)$$

8 where *n* is the number of variables and *j* is the attribute that this patient has for
 9 characteristic *i*.

10 The points of the scaled score can be split between the characteristics by dividing the
 11 parts of the score that are not characteristic-specific between them.

$$12 \quad \text{Points for characteristic } i = \left(\frac{\beta_0}{n} + \beta_i W_{ij} \right) \cdot \text{factor} + \frac{\text{offset}}{n} \quad (7)$$

13 Finally, to find the total score for a particular patient, the points for each of the patient's
 14 attributes are added up.

15

16 Models were developed in SAS Enterprise Miner 14.1 software. We developed two
 17 models: a full scorecard containing seven characteristics and a simple scorecard
 18 containing fewer characteristics. We selected variables for inclusion in the simple
 19 scorecard based on their information values, which measure how much each variable
 20 contributes to the abnormality prediction (Thomas, 2009). The simple scorecard
 21 contains only the two most predictive characteristics, “lump” and “age” which are
 22 strong and medium predictors of abnormal results, respectively (SAS, 2013). Values of
 23 the two continuous characteristics, “duration of symptoms” and “age” were grouped
 24 into attributes using the Interactive Grouping feature. This feature automatically
 25 generates groups using a decision tree algorithm aiming to maximise patient similarity
 26 (in terms of diagnostic results) within groups. For example, perfect similarity
 27 (technically, zero entropy) would mean that all patients in a group had the same

1 diagnostic result. Such grouping helps build more parsimonious and robust models that
2 can capture non-monotonic relationships. For the simple scorecard, the age groups were
3 adapted to 10-year brackets for ease of use.

4
5 Since the sample size is relatively small ($n = 179$ patients), the scorecards were
6 estimated using the entire dataset to make use of all the available data, rather than
7 removing some to use in validation. Instead, a bootstrapping technique (sampling with
8 replacement), implemented in Microsoft Excel, was used to internally validate the
9 models.

10 ***4.3 Simulation modelling***

11 The simulation models patients' visits to the clinic, from arrival to discharge (see Figure
12 1). Waits are omitted from the figure but are potentially present between each stage.
13 Patients who are predicted abnormal results (have scores at least as high as the cut-off
14 score) are sent straight for imaging tests; otherwise patients see a clinician first. In the
15 simulation, we ignored the small numbers of patients who are ineligible for the
16 scorecard (males and non-GP referrals), and assumed referral information would be
17 provided for all patients.

18 [Insert Figure 1 here]

19 Patients have the following imaging tests. It is assumed that *imaging first* patients (those
20 that are predicted an abnormal result) are given the same imaging tests as those patients
21 with actual abnormal results in our data, dependent on age (see Table 3 in the
22 Supplemental Material). For the *clinician first* patients (those patients predicted normal
23 results), we assume clinician behaviour in requesting tests remains unchanged from
24 current behaviour. That is, we assume the same test proportions as in the data,
25 dependent on age and actual result (see Tables 3 and 4 in the Supplemental Material). It
26 is assumed that patients with actual normal results never have a biopsy. Of those
27 patients with actual abnormal results who have an ultrasound, we assume 44% also have
28 a biopsy, as in our dataset.

29

1 Patients are prioritised for tests in the simulation in the following way. For
2 mammograms, patients who have had an ultrasound are prioritised in order of waiting
3 time. Then other patients are seen in order of their waiting time, regardless of whether
4 they have come from a consultation or straight for imaging. Similarly, for ultrasound,
5 patients who have had a mammogram are prioritised in order of waiting time. Then
6 other patients are seen in order of their waiting time, again regardless of whether they
7 visited a clinician or imaging first. In terms of prioritising patients to see a clinician, we
8 observed different prioritisation behaviours, so experimented with these in the
9 simulation. It was found that the choice of prioritisation rule has little effect on the
10 mean average wait for the initial consultation. However, the mean average waiting time
11 for a results consultation is about 12 minutes shorter if results consultations are
12 prioritised than if initial consultations are prioritised.

13 The objective of building the simulation model is to test alternative scenarios that differ
14 in the scorecard used (simple or full) and the cut-off score between low- and high- risk
15 patients. For the simple scorecard, the following scores are possible: 5, 10, 12, 14, 16,
16 21, 23, and 25. There are seven ways of dividing the scores into two groups at higher
17 and lower risk, for example using the cut-off scores 24, 22, 17, 15, 13, 11 and 6. It is
18 also possible to predict that all patients will have a normal result (for example with cut-
19 off 26), or that all patients will have an abnormal result (for example with cut-off 4).
20 Therefore, there are nine possible cut-off scores. For the full scorecard, possible scores
21 range from 167 to 299. We test cut-off scores at intervals of five.

22
23 At the start of the simulation, a set of initial patient labels is assigned to each new
24 patient. These patient labels are age divide, predicted result and actual result, which all
25 influence progress through the simulation of breast diagnostic services. Our novel
26 approach of using simulation to test the operational impact of alternative cut-off scores
27 involves changing the proportions of patients with abnormal and normal predicted
28 results between scenarios. Since the predicted result depends on the cut-off score, and

1 we want to test different cut-offs, we need to know how patients are likely to be
2 distributed between combinations of risk scores, age divide and actual result.

3 The age divide label takes one of two values: below 35 years old or at least 35 years old.
4 This label affects which imaging tests are required, as explained earlier. The predicted
5 result label (either normal or abnormal) is calculated from a scorecard and depends on
6 the cut-off score. When using the simple scorecard, only two referral characteristics are
7 needed to calculate the predicted result, but for the full scorecard, seven referral
8 characteristics are needed. The actual result label indicates whether the clinical
9 assessment and optional imaging tests show a normal or abnormal result.

10 In order to test the impact of using the simple scorecard to predict abnormal results and
11 route patients accordingly, we generate the patient labels in the following way. For the
12 179 patients in our sample, we calculated the age divide and actual result labels. From
13 the lump and age group characteristics, we calculated each patient's risk score according
14 to the simple scorecard. Each possible combination of the age divide, actual result and
15 risk score was present in our dataset. In the simulation, combinations of patient labels
16 are sampled from the empirical joint distributions (see Table 5 in the Supplementary
17 Material).

18 For testing the impact of using the full scorecard, we also need a joint distribution of the
19 age divide, actual result and predicted result for different cut-off scores. In this case, the
20 risk score is calculated from seven characteristics, including the age group. By splitting
21 the 29-42-year-old age group category at age 35, these amended age groups can also be
22 used to generate the age divide label value. Thus, we need the joint distribution of the
23 seven (amended) scorecard characteristics plus the actual result. There are 2880
24 possible combinations of levels (attribute scores), with only 166 unique combinations
25 present in our data sample (see Table 6 in Supplementary material for the number of
26 levels). Thus, using empirical proportions would not fully capture the likely joint
27 distribution present in the population of patients attending the breast diagnostic clinic,
28 so we use an alternative approach.

29 ***4.4 Poisson regression for generating patient labels in the simulation***

30 Instead, our approach is fitting a Poisson loglinear model (also known as Poisson

1 regression), a generalised linear model to predict counts in a contingency table with
2 these eight factors (Agresti, 2013). Next, the expected counts are converted to the
3 expected proportions of patients with each label combination, by dividing by the sample
4 size. In the simulation, when a new patient is generated, their initial label combination is
5 sampled from this distribution. A computational advantage is that only one sample is
6 drawn per patient, as opposed to one sample for each patient label.
7 We introduce some notation to represent Poisson loglinear models for ease of
8 exposition. A purely illustrative two-way contingency table for the factors *lump* (L) and
9 *urgency* (U) is shown in Table 0. Here μ_{LU} is the expected count for the cell in row i
10 and column j of the contingency table. For instance, μ_{00} is the expected number of
11 patients with no lump who are symptomatic. The Poisson loglinear model that includes
12 row effects, column effects and the row-column interaction can be summarised as
13 (L, U, LU) and is written in full as shown in Equation 8.

14 [Insert Table 0 here]

15

$$16 \quad \ln(\mu_{LU}) = \lambda + \lambda_1^L L_1 + \lambda_1^U U_1 + \lambda_2^U U_2 + \lambda_{11}^{LU} L_1 U_1 + \lambda_{12}^{LU} L_1 U_2 \quad (8)$$

17 In the above equation, λ is the offset parameter and dummy variables are used to code
18 the factor levels. Since *lump* has two levels, *lump present* and *no lump*, this is coded
19 using one dummy variable, L_1 . *Urgency* has three levels, *symptomatic*, *suspected*
20 *cancer*, and *other*, so is coded using two dummy variables U_1 and U_2 . The dummy
21 variable codings are given in Equations 9 to 11. The parameter λ_1^L represents the effect
22 of the level *lump present*, compared to *no lump*, on the expected count. Similarly λ_1^U and
23 λ_2^U estimate the effects of *suspected cancer* and *other* urgency compared to the reference
24 level, *symptomatic*. The dependence between *lump* and *urgency* is captured by the row-
25 column interaction effects, λ_{11}^{LU} and λ_{12}^{LU} .

$$26 \quad L_1 = \begin{cases} 0 & \text{if no lump,} \\ 1 & \text{if lump present} \end{cases} \quad (9)$$

$$1 \quad U_1 = \begin{cases} 0 & \text{if symptomatic or other,} \\ 1 & \text{if suspected cancer} \end{cases} \quad (10)$$

$$2 \quad U_2 = \begin{cases} 0 & \text{if symptomatic or suspected,} \\ 1 & \text{if other} \end{cases} \quad (11)$$

3 Alternative Poisson loglinear models for the same dataset differ in which interaction
 4 effects they include, and consequently also how many parameters must be estimated.
 5 Using dummy variables means that the number of parameters for each single variable
 6 effect is equal to the number of levels minus one. For each two-way interaction
 7 included, the number of parameters to estimate is equal to the product of the number of
 8 dummy variables for the two factors. Inclusion of higher-order interactions is also
 9 possible, but the feasibility of estimating the associated parameters depends on the size
 10 of the dataset (Agresti, 2013).

11
 12 We fitted two alternative models to predict counts in our 8-way contingency table (see
 13 Supplementary material Table 6) using the `glm()` command in the stats package in R.
 14 This command performs maximum likelihood estimation using iteratively reweighted
 15 least squares (Quick-R, 2017). We fitted firstly Model 1 which contained single variable
 16 effects only, (A, F, L, P, O, D, U, R), and secondly Model 2 which contained single
 17 variable effects, as well as all two-way interactions, ($A, F, L, P, O, D, U, R, AF, AL, AP,$
 18 $AO, AD, AU, AR, FL, FP, FO, FD, FU, FR, LP, LO, LD, LU, LR, PO, PD, PU, OD,$
 19 OU, OR, DU, DR, UR). Fitting Model 1 using dummy variables involved estimating 17
 20 parameters (for single effects) while Model 2 had 120 parameters (for both single effects
 21 and interactions).

22 To test how well the models fit the data, we simulated Pearson goodness-of-fit statistics
 23 for 3000 samples from the model distributions, in R software. This large number of
 24 iterations gave stable results. The simulated p-values are 0.39 for Model 1 and 0.59 for
 25 Model 2, so at the 0.05 level, we do not reject the null hypotheses that the models fit the

1 data. We want a model that provides good estimates of expected counts. Given that both
2 models fit well, we chose Model 2 because, in reality, characteristics are dependent; for
3 example, older women are more likely to have a lump. Table 7 in the Supplementary
4 Material shows the Model 2 joint distribution of label combinations for the full
5 scorecard with some different cut-off scores.

6

7 ***4.5 Measuring classification and operational performance of the scorecards***

8 We first consider the classification performance of the scorecards for different cut-off
9 scores. *True positives* (TP) and *true negatives* (TN) are patients who were correctly
10 classified with normal and abnormal results respectively. On the other hand, *false*
11 *positives* (FP) and *false negatives* (FN) are patients who were wrongly predicted normal
12 and abnormal results respectively. *Sensitivity* and *specificity* are the proportions of
13 normal and abnormal results that were correctly predicted, respectively. Finally,
14 *classification accuracy* is the proportion of all results that were correctly predicted.

15 Next we look at the operational performance of the simulated clinic when using the
16 scorecards with different cut-off scores. The goal is to maximise the clinic efficiency,
17 defined as the average proportion of patients' time at the clinic that is value-added, over
18 a set of cut-off score scenarios. Value-added activities are consultations and tests, as
19 opposed to waiting and queuing. In the case of ties, we prefer the cut-off score leading
20 to the lowest average patient total time spent in the clinic.

21

22 To define clinic efficiency and total time mathematically, first we define a patient's *start*
23 *time* as the time at which the clinic's performance for that patient begins to be measured.
24 For a patient arriving on time, the *start time* is the scheduled appointment time, which is
25 the same as the registration time. For a late patient, the *start time* is the registration time,
26 since the delay between the scheduled time and the registration time is caused by the
27 patient not the clinic. For early patients there are two possibilities. If an early patient is
28 seen early, their *start time* is the actual appointment time. If an early patient is seen late,

1 the *start time* is the scheduled appointment time. This corresponds to how waiting times
 2 for unpunctual patients are dealt with in the literature (see for example Santibáñez,
 3 Chow, French, Puterman, & Tyldesley, 2009). This allows us to define a patient's *total*
 4 *time* as the period from the *start time* until the *end time*, when the patient leaves the
 5 clinic.

$$6 \qquad \qquad \qquad \text{Total time} = \text{End time} - \text{Start time} \qquad (12)$$

7 For each day (run) of the simulation, we calculate the *average total time*. The mean
 8 *average total time* is the tie-breaker when clinic efficiency is equal for several
 9 scenarios.

10 We define the *value-added time* as the time during which a patient is in a consultation
 11 or having tests done.

$$12 \text{ Value-added time} = \text{Time in consultations} + \text{time in mammogram room} + \text{time in ultrasound room}$$

13
 14 Hence *efficiency* is the proportion of time at the clinic during which a patient is in a
 15 consultation or having tests done.

$$16 \qquad \qquad \qquad \text{Efficiency} = \frac{\text{Value-added time}}{\text{Total time}} \qquad (14)$$

17 The overall *clinic efficiency* is the average *efficiency* over all patients, so can be used as
 18 a performance measure on a particular day.

$$19 \qquad \qquad \qquad \text{Clinic efficiency} = \frac{\sum_{patients} \text{Efficiency}}{\text{Number patients}} \qquad (15)$$

20 Following the method suggested by (Banks, Carson II, Nelson, & Nicol, 2010), the
 21 simulation model was run many times to obtain a 95% confidence interval for the mean
 22 value of each operational performance measure. The trial run calculator feature of
 23 Simul8 was used to find the number of runs required for the 95% confidence limits to
 24 be within 10% of the estimate of the mean (the “precision”). Since each day is

1 independent with no patients staying overnight, the simulation was run from empty with
2 no warm-up period, and the run length was one day.

3 **5. Results**

4 **5.1 Scorecards**

5 The simple and full scorecards are shown in Tables 1 and 2. The scorecards work by
6 adding up the points corresponding to a patient's attributes (as recorded in their referral).
7 The higher the total risk score, the greater the chance that the patient will receive an
8 abnormal result.

9

10 [Insert Table 1 and 2 here]

11 **5.2 Classification performance for different cut-off scores**

12 Table 3 shows how well the simple scorecard separates normal and abnormal results in
13 the training data for each cut-off. The best classification accuracy is 68% and is
14 achieved when the cut-off is set at 17. (Currently, all patients are sent to a clinician first,
15 which corresponds to a cut-off of 4 and classification accuracy of 45%; that is, for 45%
16 of people abnormal results are found and imaging is needed). Both sensitivity and
17 specificity are greater than 60% when the cut-off is 15. Therefore, if the only
18 consideration is that sensitivity and specificity are equally important, then 15 would be
19 the best choice of cut-off. However, before finalising our choice of cut-off, we also need
20 to consider the classification performance of the full scorecard, as well as operational
21 measures for both scorecards.

22 [Insert Table 3 here]

23 For the full scorecard, there are a large number of possible scores, so in Table 4 we
24 present classification performance measures for a selection of cut-off scores only. Here,
25 the best classification accuracy among the cut-off scores considered is 70%, which is
26 achieved when the cut-off score is 220. This cut-off also achieves the best balance
27 between specificity and sensitivity, with both greater than 60%. The best classification
28 accuracy from the full scorecard offers a marginal improvement over the simple
29 scorecard (70% versus 68%), and the cut-off with the best balance between specificity
30 and sensitivity improves the specificity substantially (81% versus 72%) with a small

1 decrease in sensitivity (61% versus 62%).

2 [Insert Table 4 here]

3 ***5.3 Operational performance for different cut-off scores***

4 When choosing the cut-off score in this situation, the proportion of time patients are in
5 consultations/tests and how long they spend in the clinic are also of importance. Initial
6 experiments showed that under the current appointment schedule, with two clinicians
7 working simultaneously, patients arrive at the ultrasound area at a faster rate than they
8 can be processed. Therefore, we experimented with alternative set-ups. The best results
9 from preliminary simulation experiments (see Supplementary Material, Table 9 and 10)
10 were achieved using 15-minute gaps between appointments and one clinician working
11 at a time, so we use that set-up in our experiments below.

12

13 Simulation results using the simple scorecard and a subset of cut-offs are shown in
14 Table 5. The best clinic efficiency is 0.27 and is achieved with cut-off scores 15 (where
15 both sensitivity and specificity are higher than 60%), 17 (where the classification
16 accuracy is highest) and 22. The worst clinic efficiency is achieved when all patients are
17 sent straight to imaging (0% sensitivity and 100% specificity). The shortest average
18 total time is 107 minutes, when the cut-off score is 15, which corresponds to sending
19 53% of patients straight to imaging. A further benefit of using 15 as the cut-off would
20 be that it simplifies use of the scorecard, since it is equivalent to sending patients with a
21 lump recorded straight to imaging, and patients without a lump recorded to a clinician
22 first. A cut-off score of 15 was also the best choice to ensure that both sensitivity and
23 specificity were over 60%.

24 [Insert Table 5 here]

25 Next, we investigate whether by using the full scorecard we could further improve
26 clinic efficiency. The results are shown in Table 6. The highest average clinic efficiency
27 is 0.28 compared to 0.27 with the simple scorecard. This is achieved with a cut-off
28 score of 235, which corresponds to sending 32% of patients straight to imaging,
29 compared to 53% with the simple scorecard. This approach has classification accuracy
30 of 68% compared to 66% with the simple scorecard. The average total time at the clinic
31 is 113 minutes, slightly longer than for the simple scorecard (107 minutes). Since the
32 full scorecard is more complicated to use in practice, involving assessing seven

1 characteristics per patient rather than two, the simple scorecard is more promising for
2 practical use, particularly while referral forms are on paper rather than electronic.

3 [Insert Table 6 here]

4 **6. Discussion**

5 We have demonstrated the construction of a risk classification tool and shown how
6 discrete-event simulation could be used to guide decisions on where to set a cut-off
7 score between low- and high-risk patients. This enables the best cut-off to be chosen
8 from both classification accuracy and operational perspectives. Our unique approach
9 combining logistic regression for classification and simulation to select a cut-off score
10 allows decision makers to consider the wider implications of their choice of cut-off.
11 This approach is more versatile than considering solely predictive performance
12 measures; any measure that can be simulated can be used to compare cut-off scores. The
13 practical impact of the cut-off score may be more important than the predictive
14 accuracy, particularly where the classification model is being used to sequence or
15 prioritise services rather than deciding whether to offer a service.

16 The Poisson approach for generating combinations of categorical patient characteristics
17 that we propose in this paper provides a statistically sound solution to a general
18 problem. When there are many characteristics, small samples will not contain all the
19 combinations that may be present in the population. Poisson regression is a well-
20 established tool for modelling count data, which we apply to model counts of
21 combinations of characteristics. It enables the inclusion of inter-dependencies between
22 characteristics, and is appropriate for categorical data (by using dummy variables).
23 Unlike using empirical distributions, the Poisson approach is able to generate unseen
24 combinations.

25
26 Comparing our results to previous studies, we have generated a scorecard that allows
27 clinicians to add up each patients' risk score based on a small number of characteristics.
28 Alternative classification models would have provided results in a different format, but
29 the best model depends on the context. Classification and regression trees, used for
30 example by Bhattacharjee and Ray (2016) and Harper (2002) are useful in situations

1 where patients are grouped based on a continuous variable, e.g. service time or
2 operation time, rather than a binary variable, e.g. normal or abnormal result.

3 There were several limitations to this study which could point the way to future
4 research. When simulating the impacts of using scorecards to triage patients, we
5 focused our attention on just two operational performance measures: clinic efficiency
6 and average total time. The research could be extended by performing a cost-
7 effectiveness analysis and by considering additional performance measures, perhaps in a
8 weighted function. Another extension could be looking at the range in performance
9 across different patients and across different days, for example by calculating
10 percentiles rather than average measures.

11 A limitation of the Poisson approach is that the number of parameters to estimate
12 increases substantially when additional interactions are included. In our example we
13 were limited to including two-way interactions since there were insufficient data to
14 estimate three-way interactions. It would be useful to know in what situations applying
15 the Poisson approach is worthwhile and in which it makes little difference to estimates
16 of operational performance. Future work could compare recommendations obtained
17 from using the Poisson model distribution to the empirical distribution, for a series of
18 case studies. We suspect that situations where the rarer combinations of characteristics
19 correspond to higher service use will particularly benefit from the Poisson approach.

20

21 The problem under study can be generalised to other contexts outside diagnostic clinics,
22 and even outside healthcare: Is information provided by a non-specialist (or a patient or
23 customer themselves) complete and accurate enough to make decisions related to the
24 patient or customer (e.g. assign resources or allow access), without first performing a
25 specialist assessment? The classification-discrete-event simulation approach allows
26 different operational measures to be considered depending on the context.

27 ***7. Conclusion***

28 We have demonstrated the evaluation of classification tools within a discrete-event
29 simulation model and choice of a cut-off score based on operational performance

1 measures. Moreover, we have proposed the use of Poisson regression to generate patient
2 labels for simulation when data is limited.

3 A simple scorecard based on just the two most predictive patient characteristics, lump
4 and age, has the advantage of simplicity of use in the clinic, and improved accuracy and
5 efficiency compared to current practice. The full scorecard based on seven characteristics
6 only slightly improved accuracy and efficiency, but with more complexity in use. Using
7 a scorecard, patients (in the simulation) spend a higher proportion of their time in tests
8 and consultations, rather than waiting, compared to current practice. Also, overall, the
9 total time that patients spend at the clinic is reduced because high-risk patients have one
10 less consultation, being expedited straight to test, and waiting times are reduced.

11
12 The feasibility of using GP referral information to plan breast clinic diagnostics has also
13 been demonstrated for the first time to the best of our knowledge. Based on this
14 analysis, a larger-scale study is recommended to validate the accuracy and efficiency of
15 using a classification tool, namely a scorecard, to direct patients on risk-based pathways
16 through the clinic. As part of this, the best cut-off in different clinics and situations
17 should be compared to find out whether and how this varies, and patients should be
18 involved in discussions about the fairness of routing patients differently.

19
20 Note however that using a scorecard to route patients does not prevent them having
21 tests: if the scorecard misclassifies a patient as likely to have an abnormal result, they
22 may have tests that were unnecessary, and vice versa, if the scorecard misclassifies a
23 patient as likely to have a normal result, the clinician can correct the assessment and
24 send them for tests (as is the case for all patients today). Simulation offers a starting
25 point for the discussion as different scenarios can be compared without affecting real
26 patients. Some clinics already operate a split system with some patients being sent
27 straight to test and others to a clinician first. The potential benefits are better use of

1 resources (GP, clinician and imaging), as well as reductions in patients' non-value-
2 added time, such as on-the-day waiting and answering questions for a second time.
3
4 There has already been practical impact for the Whittington clinic from this research
5 project, which has benefitted from close collaboration with clinic staff from initial
6 stages onwards. Suggestions for a more efficient discharge system were tested in
7 preliminary simulation modelling and found to be efficient for all types of patients
8 (those with both normal and abnormal results). This change has already been fully
9 embedded by clinic staff. Practical suggestions on numbers of appointments offered
10 each day have also been made, and on balancing numbers of clinicians working at any
11 time with the availability of diagnostic tests.

12 **Acknowledgements**

13 We acknowledge all the patients who kindly allowed us to use their records in the study, as well
14 as the Whittington clinic and imaging staff for their help and patience. This article follows
15 STRESS-DES reporting guidelines (Monks et al., 2019).

16 Funding: This research was funded as part of an EPSRC PhD studentship from a block grant.
17 EPSRC had no role in the design of the study nor in the collection, analysis, and the
18 interpretation of data, nor in writing the manuscript.

19 Ethics approval and consent to participate: Ethical approval was granted by the London-
20 Bromley NRES Committee (reference number 15/LO/1335). The linked data used in this study
21 were obtained from patients who provided informed consent.

22 Availability of data and material: The dataset generated and analysed during the current study is
23 not publicly available because this was a condition of patient consent. However, an anonymised
24 version is available from the corresponding author on reasonable request.

25

1 **References**

- 2 Agresti, A. (2013). *Categorical data analysis* (3rd ed.). Hoboken, New Jersey: John
3 Wiley & Sons.
- 4 Ayer, T., Chhatwal, J., Alagoz, O., & Al, E. (2010). Comparison of logistic regression
5 and artificial neural network models in breast cancer risk estimation. *Informatics in*
6 *Radiology*, 30(1), 13–22.
- 7 Banks, J., Carson II, J. S., Nelson, B. L., & Nicol, D. M. (2010). Introduction to
8 simulation. In W. J. Fabrycky & J. H. Mize (Eds.), *Discrete-event system*
9 *simulation* (5th ed., pp. 1–22). New Jersey: Pearson Prentice Hall.
- 10 Bayer, S., Petsoulas, C., Cox, B., Honeyman, A., & Barlow, J. (2010). Facilitating
11 stroke care planning through simulation modelling. *Health Informatics Journal*,
12 16(2), 129–143. <https://doi.org/10.1177/1460458209361142>
- 13 Bhattacharjee, P., & Ray, P. K. (2016). Simulation modelling and analysis of
14 appointment system performance for multiple classes of patients in a hospital: A
15 case study. *Operations Research for Health Care*, 8, 71–84.
16 <https://doi.org/10.1016/j.orhc.2015.07.005>
- 17 Burr, J. M., Botello-Pinzon, P., Takwoingi, Y., Hernandez, R., Vazquez-Montes, M.,
18 Elders, A., ... Cook, J. (2012). Surveillance for ocular hypertension: An evidence
19 synthesis and economic evaluation. *Health Technology Assessment*, 16(29), 1–271.
- 20 Cancer Research UK. (2011). People fear cancer more than other serious illness.
- 21 Cancer Research UK. (2016a). Breast cancer diagnosis and treatment statistics.
22 Retrieved October 11, 2016, from [http://www.cancerresearchuk.org/health-](http://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/breast-cancer/diagnosis-and-treatment#heading-Zero)
23 [professional/cancer-statistics/statistics-by-cancer-type/breast-cancer/diagnosis-and-](http://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/breast-cancer/diagnosis-and-treatment#heading-Zero)
24 [treatment#heading-Zero](http://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/breast-cancer/diagnosis-and-treatment#heading-Zero)
- 25 Cancer Research UK. (2016b). Breast cancer statistics. Retrieved October 11, 2016,
26 from [http://www.cancerresearchuk.org/health-professional/cancer-](http://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/breast-cancer#heading-Zero)
27 [statistics/statistics-by-cancer-type/breast-cancer#heading-Zero](http://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/breast-cancer#heading-Zero)
- 28 Cancer Research UK. (2018). Breast cancer survival by stage at diagnosis.
- 29 Cannon, J. W., Mueller, U. A., Hornbuckle, J., Larson, A., Simmer, K., Newnham, J. P.,
30 & Doherty, D. A. (2013). Economic implications of poor access to antenatal care

- 1 in rural and remote Western Australian Aboriginal communities: An individual
2 sampling model of pregnancy. *European Journal of Operational Research*, 226(2),
3 313–324. <https://doi.org/10.1016/j.ejor.2012.10.041>
- 4 Ceglowski, R., Churilov, L., & Wasserthiel, J. (2006). Combining data mining and
5 discrete event simulation for a value-added view of a hospital emergency
6 department. *Journal of the Operational Research Society*, 58(2), 246–254.
7 <https://doi.org/10.1057/palgrave.jors.2602270>
- 8 Chemweno, P., Thijs, V., Pintelon, L., & van Horenbeek, A. (2014). Discrete event
9 simulation case study: Diagnostic path for stroke patients in a stroke unit.
10 *Simulation Modelling Practice and Theory*, 48, 45–57.
11 <https://doi.org/10.1016/j.simpat.2014.07.006>
- 12 Cooper, K., Davies, R., Roderick, P., Chase, D., & Raftery, J. (2002). The development
13 of a simulation model of the treatment of coronary heart disease. *Health Care
14 Management Science*, 5(4), 259–267. <https://doi.org/10.1023/A:1020378022303>
- 15 Crane, G. J., Kymes, S. M., Hiller, J. E., Casson, R., & Karnon, J. D. (2013).
16 Development and calibration of a constrained resource health outcomes simulation
17 model of hospital-based glaucoma services. *Health Systems*, 2(3), 181–197.
18 <https://doi.org/10.1057/hs.2013.5>
- 19 Crawford, E. A., Parikh, P. J., Kong, N., & Thakar, C. V. (2014). Analyzing discharge
20 strategies during acute care: A discrete-event simulation study. *Medical Decision
21 Making*, 34(2), 231–241. <https://doi.org/10.1177/0272989X13503500>
- 22 Eatock, J., Clarke, M., Picton, C., & Young, T. (2011). Meeting the four-hour deadline
23 in an A&E department. *Journal of Health Organization and Management*, 25(6),
24 606–624. <https://doi.org/10.1108/14777261111178510>
- 25 Gillespie, J., McClean, S., Garg, L., Barton, M., Scotney, B., & Fullerton, K. (2016). A
26 multi-phase DES modelling framework for patient-centred care. *Journal of the
27 Operational Research Society*, 67(10), 1239–1249.
28 <https://doi.org/10.1057/jors.2016.8>
- 29 Hanna, T. P., Aggarwal, A., Booth, C. M., & Sullivan, R. (2020). Counting the invisible
30 costs of covid-19: the cancer pandemic. *The BMJ Opinion*.
- 31 Harper, P. R. (2002). A framework for operational modelling of hospital resources.

- 1 *Health Care Management Science*, 5(3), 165–173.
2 <https://doi.org/10.1023/A:1019767900627>
- 3 Harper, P. R., Sayyad, M. G., De Senna, V., Shahani, A. K., Yajnik, C. S., & Shelgikar,
4 K. M. (2003). A systems modelling approach for the prevention and treatment of
5 diabetic retinopathy. *European Journal of Operational Research*, 150(1), 81–91.
6 [https://doi.org/10.1016/S0377-2217\(02\)00787-7](https://doi.org/10.1016/S0377-2217(02)00787-7)
- 7 Harvey, J., Down, S., Bright-Thomas, R., Winstanley, J., & Bishop, H. (2014). *Breast*
8 *disease management: A multidisciplinary manual*. Oxford: Oxford University
9 Press.
- 10 Huang, Y.-L., & Hanauer, D. A. (2016). Time dependent patient no-show predictive
11 modelling development. *International Journal of Health Care Quality Assurance*,
12 29(4), 475–488. <https://doi.org/10.1108/09526860710819440>
- 13 Isken, M. W., & Rajagopalan, B. (2002). Data mining to support simulation modeling
14 of patient flow in hospitals. *Journal of Medical Systems*, 26(2), 179–197.
- 15 Keogh, B. (2009). Operational Standards for the Cancer Waiting Times Commitments.
- 16 Khanna, S., Sier, D., Boyle, J., & Zeitz, K. (2016). Discharge timeliness and its impact
17 on hospital crowding and emergency department flow performance. *Emergency*
18 *Medicine Australasia*, 28(2), 164–170. <https://doi.org/10.1111/1742-6723.12543>
- 19 Lord, J., Willis, S., Eatock, J., Tappenden, P., Trapero-Bertran, M., Miners, A., ... Ruiz,
20 F. (2013). Economic modelling of diagnostic and treatment pathways in National
21 Institute for Health and Care Excellence clinical guidelines: The Modelling
22 Algorithm Pathways in Guidelines (MAPGuide) project. *Health Technology*
23 *Assessment*, 17(58), 1–150. <https://doi.org/10.3310/hta17580>
- 24 Mangasarian, O. L., Street, W. N., & Wolberg, H. (1995). Breast cancer diagnosis and
25 prognosis via linear programming. *Operations Research*, 43(4), 570–577.
- 26 McCowan, C., Donnan, P. T., Dewar, J., Thompson, A., & Fahey, T. (2011). Identifying
27 suspected breast cancer: Development and validation of a clinical prediction rule.
28 *The British Journal of General Practice*, 61(586), e205–e214.
29 <https://doi.org/10.3399/bjgp11X572661>
- 30 Monks, T., Currie, C. S. M., Onggo, B. S., Robinson, S., Kunc, M., & Taylor, S. J. E.
31 (2019). Strengthening the reporting of empirical simulation studies: Introducing

1 the STRESS guidelines. *Journal of Simulation*, 13(1), 55–67.
2 <https://doi.org/10.1080/17477778.2018.1442155>

3 Monks, T., Worthington, D., Allen, M., Pitt, M., Stein, K., & James, M. A. (2016). A
4 modelling tool for capacity planning in acute and community stroke services. *BMC*
5 *Health Services Research*, 16, 1–8. <https://doi.org/10.1186/s12913-016-1789-4>

6 Pazzani, M., Merz, C., Murphy, P., Ali, K., Hume, T., & Brunk, C. (1994). Reducing
7 misclassification costs. In *Proceedings of the Eleventh International Conference*
8 *on Machine Learning* (pp. 217–225). New Brunswick, NJ, USA.

9 Pendharkar, P. C., Rodger, J. A., Yaverbaum, G. J., Herman, N., & Benner, M. (1999).
10 Association, statistical, mathematical and neural approaches for mining breast
11 cancer patterns. *Expert Systems with Applications*, 17(3), 223–232.

12 Pilgrim, H., Tappenden, P., Chilcott, J., Bending, M., Trueman, P., Shorthouse, A., &
13 Tappenden, J. (2008). The costs and benefits of bowel cancer service
14 developments using discrete event simulation. *Journal of the Operational*
15 *Research Society*, 60(10), 1305–1314. <https://doi.org/10.1057/jors.2008.109>

16 Revankar, N., Ward, A. J., Pelligra, C. G., Kongnakorn, T., Fan, W., & LaPensee, K. T.
17 (2014). Modeling economic implications of alternative treatment strategies for
18 acute bacterial skin and skin structure infections. *Journal of Medical Economics*,
19 17(10), 730–740. <https://doi.org/10.3111/13696998.2014.941065>

20 Santibáñez, P., Chow, V. S., French, J., Puterman, M. L., & Tyldesley, S. (2009).
21 Reducing patient wait times and improving resource utilization at British Columbia
22 Cancer Agency’s ambulatory care unit through simulation. *Health Care*
23 *Management Science*, 12(4), 392–407. <https://doi.org/10.1007/s10729-009-9103-1>

24 SAS. (2013). *Developing Credit Scorecards Using Credit Scoring for SAS® Enterprise*
25 *Miner™ 13.1*. Cary, NC, USA.

26 Saville, C., Smith, H., & Bijak, K. (2019). Operational research techniques applied
27 throughout cancer care services: a review. *Health Systems*, 8(1), 52–73.
28 <https://doi.org/10.1080/20476965.2017.1414741>

29 Sullivan, L. M., Massaro, J. M., & D’Agostino, R. B. (2004). Presentation of
30 multivariate data for clinical use: The Framingham Study risk score functions.
31 *Statistics in Medicine*, 23(10), 1631–1660. <https://doi.org/10.1002/sim.1742>

- 1 Thomas, L. C. (2009). Using logistic regression to build scorecards. In *Consumer credit*
2 *models: Pricing, profit and portfolios* (1st ed., pp. 79–84). Oxford: Oxford
3 University Press.
- 4 Tran-Duy, A., Boonen, A., Kievit, W., van Riel, P. L. C. M., van de Laar, M. A. F. J., &
5 Severens, J. L. (2014). Modelling outcomes of complex treatment strategies
6 following a clinical guideline for treatment decisions in patients with rheumatoid
7 arthritis. *PharmacoEconomics*, *32*(10), 1015–1028.
8 <https://doi.org/10.1007/s40273-014-0184-4>
- 9 Vataire, A.-L., Aballea, S., Antonanzas, F., Hakkaart-van Roijen, L., Lam, R. W.,
10 McCrone, P., ... Toumi, M. (2014). Core discrete event simulation model for the
11 evaluation of health care technologies in major depressive disorder. *Value in*
12 *Health*, *17*(2), 183–195. <https://doi.org/10.1016/j.jval.2013.11.012>
- 13 Wang, H.-I., Smith, A., Aas, E., Roman, E., Crouch, S., Burton, C., & Patmore, R.
14 (2017). Treatment cost and life expectancy of diffuse large B-cell lymphoma
15 (DLBCL): A discrete event simulation model on a UK population-based
16 observational cohort. *European Journal of Health Economics*, *18*(2), 255–267.
17 <https://doi.org/10.1007/s10198-016-0775-4>
- 18 Whittington Health NHS. (2019a). About us. Retrieved August 9, 2019, from
19 <http://www.whittington.nhs.uk/default.asp?c=3920>
- 20 Whittington Health NHS. (2019b). Breast cancer. Retrieved August 9, 2019, from
21 <https://www.whittington.nhs.uk/default.asp?c=27104>
- 22 Willett, A. M., Michell, M. J., & Lee, M. J. R. (2010). *Best practice diagnostic*
23 *guidelines for patients presenting with breast symptoms*.
- 24 Zhao, H. (2008). Instance weighting versus threshold adjusting for cost-sensitive
25 classification. *Knowledge and Information Systems*, *15*(3), 321–334.
- 26
27

1 Table 0. Illustrative two-way contingency table of lump (L) and urgency (U)

	Column j	0	1	2
Row i		Symptomatic	Suspected cancer	Other
0	No lump	μ_{00}	μ_{01}	μ_{02}
1	Lump present	μ_{10}	μ_{11}	μ_{12}

2

3 Table 1. Simple scorecard

	Scorecard points
Age	
Age < 30	10
30 <= Age < 40	3
40 <= Age < 50	8
50 <= Age	12
Lump	
Yes	13
No or not recorded	2

4

5

1 Table 2. Full scorecard

	Scorecard points
Age	
Age < 29	43
29 <= Age < 42	16
42 <= Age < 47	33
47 <= Age < 52	21
52 <= Age	64
Duration of symptoms	
Not applicable or not recorded	29
Less than 2 weeks	42
2 weeks - 2 months	35
2 - 5 months	26
Over 5 months	35
Family history of cancer	
No or not recorded	29
Yes	38
Lump	
No or not recorded	5
Yes	54
Other symptom	
No or not recorded	29
Yes	36
Unilateral pain	
No or not recorded	31
Yes	33
Urgency	
Suspected cancer	32
Symptomatic	31
Other or not recorded	32

2

3

4

5

6

7

8

1

2 Table 3: Classification performance of simple scorecard for different cut-off scores

Cut-off score	True Positives	True Negatives	False Positives	False Negatives	Sensitivity	Specificity	Classification accuracy
26	98	0	81	0	100%	0%	55%
24	93	11	70	5	95%	14%	58%
22	83	31	50	15	85%	38%	64%
17	77	45	36	21	79%	56%	68%
15	61	58	23	37	62%	72%	66%
13	50	68	13	48	51%	84%	66%
11	37	72	9	61	38%	89%	61%
8	19	78	3	79	19%	96%	54%
4	0	81	0	98	0%	100%	45%

3

4

1 Table 4: Classification performance of full scorecard for different cut-off scores

Cut-off score	True Positives	True Negatives	False Positives	False Negatives	Sensitivity	Specificity	Classification accuracy
290	98	0	81	0	100%	0%	55%
280	96	3	78	2	98%	4%	55%
270	95	7	74	3	97%	9%	57%
260	94	15	66	4	96%	19%	61%
250	90	24	57	8	92%	30%	64%
240	85	37	44	13	87%	46%	68%
230	76	45	36	22	78%	56%	68%
220	60	66	15	38	61%	81%	70%
210	51	69	12	47	52%	85%	67%
200	38	73	8	60	39%	90%	62%
190	28	74	7	70	29%	91%	57%
180	16	79	2	82	16%	98%	53%
170	0	81	0	98	0%	100%	45%
160	0	81	0	98	0%	100%	45%

2

3

1 Table 5: Operational performance of simple scorecard for different cut-off scores

Cut-off score	Clinic efficiency		Average total time (minutes)	
	Number runs for 5% precision	Mean [95% confidence interval]	Number runs for 10% precision	Mean [95% confidence interval]
26	48	0.24 [0.23,0.25]	18	143 [129,156]
24	60	0.26 [0.24,0.27]	29	134 [121,147]
22	59	0.27 [0.25,0.28]	35	124 [111,136]
17	69	0.27 [0.26,0.29]	42	116 [105, 128]
15	106	0.27 [0.26,0.28]	47	107 [97,118]
13	112	0.26 [0.24, 0.27]	46	109 [98, 119]
11	132	0.25 [0.24,0.26]	52	110 [99,121]
6	150	0.23 [0.22,0.24]	65	116 [104,127]
4	166	0.2 [0.19,0.22]	61	123 [110, 135]

2
3
4
5

1 Table 6: Operational performance of full scorecard for different cut-off scores

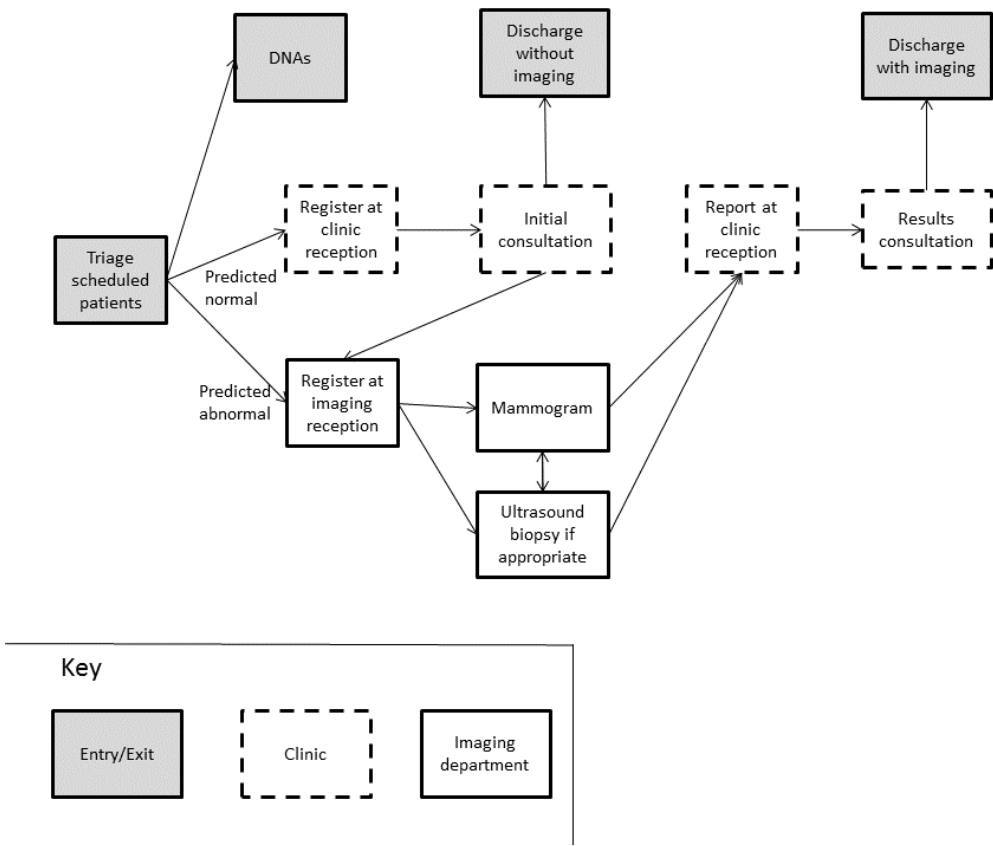
Cut-off score	Clinic efficiency		Average total time (minutes)	
	Number runs for 5% precision	Mean [95% confidence interval]	Number runs for 10% precision	Mean [95% confidence interval]
165	166	0.20 [0.19,0.22]	61	123 [110,135]
170	162	0.21 [0.20,0.22]	61	123 [111,135]
175	158	0.21 [0.20,0.22]	61	123 [111,135]
180	147	0.22 [0.21,0.23]	66	117 [105,128]
185	140	0.23 [0.22,0.24]	59	115 [104,127]
190	141	0.24 [0.23,0.25]	66	113 [102,124]
195	132	0.25 [0.24,0.26]	57	113 [102,125]
200	123	0.25 [0.24,0.26]	47	105 [95,116]
205	120	0.26 [0.24,0.27]	47	107 [96,117]
210	116	0.26 [0.25,0.27]	46	106 [95,116]
215	112	0.26 [0.25,0.27]	46	108 [97,118]
220	112	0.26 [0.25,0.28]	52	109 [98,119]
225	83	0.27 [0.25,0.28]	43	110 [99,121]
230	71	0.27 [0.26,0.29]	35	114 [103,126]
235	63	0.28 [0.26,0.29]	41	113 [102,124]
240	64	0.27 [0.26,0.29]	37	120 [108,131]
245	54	0.27 [0.26,0.28]	34	125 [112,137]
250	51	0.26 [0.25,0.28]	33	131 [118,144]
255	52	0.26 [0.25,0.28]	24	127 [115,139]
260	62	0.26 [0.25,0.27]	24	132 [119,144]
265	45	0.25 [0.24,0.27]	24	134 [121,147]
270	49	0.25 [0.24,0.27]	20	138 [124,151]
275	48	0.25 [0.23,0.26]	19	141 [127,155]
280	43	0.24 [0.23,0.25]	16	142 [128,156]
285	44	0.24 [0.23,0.25]	18	143 [129,156]
290	44	0.24 [0.23,0.25]	18	143 [129,156]
295	48	0.24 [0.23,0.25]	18	143 [129,156]
300	48	0.24 [0.23,0.25]	18	143 [129,156]

2

3

4

1 Figure 1: Simulation process map



2

3 *DNA=Did not attend*