

# Machine learning to predict early recurrence after oesophageal cancer surgery

Rahman, Saqib A BSc MRCS<sup>1</sup>

Walker, Robert C MRCS<sup>1</sup>

Lloyd, Megan A MMedSc BMBS<sup>1</sup>

Grace, Ben L BA<sup>1</sup>

van Boxel, Gijs I PhD FRCS<sup>2</sup>

Kingma, B.Feike MD<sup>2</sup>

Ruurda, Jelle P PhD MD<sup>2</sup>

van Hillegersberg, Richard PhD MD<sup>2</sup>

Harris, Scott MSc<sup>3</sup>

Parsons, Simon DM FRCS<sup>4</sup>

Mercer, Stuart DM FRCS<sup>5</sup>

Griffiths, Ewen A MD FRCS<sup>6</sup>

O'Neill, J.Robert PhD FRCSEd<sup>7</sup>

Turkington, Richard PhD MRCP<sup>8</sup>

Fitzgerald, Rebecca C PhD MRCP<sup>9</sup>

Underwood, Timothy J PhD FRCS<sup>1</sup>

On behalf of the OCCAMS Consortium, the full list of contributors is displayed in acknowledgements

1. Cancer Sciences Unit, University of Southampton, Southampton, UK
2. University Medical Centre, Utrecht, The Netherlands
3. Public Health Sciences & Medical Statistics Department, University of Southampton, Southampton, UK
4. Department of Surgery, Nottingham University Hospitals NHS Trust, Nottingham, UK
5. Department of Surgery, Portsmouth Hospitals NHS Trust, Portsmouth, UK
6. Department of Upper Gastrointestinal Surgery, University Hospitals Birmingham NHS Foundation Trust, Birmingham, UK
7. Cambridge Oesophagogastric Centre, Addenbrookes Hospital, Cambridge University Hospitals Foundation Trust
8. Centre for Cancer Research and Cell Biology, Queen's University Belfast, Belfast, UK
9. Hutchison/Medical Research Council Cancer Unit, University of Cambridge, Cambridge UK

## Keywords

oesophagus, adenocarcinoma, machine-learning, risk-prediction

Corresponding Author and Reprints: Professor T J Underwood

Email: [tju@soton.ac.uk](mailto:tju@soton.ac.uk)

Tel: +44 (0)2381206923

Fax: +44 (0)2381205152

#### Disclosures and Funding

The authors present no conflicts of interest.

TJU is supported by a Cancer Research UK and Royal College of Surgeons of England Advanced Clinician Scientist Fellowship, ID:A23924

OCCAMS2 was funded by a Programme Grant from Cancer Research UK (RG81771/84119)

#### Research Category

Original Article

#### Presentation

This research was presented at the meeting of the Association of Upper Gastrointestinal Surgeons of Great Britain and Ireland, Liverpool, 26/09/2019 and was the winner of the best Upper GI Cancer Parallel Papers Oral Presentation.

It has also been accepted for oral presentation to the BASO-ACS Annual Scientific Conference, London, 17<sup>th</sup> November 2019 in the Raven Proffered Prize Papers Session

#### Acknowledgements

Author contributions are listed in the attached files

Rogier van der Sluijs, UMC Utrecht, provided advice on predictive model methodology

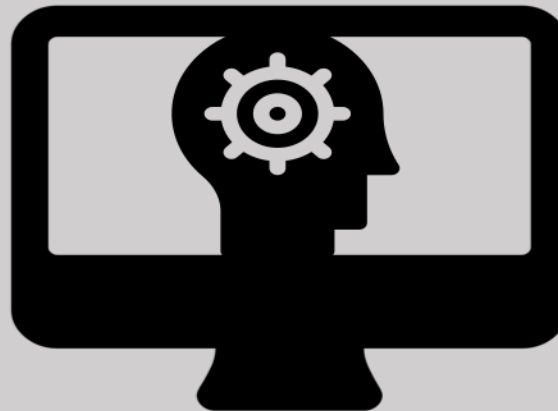
# Machine learning to predict early recurrence after oesophageal cancer surgery

International Cohort, Surgery after neoadjuvant treatment for oesophageal cancer



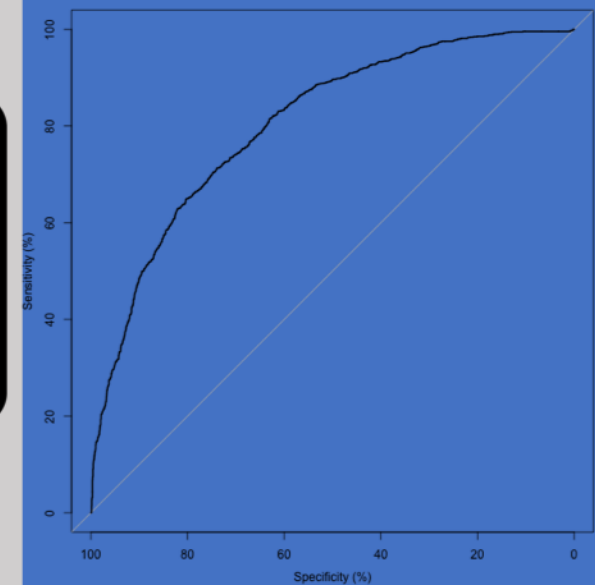
**812 Patients**

Methods = Elastic Net,  
Random Forest, XG Boost



**Outcome = Recurrence < 1year**

**AUC = 0.81**



Icons taken from [www.flaticon.com](http://www.flaticon.com), made by 'Freepik', 'smashicons', and 'prettycons'. Reproduced under creative commons attribution license

## **ABSTRACT**

### *Background*

Early cancer recurrence after oesophagectomy is a common problem with an incidence of 20-30% despite the widespread use of neoadjuvant treatment. Quantification of this risk is difficult and existing models perform poorly. This study aimed to develop a predictive model for early recurrence after surgery for oesophageal adenocarcinoma using a large multi-national cohort and machine learning approaches.

### *Methods*

Consecutive patients who underwent oesophagectomy for adenocarcinoma and had neoadjuvant treatment in 6 UK and 1 Dutch oesophago-gastric units were analysed. Using clinical characteristics and post-operative histopathology, models were generated using elastic net regression (ELR) and the machine learning methods random forest (RF) and XG boost (XGB). Finally, a combined (Ensemble) model of these was generated. The relative importance of factors to outcome was calculated as a percentage contribution to the model.

### *Results*

In total 812 patients were included. The recurrence rate at less than 1 year was 29.1%. All of the models demonstrated good discrimination. Internally validated AUCs were similar, with the Ensemble model performing best (ELR=0.791, RF=0.801, XGB=0.804, Ensemble=0.805). Performance was similar when using internal-external validation (validation across sites, Ensemble AUC=0.804). In the final model the most important variables were number of positive lymph nodes (25.7%) and **lymphovascular** invasion (16.9%).

### *Conclusions*

The derived model using machine learning approaches and an international dataset provided excellent performance in quantifying the risk of early recurrence after surgery and will be useful in prognostication for clinicians and patients.

## **INTRODUCTION**

Oesophageal adenocarcinoma carries a poor prognosis. Of the <40% of patients who are candidates for curative treatment(1), the 5-year survival rate remains approximately 25-50% in randomised trials(2–4) and rarely in excess of 50% in case series.

Early recurrence (less than 1 year) after surgery is a feared outcome with rates of 20-30%(3–5) frequently reported, despite the increasing uptake of neoadjuvant chemotherapy (NACT) and chemoradiotherapy (NACRT). This is particularly concerning because recovery from oesophagectomy is often long and the risk of major complications (Clavien-Dindo III-V) is as much as 31.1%(6). Many patients have not recovered from their primary cancer treatment when they experience cancer recurrence.

In an ideal setting prediction of early recurrence before embarking on a multimodal surgical pathway would provide the most useful information for patients and clinicians. However, staging information correlates poorly between pre- and post-operative settings(7), and genomic information is not yet able to predict outcome. Even the most robust preoperative models for prediction have an average performance at best(8). In contrast, postoperative information, although not able to influence surgical treatment decisions, is more prognostic and potentially informative to patients. It may also be helpful in decisions on the merits of adjuvant therapy, further refining the “high risk” group of patients where novel adjuvant treatments are currently being considered.

Naïve logistic regression (LR) has been the dominant approach to binary outcome prediction in clinical medicine for decades. Adoption of modern modified regression and 'machine learning' (ML) techniques has been limited, in part due to concerns over computational complexity and reliability. However, an increasing body of evidence demonstrates that they outperform traditional techniques in predictive performance(9,10), although this is debated(11). In part, the appeal of these approaches lies in their ability to model complex non-linear relationships which are common in cancer data, and which are challenging to model effectively with logistic/linear approaches. The increasing accessibility of software design now also allows the relatively straightforward deployment of these 'black-box' techniques.

Our group has previously published a multicentre UK cohort study which assessed survival according to Mandard Tumour Regression Grade (TRG)(12). This study included patients who had undergone oesophagectomy for adenocarcinoma of the oesophagus or gastro-oesophageal junction (GOJ) preceded by NACT as part of the Oesophageal Cancer Clinical and Molecular Stratification (OCCAMS) consortium. A clinically meaningful response to NACT was limited to TRG 1-2 only, which represented ~15% of patients. In the current study we set out to use this database, supplemented with an international cohort from the Netherlands, and machine learning techniques to develop and validate a clinically useful predictive model for early recurrence in oesophageal adenocarcinoma.

## **METHODS**

### *Ethics*

The OCCAMS consortium is a UK-wide multicentre consortium to facilitate clinical and molecular stratification of oesophagogastric cancer with ethical approval for biological sample collection and analysis in conjunction with detailed clinical annotation (Research Ethics Committee number: 10/H0305/1). Data collection and participation in research was approved by Institutional ethics committees at each OCCAMS site and UMC Utrecht.

### *Source of Data*

Data was sourced from 6 tertiary oesophago-gastric centres in the UK, as previously described(12). Briefly, the records of consecutive patients from each centre between 2000 and 2013 who underwent a planned curative oesophagectomy for adenocarcinoma and also received NACT (platinum-based triplet or cisplatin and 5-Fluorouracil) were reviewed and collated. Treatment was decided by a Multi-Disciplinary Team in individual institutions. Neoadjuvant treatment was considered for patients with locally advanced (cT2+) or node positive disease according to local and national guidelines. Clinical, pathological, recurrence and survival data were recorded. Data from one of the original centres was incomplete to the extent that modelling could not take place and was excluded *a priori*. In order to include NACT as a factor in the model further patients were identified from University Hospitals Southampton (UHS) and University Medical Centre Utrecht (UMCU), where CROSS type NACT(4) has been standard of care for oesophageal adenocarcinoma for a number of years. Patients who were deemed irresectable at the time of surgery or who had metastatic disease on the postoperative histology (i.e. pM1) were excluded from analysis.

The primary outcome measure was early recurrence, defined as confirmed local, regional or distant recurrence at less than 1 year from the date of surgery(5,8,13). Missing Data was treated as being missing completely at random and handled by list wise deletion. Modelling was based on a complete case analysis.

### *Predictor Characteristics*

Univariate statistics were calculated using Mann Whitney U and Chi-Square test for non-parametric data. The predictive models were generated on the whole dataset (n=812). All available variables were included in the analysis. The circumferential resection margin (CRM) was considered to be involved (and hence R1) in line with Royal College of Pathologists guidelines (i.e. CRM <1mm is positive)(14). Tumour grade and TRG(15) were assessed by dedicated gastrointestinal histopathologists who were blinded to clinical data. TRG was considered as responder (TRG1-2) vs non-responder (TRG 3-5) in line with our previous publication using this dataset(12). To increase the yield of information from lymph node data, both the number of positive lymph nodes and total lymph node harvest were considered as absolute number. For the regression model, linearity was assumed for continuous variables. Explicitly, the variables used to predict outcome were; gender, age, location of tumour, type of neoadjuvant therapy, response to neoadjuvant therapy (TRG), ypT, **lymphovascular** invasion, completeness of resection, grade of differentiation, number of positive lymph nodes and total number of lymph nodes examined.

### *Model Building and Validation*

We elected to use elastic-net regularized logistic regression (ELR)(16) along with two machine learning techniques; Random Forest (RF)(17) and Extreme Gradient Boosting

The EROC Model

(XGboost, XGB)(18). ELR applies a combination of the 'ridge' and 'lasso' penalties(19,20) with the benefits of both (partly minimisation of overfitting and variable selection). RF combines a specified number of decision trees (typically around 1000) created on random subsets of the dataset and is probably the most widely used machine learning approach in medical literature. XGB attempts to improve sequentially by generating models to explain where the original model fails and then repeating this process (typically around 1000 times), while simultaneously applying regularisation to minimise overfitting. Having generated individual models, we then combined them to generate overall predictions(21), an approach which theoretically is particularly beneficial when using diverse model types (such as those described above) that capture different elements of patients' risk profiles.

For ELR, the optimal alpha and lambda hyperparameters (penalty severities) were selected by grid-search using 10-fold cross validation with 5 repeats during model generation and 'log-loss' as the metric for optimisation. The RF model was derived from 1000 decision trees and hyperparameter tuning was conducted in a similar fashion (for number of variables per tree, split rule and minimum node size). The XGB model was again derived by cross validation of hyperparameters (number of optimisation rounds, maximum tree depth, minimum weight in each child node, minimum loss reduction (gamma), regularization penalty (eta) and subsampling for regularization). Full details of hyperparameter tuning is given in the supplementary materials (S7). These three models were then combined to generate the final (ensemble) model by generating a linear blend of predicted probabilities using logistic regression.

Discrimination of the models was assessed using the area under the receiver operator characteristic (ROC) curve (AUC). In the context of this paper, if two random patients were selected, one with a recurrence of cancer at less than 1 year and one disease free at 1 year, the AUC is equivalent to the probability the model will score the patient with recurrence higher than the patient without. Internal validation was performed using 0.632 bootstrapping, with 1000 resampled datasets. Bootstrapping was preferred for internal validation over splitting the cohort into derivation and validation sets, as this has been shown to reduce bias and improve overall model performance, particularly with moderately sized datasets(22–24). Calibration was assessed visually and formally with the Hosmer-Lemeshow Test. As our dataset contains multiple centres with small numbers of patients, we also opted for an internal-external validation procedure, as advocated by Steyerberg and Harrell(25). This entails generating models on all centres apart from one and validating the model on the remaining centre. This process is then repeated leaving each centre out sequentially and an average calculated. This method demonstrates how the model performs in external data while also allowing the whole dataset to be used for training.

Unadjusted tree models (such as RF, which is included in the Ensemble) and other maximum margin methods typically calibrate poorly as a consequence of their methodology, with predicted probabilities biased towards the centre. To allow meaningful interpretation of probability, Isotonic regression was used to scale probabilities on the final model, as has been previously described(26,27).

In contrast to logistic regression, assessing global variable importance is challenging using machine learning techniques and to an extent they are ‘black-boxes’. As coefficients, as

would be seen in a logistic regression, an alternative method is required. We used the 'VarImp' function of the caret R package, where ROC curves are generated for the outcome for each individual predictor and the contribution to the global ROC curve calculated as a percentage. Due to the nature of higher-order interactions present in the model, variable importance in individual predictions must be calculated independently. We calculated the average marginal contribution of each variable (change from the mean prediction i.e. the Shapley value(28)) for individual predictions. A similar approach was used by Nanayakkara et al. for analysing in-hospital mortality following cardiac arrest(29).

Data analysis was conducted using R (Version 3.5.3, The R Foundation for Statistical Computing). Models were trained using the 'caret'(30) and 'caretEnsemble'(31) packages. Individual variable importance was calculated using 'iml'(32). All are available at <https://CRAN.R-project.org/>. Full R code to train the models as described is given in the supplementary materials (S7), along with a list of packages used.

The calibrated final model was designed using R Shiny(33) and is available freely at: <https://uoscancer.shinyapps.io/EROC/>. No data entered into the model is collected or stored.

## **RESULTS**

A total of 812 patients from 7 centres were included in model training. A consort diagram detailing patient numbers and the final sample size is presented in Figure 1.

Most patients were male (84.6%), with a median age of 64 years. The majority of tumours were at the GOJ (55.5%), with a high proportion of locally advanced (ypT3-4 – 66.8%) and node positive disease (61.0%). First recurrence of cancer within 1 year of surgery was identified in 236 patients (29.1%). The early recurrence group were significantly less likely to have responded to neoadjuvant treatment (8.5% vs 21.7%), and had worse ypT, ypN, **lymphovascular** invasion, R1 resection rate and grade of differentiation (all  $p < 0.001$ ). Detailed group clinicopathological data is shown in Table 1.

### *Model performance (discrimination)*

Discrimination was assessed in the training set, internally (via bootstrapping) and internally-externally (across centres). ROC curves of the internal validation of each model are shown in Figure 2. All models demonstrated excellent discrimination on the training set (apparent discrimination), with the Random Forest Model performing the best (AUC 0.98), followed by the Ensemble model (0.90), XGB (0.85) and ELR (0.81). On internal validation, the Ensemble model had the best performance (AUC 0.81) and the ELR the worst (AUC 0.79). Overall discrimination for each model is summarised in table 2.

### *Model performance (calibration)*

Calibration on the training set was visually best in the ELR, and worst in the RF and Ensemble models (supplementary materials). This was corroborated by the Hosmer Lemeshow test (p value ELR=0.806, RF=<0.001, XGB=0.030, Ensemble=<0.001). Probabilities generated by the final model were scaled using isotonic regression. Calibration before and after scaling is shown in Figure 3 (shaded area represents two standard errors, calibration tables in supplementary materials). The Hosmer Lemeshow test before scaling gives a Chi2 of 38.0 and  $p<0.001$ , and after gives Chi2 of 4.5 and  $p=0.806$ . Similarly, the Brier score (a measure of overall model performance) also improves from 0.119 to 0.114.

### *Variable Importance*

Coefficients and odds-ratios cannot be generated for these models. We therefore computed variable importance as a percentage contribution to the model. The results are displayed in Table 3.

Overall the most influential predictor variable is number of positive lymph nodes (25.7%), followed by **lymphovascular** Invasion (16.9%). There is considerable variability in importance across models. For example, Age contributes 0.3% to the ELR model, 18.2% to the RF model, 10.2% to the XGB model and 9.6% to the Final model.

It is important to restate that the relationships between the variables and outcome are non-linear and their importance varies considerably according to other variables due to higher order interactions. As an example, even though lymph node status is the most influential marker overall, there are combinations of other variables that would make other variables most important in individual patients. To illustrate this and demonstrate how variables interact, three example patients are considered below. The technique used measures the change in the prediction from the mean prediction (27.1%) that can be attributed to each predictor variable. This approach (calculation of the Shapley value) originates from cooperative game theory.

#### **Example 1: Low Risk Patient (AJCC ypTONOMO: Stage 1)**

*50 year old Male with a G0J adenocarcinoma who undergoes neoadjuvant chemoradiotherapy. On postoperative pathology he is a responder with ypT0, negative **lymphovascular** invasion, R0 resection and a well differentiated tumour. He has 0 positive lymph nodes out of 30 sampled.*

The EROC Model

**Example 2: Medium Risk Patient (AJCC ypT3N0M0: Stage 2)**

66 year old Male, with an Oesophageal adenocarcinoma who undergoes neoadjuvant chemoradiotherapy. On postoperative pathology he is a non-responder with ypT3, positive **lymphovascular** invasion, R0 resection and a moderately differentiated tumour. He has 0 positive lymph nodes out of 30 sampled.

**Example 3: High Risk patient (AJCC ypT3N2M0: Stage 3b)**

70 year old Female with an Oesophageal adenocarcinoma who undergoes neoadjuvant chemotherapy. On postoperative pathology she is a non-responder with ypT3, positive **lymphovascular** invasion, R1 resection, Poor differentiation with 5 positive lymph nodes out of 30 sampled.

## **DISCUSSION**

In this study we have derived an easy to use and robust clinical model for predicting the risk of early recurrence after surgery for oesophageal adenocarcinoma. It uses routinely collected clinical and pathological data which should be available for every patient, which together allow considerably more precision in risk estimation than would be possible using individual variables which are known to be influential such as pathological lymph node involvement. The final model demonstrated excellent discrimination, and validation techniques supported the generalisability of the approach.

In addition to prognostication, this model may be useful for planning adjuvant therapy. Early recurrence after oesophagectomy, often before recovery from surgery is complete, is a devastating outcome for patients. Targeting existing and emerging treatment combinations in this patient group to prolong time to recurrence or prevent recurrence is vital, however can only happen with accurate predictions of the likelihood of relapse. **The starting point for the consideration of treatment escalation or novel combinations (e.g. immunotherapy) after surgery is the identification of patients who are at high risk of recurrence.** We have purposefully avoided dichotomization/stratification based on outcome and presented raw probability in preference to this. This will allow full discussions between surgeons/oncologists and patients to take place regarding the benefits of adjuvant therapy and tailored to individual patient's post-operative recovery and wishes. It may also allow stratification of adjuvant trials based on layered levels of risk.

This cohort exhibited an early recurrence rate of 29.1%, which is similar to previous reports where this outcome was explicitly specified(3–5,8). There was also an R1 resection rate of 29.1%, in line with previously reported data(34,35) with an RCP definition of CRM positivity (CRM<1mm involved). On univariate analysis all factors expected to correlate with worse prognosis (including ypT, ypN, **lymphovascular** invasion, R1 resection and grade of differentiation) were significantly worse in those patients who developed an early recurrence. This validates our cohort as a true representation of contemporary practice and a sensible place to begin building more complex models.

Discrimination of the different models was similar, with minimal variability of AUC between models on validation. However, the ensemble model consistently performed the best and is a suitable choice for the final model. The decline in performance from the training set to validation, which was particularly marked in the RF and ensemble models, is a consequence of the tuning process, whereby the optimum values are chosen from a grid of thousands after repeated tests (in this case repeated 10-fold cross validation). In this setting, the apparent performance of the model on the training set is over-estimated and should be disregarded.

There was marked heterogeneity in variable importance between models. This is interesting, particularly in the context of the models performing so similarly overall and supports the idea of combining them to capture different patient information. The most important variables overall were number of positive lymph nodes and **lymphovascular** invasion, accounting for 42.6% of performance. This is not only biologically sensible, but the subject of several recent publications and ongoing translational work(12,36,37). Although

not available for this study, more detail regarding lymphadenopathy – e.g. downstaging and anatomical location would likely be informative. Firm conclusions regarding variables are difficult considering the nature of the study. However, we would draw attention to two facets of the model. Firstly, TRG was the least influential variable across the board, with an importance of almost 0%. This suggests that in itself TRG adds no information over the other measured variables in predicting early outcomes. This is in keeping with emerging data regarding the genomic disparity between primary tumours and their metastasis (lymph node or distant) and our previous report of the importance of lymph node downstaging to clinical outcome(12,38). Secondly, modality of treatment was the third most important determinant of outcome, with NACT conferring an advantage over NACRT. In this cohort, despite having considerably **more favourable** postoperative pathology after NACRT, the rate of early recurrence was no less, and borderline higher (**NACRT 35.5%, NACT 27.5%, p = 0.061, Supplement 4**). This suggests that although there is **more favourable** post-operative pathology seen with NACRT, this does not translate to better outcome(39–41) and hence a ypT3N1R0 after NACT does not have the same meaning as a ypT3N1R0 result after NACRT, at least in the early period after treatment. This is important in postoperative discussions with patients. **As the machine learning approaches detailed here allow interactions between variables, the model suggests that NACRT confers a greater risk – but this increased risk is conditional on the other variables being static rather than an overall increase in risk from having NACRT.**

**To further explore this, details of recurrence location (i.e. loco-regional vs distant) would be informative, however due to the historical nature of the majority of the patients (data collected for the first study) we were unable to reliably ascertain this for the majority of**

the cohort. The concern with NACRT is that improved locoregional control is at the expense of undertreatment of microscopic distant disease, particularly where the radiotherapy field is limited anatomically (e.g. with GOJ tumours). The expected consequence of this would be fewer loco-regional recurrences and more distant recurrences, although this has not been demonstrated in other comparative studies and a recently published RCT(41).

This study lacks the number of patients to discretely analyse this relationship, however using individual variable importance calculation (available in the web app), the relative negative influence of NACRT (i.e. increased risk of recurrence compared to NACT) is on the whole more pronounced in GOJ tumours compared to Oesophageal tumours (an example of a 2<sup>nd</sup> order interaction), despite the recurrence rate being higher in Oesophageal compared to GOJ tumours.

Other risk factors for early recurrence including perioperative blood transfusion(42), complications of surgery(43) and preoperative staging were not available for this study, but are less discriminatory. Precise neoadjuvant regimens were not available for all patients in this study. It is therefore unclear if these results would be influenced by completion of treatment as prescribed, or indeed any adjuvant therapy given. This seems to have minimal effect on the model and suggests a small margin of effect on outcomes. Combining these factors could potentially increase the performance of our model if incorporated in the future. Ultimately, differential gene expression and mutation(44,45) may well determine prognostication and treatment pathways(46), but we are likely years from this being

universally available. Until then clinical and histopathological data remains the gold standard.

In that context, gains from mathematical and computer-based techniques are key to precision in delivery of cancer care. Here we have demonstrated several modern approaches that produce viable models. This study uses a dataset which is relatively small and simple in a ML context, and the improvement in performance over a standard LR is small (internal validation AUC 0.781). This is none-the-less important as this improvement is in effect 'free'. The strengths of this study lie in its multi-centre nature and heterogeneity of the cohort. This approach should maximise the utility of the model on external populations. All the data points used should be collected routinely at the majority of institutions, which should allow uptake without change in practice. The College of American Pathologists (CAP) definition of CRM positivity (i.e. CRM positive if tumour at the resection margin) was derivable for Centre G and performance was preserved in this subgroup if used instead of RCP definition (AUC of 0.813 with model generated on centres A-F (n=650) and validated on centre G (n=162), supplementary materials 5), supporting utility in both settings. We have also focussed on predictive model study design and reporting as suggested by the AJCC(47) and TRIPOD statements(48).

The training set was limited to patients undergoing neoadjuvant therapy for adenocarcinoma of the oesophagus. We have made no attempt to apply the model to a chemotherapy naïve population, and it is unlikely to calibrate well in this group due to the differing influence on survival of 'yp' compared to 'p' staging(49). It is also unclear if the model would be valid in patients with squamous cell carcinoma and we would advocate an

early external validation exercise using this patient group. A formal prospective validation/recalibration using the CAP definitions of CRM positivity would also be beneficial. Simulation studies have suggested that 100 – 200 cases (i.e. positives) are required for accurate validation(50), which assuming a stable incidence would require approximately 380 – 760 patients. A further limitation was the significant proportion of the original cases with missing data, which will have introduced a degree of selection bias. Multiple imputation is possible as a means of addressing this, however, was felt less appropriate in this study due to the high proportion of missing data being in the outcome measure and the lack of an external validation set.

## Conclusion

This large, multicentre cohort of patients who underwent oesophagectomy has been used to derive an accurate prediction model for early cancer recurrence, with excellent performance on validation. Machine learning techniques represent an attractive proposition for maximising performance of predictive models. The model is presented for use at <https://uoscancer.shinyapps.io/EROC/>.

## REFERENCES

1. Maynard M, Chadwick G, Varagunam M, Brand C, Cromwell D, Riley S, et al. National Oesophago-Gastric Cancer Audit 2017. *R Coll Surg Engl*. 2017;103.
2. Medical Research Council Oesophageal Cancer Working Group. Surgical resection with or without preoperative chemotherapy in oesophageal cancer: a randomised controlled trial. *Lancet*. 2002;359(9319):1727-1733.
3. Cunningham D, Allum WH, Stenning SP, Thompson JN, Van De Velde CJH, Nicolson M, et al. Perioperative Chemotherapy versus Surgery Alone for Resectable Gastroesophageal Cancer. *N Engl J Med*. 2006;355(1):11–20.
4. Shapiro J, van Lanschot JJB, Hulshof MCCM, van Hagen P, van Berge Henegouwen MI, Wijnhoven BPL, et al. Neoadjuvant chemoradiotherapy plus surgery versus surgery alone for oesophageal or junctional cancer (CROSS): long-term results of a randomised controlled trial. *Lancet Oncol*. 2015 Sep;16(9):1090–1098.
5. Davies AR, Pillai A, Sinha P, Sandhu H, Adeniran A, Mattsson F, et al. Factors associated with early recurrence and death after esophagectomy for cancer. *J Surg Oncol*. 2014;109(5):459–464.
6. Low DE, Kuppusamy MK, Alderson D, Cecconello I, Chang AC, Darling G, et al. Benchmarking Complications Associated with Esophagectomy. *Ann Surg*. 2019 Feb;269(2):291–298.
7. Shapiro J, Biermann K, Van Klaveren D, Offerhaus GJA, Ten Kate FJW, Meijer SL, et al. Prognostic value of pretreatment pathological tumor extent in patients treated with neoadjuvant chemoradiotherapy plus surgery for esophageal or junctional cancer. *Ann Surg*. 2017;265(2):356–362.
8. Goense L, van Rossum PSN, Xi M, Maru DM, Carter BW, Meijer GJ, et al. Preoperative

- Nomogram to Risk Stratify Patients for the Benefit of Trimodality Therapy in Esophageal Adenocarcinoma. *Ann Surg Oncol*. 2018;25(6):1598–1607.
9. Caruana R. An Empirical Comparison of Supervised Learning Algorithms. In: *Proceedings of the 23rd International Conference on Machine Learning*. Pittsburgh; 2006.
  10. Fernández-Delgado M, Cernadas E, Barro S, Amorim D, Amorim Fernández-Delgado D. Do we Need Hundreds of Classifiers to Solve Real World Classification Problems? *J Mach Learn Res*. 2014;15:3133–3181.
  11. Christodoulou E, Jie MA, Collins GS, Steyerberg EW, Verbakel JY, van Calster B. A systematic review shows no performance benefit of machine learning over logistic regression for clinical prediction models. *J Clin Epidemiol*. 2019;
  12. Noble F, Lloyd MA, Turkington R, Griffiths E, O'Donovan M, O'Neill JR, et al. Multicentre cohort study to define and validate pathological assessment of response to neoadjuvant therapy in oesophagogastric adenocarcinoma. *Br J Surg*. 2017 Dec;104(13):1816–1828.
  13. Stiles BM, Salzler GG, Nasar A, Paul S, Lee PC, Port JL, et al. Clinical predictors of early cancer-related mortality following neoadjuvant therapy and oesophagectomy. *Eur J Cardio-thoracic Surg*. 2015;48(3):455–460.
  14. RCPATH Cancer Services Working Group. Dataset for the histopathological reporting of oesophageal carcinoma (2nd edition). *R Coll Pathol*. 2013;(261035):1–27.
  15. Mandard AM, Dalibard F, Mandard JC, Marnay J, Henry-Amar M, Petiot JF, et al. Pathologic assessment of tumor regression after preoperative chemoradiotherapy of esophageal carcinoma. Clinicopathologic correlations. *Cancer*. 1994 Jun 1;73(11):2680–2686.

16. Zou H, Hastie T. Regularization and variable selection via the elastic net. *J R Stat Soc Ser B Stat Methodol.* 2005;67(2):301–320.
17. Breiman L. Random forests. *Mach Learn.* 2001;45(1):5–32.
18. Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining: KDD 2016.* San Francisco, CA; 2016. p. 785–794.
19. Pavlou M, Ambler G, Seaman SR, Guttman O, Elliott P, King M, et al. How to develop a more accurate risk prediction model when there are few events. *BMJ.* 2015;351:7–11.
20. Ranstam J, Cook JA. LASSO regression. *Br J Surg.* 2018;105(10):1348–1348.
21. Caruana R, Niculescu-Mizil A, Crew G, Ksikes A. Ensemble Selection from Libraries of Models. In: *Proceedings of the 21st International Conference on Machine Learning.* Banff, Canada; 2004.
22. Steyerberg EW, Harrell FE, Borsboom GJ, Eijkemans MJ, Vergouwe Y, Habbema JD. Internal validation of predictive models: efficiency of some procedures for logistic regression analysis. *J Clin Epidemiol.* 2001 Aug;54(8):774–781.
23. Steyerberg E, Moons KGM, van der Windt D, Hayden J, Perel P, Schroter S, et al. Prognosis research strategy (PROGRESS) series 3: prognostic model research. *PLoS Med.* 2013;10(2):e1001381.
24. Harrell FJ. *Regression Modeling Strategies: With Applications to Linear Models, Logistic and Ordinal Regression, and Survival Analysis.* Second Edi. Springer; 2015.
25. Steyerberg EW, Harrell FE, Jr. Prediction models need appropriate internal, internal-external, and external validation. *J Clin Epidemiol.* 2016;69:245.
26. Niculescu-Mizil A, Caruana R. Predicting Good Probabilities With Supervised Learning.

- In: Proceeding of the 22 international Conference on Machine Learning, Bonn, Germany. 2005. p. 625–632.
27. Chen W, Sahiner B, Samuelson F, Pezeshk A, Petrick N. Calibration of medical diagnostic classifier scores to the probability of disease. *Stat Methods Med Res.* 2018;27(5):1394–1409.
  28. Lundberg S, Lee S-I. A Unified Approach to Interpreting Model Predictions. In: *NIPS'17 Proceedings of the 31st International Conference on Neural Information Processing Systems*. California, USA; 2017. p. 4768–4777.
  29. Nanayakkara S, Fogarty S, Tremeer M, Ross K, Richards B, Bergmeir C, et al. Characterising risk of in-hospital mortality following cardiac arrest using machine learning: A retrospective international registry study. *PLoS Med.* 2018;15(11):1–16.
  30. Kuhn M. caret: Classification and Regression Training (Ver 6.0-81). 2018.
  31. Deane-Mayer Z, Knowles J. CaretEnsemble: Ensembles of Caret Models (ver 2.0.0). 2016.
  32. Molnar C, Bischl B, Casalicchio G. iml: An R Package for Interpretable Machine Learning. *J Open Source Softw.* 2018;3(26):786.
  33. Chang W, Cheng J, Xie Y, McPherson J. Shiny: Web Application Framework for R (ver 1.2.0).
  34. Reid TD, Chan DSY, Roberts SA, Crosby TDL, Williams GT, Lewis WG. Prognostic significance of circumferential resection margin involvement following oesophagectomy for cancer and the predictive role of endoluminal ultrasonography. *Br J Cancer.* 2012;107(12):1925–1931.
  35. Knight WRC, Zylstra J, Wulaningsih W, Van Hemelrijck M, Landau D, Maissey N, et al. Impact of incremental circumferential resection margin distance on overall survival

- and recurrence in oesophageal adenocarcinoma. *BJS Open*. 2018;2(4):229–237.
36. Smyth EC, Fassan M, Cunningham D, Allum WH, Okines AFC, Lampis A, et al. Effect of Pathologic Tumor Response and Nodal Status on Survival in the Medical Research Council Adjuvant Gastric Infusional Chemotherapy Trial. *J Clin Oncol*. 2016;34(23):2721–2727.
37. Davies AR, Myoteri D, Zylstra J, Baker CR, Wulaningsih W, Van Hemelrijck M, et al. Lymph node regression and survival following neoadjuvant chemotherapy in oesophageal adenocarcinoma. *Br J Surg*. 2018;105(12):1639–1649.
38. Noorani A, Goddard M, Crawte J, Alexandrov LB, Li X, Secrier M, et al. Genomic evidence supports a clonal diaspora model for metastases of esophageal adenocarcinoma. *bioRxiv*. 2018 Oct 30;454306.
39. Klevebro F, von Döbeln GA, Wang N, Johnsen G, Jacobsen AB, Friesland S, et al. A randomized clinical trial of neoadjuvant chemotherapy versus neoadjuvant chemoradiotherapy for cancer of the oesophagus or gastro-oesophageal junction. *Ann Oncol*. 2016;27(4):660–667.
40. Anderegg MCJ, van der Sluis PC, Ruurda JP, Gisbertz SS, Hulshof MCCM, van Vulpen M, et al. Preoperative Chemoradiotherapy Versus Perioperative Chemotherapy for Patients With Resectable Esophageal or Gastroesophageal Junction Adenocarcinoma. *Ann Surg Oncol*. 2017;24(8):2282–2290.
41. von Döbeln GA, Klevebro F, Jacobsen AB, Johannessen HO, Nielsen NH, Johnsen G, et al. Neoadjuvant chemotherapy versus neoadjuvant chemoradiotherapy for cancer of the esophagus or gastroesophageal junction: long-term results of a randomized clinical trial. *Dis Esophagus*. 2019;32(2):1–11.
42. Dresner SM, Lamb PJ, Shenfine J, Hayes N, Griffin SM. Prognostic significance of peri-

- operative blood transfusion following radical resection for oesophageal carcinoma. Eur J Surg Oncol. 2000;26:492–497.
43. Booka E, Takeuchi H, Suda K, Fukuda K, Nakamura R, Wada N, et al. Meta-analysis of the impact of postoperative complications on survival after oesophagectomy for cancer. BJS Open. 2018;2(5):276–284.
  44. Ueda M, Iguchi T, Masuda T, Nakahara Y, Hirata H, Uchi R, et al. Somatic mutations in plasma cell-free DNA are diagnostic markers for esophageal squamous cell carcinoma recurrence. Oncotarget. 2016;7(38):62280–62291.
  45. Lv H, He Z, Wang H, Du T, Pang Z. Differential expression of miR-21 and miR-75 in esophageal carcinoma patients and its clinical implication. Am J Transl Res. 2016;8(7):3288–3298.
  46. Walker RC, Underwood TJ. Molecular pathways in the development and treatment of oesophageal cancer. Best Pract Res Clin Gastroenterol. 2018;36–37:9–15.
  47. Moons KGM, Weiser MR, Lu Y, Halabi S, Gershengwald JE, Gimotty PA, et al. American Joint Committee on Cancer acceptance criteria for inclusion of risk models for individualized prognosis in the practice of precision medicine. CA Cancer J Clin. 2016;66(5):370–374.
  48. Collins GS, Reitsma JB, Altman DG, Moons KGM. Transparent Reporting of a Multivariable Prediction Model for Individual Prognosis or Diagnosis (TRIPOD): The TRIPOD Statement. Ann Intern Med. 2015 Jun;162(1):55–63.
  49. Rice TW, Patil DT, Blackstone EH. 8th edition AJCC/UICC staging of cancers of the esophagus and esophagogastric junction: application to clinical practice. Ann Cardiothorac Surg. 2017;6(2):119–130.
  50. Collins GS, Ogundimu EO, Altman DG. Sample size considerations for the external

validation of a multivariable prognostic model: A resampling study. Stat Med.  
2016;35(2):214–226.