

BMJ Open Machine-learning prediction of cancer survival: a retrospective study using electronic administrative records and a cancer registry

Sunil Gupta,¹ Truyen Tran,^{1,2} Wei Luo,¹ Dinh Phung,¹ Richard Lee Kennedy,³ Adam Broad,⁴ David Campbell,⁴ David Kipp,⁴ Madhu Singh,⁴ Mustafa Khasraw,^{3,4} Leigh Matheson,⁵ David M Ashley,^{3,4,5} Svetha Venkatesh¹

To cite: Gupta S, Tran T, Luo W, *et al*. Machine-learning prediction of cancer survival: a retrospective study using electronic administrative records and a cancer registry. *BMJ Open* 2014;**4**:e004007. doi:10.1136/bmjopen-2013-004007

► Additional material is available. To view please visit the journal (<http://dx.doi.org/10.1136/bmjopen-2013-004007>).

Received 13 September 2013
Revised 17 February 2014
Accepted 21 February 2014



CrossMark

For numbered affiliations see end of article.

Correspondence to

Professor Svetha Venkatesh; svetha.venkatesh@deakin.edu.au

ABSTRACT

Objectives: Using the prediction of cancer outcome as a model, we have tested the hypothesis that through analysing routinely collected digital data contained in an electronic administrative record (EAR), using machine-learning techniques, we could enhance conventional methods in predicting clinical outcomes.

Setting: A regional cancer centre in Australia.

Participants: Disease-specific data from a purpose-built cancer registry (Evaluation of Cancer Outcomes (ECO)) from 869 patients were used to predict survival at 6, 12 and 24 months. The model was validated with data from a further 94 patients, and results compared to the assessment of five specialist oncologists. Machine-learning prediction using ECO data was compared with that using EAR and a model combining ECO and EAR data.

Primary and secondary outcome measures: Survival prediction accuracy in terms of the area under the receiver operating characteristic curve (AUC).

Results: The ECO model yielded AUCs of 0.87 (95% CI 0.848 to 0.890) at 6 months, 0.796 (95% CI 0.774 to 0.823) at 12 months and 0.764 (95% CI 0.737 to 0.789) at 24 months. Each was slightly better than the performance of the clinician panel. The model performed consistently across a range of cancers, including rare cancers. Combining ECO and EAR data yielded better prediction than the ECO-based model (AUCs ranging from 0.757 to 0.997 for 6 months, AUCs from 0.689 to 0.988 for 12 months and AUCs from 0.713 to 0.973 for 24 months). The best prediction was for genitourinary, head and neck, lung, skin, and upper gastrointestinal tumours.

Conclusions: Machine learning applied to information from a disease-specific (cancer) database and the EAR can be used to predict clinical outcomes. Importantly, the approach described made use of digital data that is already routinely collected but underexploited by clinical health systems.

INTRODUCTION

Over the past two decades, there has been an explosion in the use of digital footprints

Strengths and limitations of this study

- This is the first study using machine learning of administrative and registry data for cancer survival prediction.
- A single prognosis model is produced across all cancers, improving prediction accuracy on rare cancers.
- This is a retrospective study in a single centre.

to monitor and predict human behaviours. The source of data used for this purpose is our online use of the internet, the emails we send and transactions we make. Analysis of these footprints through machine-learning techniques (MLT) has been exploited in the public domain by government and business to predict behaviours and inform investment decisions. In research, MLT have also been used to analyse gene expression data^{1 2} and for medical image analysis.^{3 4} However to date, there has been little exploration of these methodologies in the clinical setting. We hypothesised that MLT may offer a paradigm shift in clinical medicine that can address core issues with large and complex data sets. These techniques offer the potential to derive adaptive systems from diverse data sets, discover latent connections between data items and to predict outcomes.

Most hospitals routinely collect large digital electronic administrative records (EAR). These are primarily used for organisational financial management. Historically, they have not been used extensively for clinical or research purposes. If these large data sets are able to be exploited using MLT, it may open the way to optimise the use of collected administrative data to assist in predicting patients' outcome, planning individualised patient care, monitoring resource utilisation and improving

institutional performance.^{5 6} The accurate assessment of comorbid status would improve assessment of prognosis and guide treatment decisions.^{7–10} Other important information that may be contained or inferred from an EAR includes geographical and demographic data, socio-economic status and history of healthcare facility utilisation.^{2 11 12}

In this study, using cancer outcome prediction as a model, we wished to test the hypothesis that routinely collected digital health data, if analysed by state-of-the-art, validated, MLT could be used to assist conventional tools in predicting clinical outcomes.

Accurate prediction of survival in patients with cancer remains a challenge due to the ever-increasing heterogeneity and complexity of cancer, treatment options and patient populations. If achieved, reliable predictions could assist personalised care and treatment, and improve institutional performance in cancer management. In current practice, clinicians use data collected at the bedside in consultations, medical records or purpose-built cancer registries to aid prognostication and decision-making.

The notion of using MLT to predict cancer prognosis from clinical and pathological data is not a new one.^{13 14} However, with the advent of more sophisticated and better validated techniques, not only is more accurate prediction possible, but the range of data incorporated into decision aids can be increased.^{15–17} The need to improve cancer care systems by creating linkages between registries and epidemiological surveillance through analysis of complex and large clinical databases has recently been highlighted.^{18 19}

In this study, we tested the capability of MLT to predict patient outcomes in a heterogeneous cohort of patients with cancer. We have interrogated two data sets: first, a purpose-built cancer-specific registry (Evaluation of Cancer Outcomes, ECO, from Victorian Cancer Outcomes Network in partnership with the Barwon South Western Regional Integrated Cancer Service) containing demographic and tumour-related data items according to an Australian nationally agreed protocol; and second, a hospital digital data set containing information about the patient's previous admissions and presentations (EAR). Finally, in a test group of 94 patients, we examined the performance of machine-learning methods in aiding a panel of expert clinicians in predicting patient survival.

PATIENTS AND METHODS

Study design

This is a retrospective study using the EAR and a specialised cancer registry (ECO) from Barwon Health, the only public tertiary institution in a region of Australia with more than 350 000 residents. With a unified hospital identity number in use across the region, Barwon Health's EAR provides a single point of access for information on patient encounters with the health system,

including hospitalisations, ED visits, medications and treatments. In addition, the Andrew Love Cancer Centre at Barwon Health has a specialised cancer registry called ECO, which captures clinical data for patients in the region. ECO records information on demographics, primary tumour and metastatic tumour, cancer stage, tumour size, lymph nodes and breast tumour-specific information. Treatment type, outcomes, including death, and recurrence information (primary and metastatic) are also recorded. **Box 1** shows the variables used for survival prediction. The cohort for this study consists of 963 patients identified in ECO who were first diagnosed in year 2009. The study completion date was 31 October 2012; therefore, all patients had at least 2 year and 10 months follow-up. Among these patients, 736 patients also had records in the EAR.

Analyses

The analyses centred on predicting cancer survival since the date of diagnosis, defined as the date of tumour resection. Each patient was a unit of observation in the predictive problem: patient data collected prior to the diagnosis date were used to construct the independent variables; survival status in a period following the assessment was the dependent variable. Two analyses were performed: the first compared survival prediction made by machine-learning models and the clinician panel, based on only information from ECO. The second analysis evaluated the added discriminative power provided by EAR, by comparing the best machine-learning models using three sets of predicting variables: variables from

Box 1 Evaluation of Cancer Outcomes (ECO) variables used for survival prediction

Patient demographics
 Post code
 Gender
 Age
 Tumour characteristics
 Primary site (in International Classification of Diseases (ICD)-10 code)
 Tumour stream
 Morphology (in ICD-O-3 code)
 Histological grade
 Metastatic sites
 Most valid basis of diagnosis
 Performance status diagnosis
 Stage basis (pathological or clinical)
 Stage (TNM)
 Tumour size
 Nodes taken
 Positive nodes
 Breast cancer related variables
 Oestrogen receptor
 Progesterone receptor
 Human epidermal growth factor receptor 2 (HER2)

ECO (box 1), variables from EAR (see online supplementary appendix) and the union of the two.

Although a survival analysis model (eg, a proportional hazards model²⁰) is commonly used in modelling risk factors, such models are not designed to predict events. In this study, survival was directly modelled using classification models to optimise prediction accuracy.

Comparing predictions by machine-learning models and clinician

In the first analysis, all 963 patients in the ECO registry were randomly divided into a derivation cohort of 869 patients and a validation cohort of 94 patients (table 1). To collect clinician prediction, patients in the validation cohort were assigned to a panel of five oncologists for survival prediction. For each patient, the oncologist was asked to estimate the survival probabilities based on the independent variables in box 1. All clinicians estimated the patient's survival status by producing a probability for each of the three time periods—6 months, 1 year and 2 years. When making this assessment, the clinicians did not have knowledge of the treatment type offered or given to the patient. Three machine-learning models were trained on the derivation cohort using the same set of independent variables, one for each prediction period. Each of the machine-learning models was an ensemble of 400 support vector machines (SVMs)²¹ with linear kernel (ie, the output of the model was the average of 400 SVM outputs in Platt's a posteriori probabilities²²). Ensemble was used to control the variability introduced by L1 feature selection. Each of the SVMs was trained using a random 80% subsampling (without replacement) of the derivation cohort.²³ The soft margin parameter (C) of SVM was selected through cross-validation. Two measures were taken to improve the training process. First, to compensate for the imbalance between the two outcomes (there were more

survivals than deaths), we oversampled the non-surviving cases by 50% in each training subsample. Next, variable selection was performed through fitting a generalised linear model with elastic net regularisation²⁴ (α parameter set to 0.1 and λ parameter selected using fivefold internal cross-validation), and variables with zero coefficients were removed. After the machine-learning models were constructed, they were applied to predict survival probabilities for each patient in the validation cohort. The clinician and model predictions were validated with the actual outcomes in the ECO registry. Prediction performance was measured using the area under the receiver operating characteristic curve (AUC), also known as the C-statistic,²⁵ and 95% CIs of AUCs were computed using 1000 bootstrap samples of validation cohort.

Comparing discriminative information from specialised registry and routine data

The second analysis compared the discriminative power of two data sources (ECO and EAR). In this analysis, clinician predictions were not solicited. Among the 869 patients in the derivation subset of cohort 1, only 664 had records in the EAR and these patients were included in the second analysis (cohort 2, table 1). Survival prediction models were derived based on three sets of independent variables: (1) independent variables from EAR (EAR only); (2) independent variables from ECO (ECO only) and (3) the union of the two sets (EAR+ECO). Similar to the previous analysis, the models were trained using 400 random subsamples comprising 80% data of the cohort 2, and the modelling process was identical. However, the models were evaluated not using the validation cohort. Instead, for each 80% subsample, the remaining 20% was used to compute the AUC and its 95% CI.

Table 1 Characteristics of derivation and validation cohorts

	Cohort 1: ECO		Cohort 2: ECO and EAR (n=664)
	Derivation (n=869)	Validation (n=94)	
Age (SD)	67.6 (14.6)	68.4 (13.6)	66.3 (14.9)
Gender: male	487*	48	381
Tumour stream			
Genitourinary	172	21	135
Colorectal	140	14	115
Lung	121	18	96
Breast	122	15	74
Haematological	99	7	85
Upper gastrointestinal	83	9	57
Skin	36	1	28
Head and neck	35	0	30
Gynaecological	19	4	17
CNS	15	1	9
Unknown primary	38	9	26

*Two unspecified.

CNS, central nervous system; EAR, electronic administrative records; ECO, Evaluation of Cancer Outcomes.

Table 2 Performance of survival prediction: comparison between machine-learning method and clinicians

Survival period	AUC (95% CI)	
	Clinician panel	Machine-learning model
6 months	0.79 (0.76 to 0.81)	0.87 (0.85 to 0.89)
1 year	0.79 (0.76 to 0.81)	0.80 (0.77 to 0.82)
2 years	0.75 (0.73 to 0.78)	0.76 (0.74 to 0.79)

AUC, area under the receiver operating characteristic curve.

The Wilcoxon rank-sum test was applied to answer the following comparison problems:

1. Does *ECO only* provide more discriminative power than *EAR only*?
2. Does *EAR+ECO* provide more discriminative power than *EAR only*?
3. Does *EAR+ECO* provide more discriminative power than *ECO only*?

Details of the machine-learning model and the predictor variables can be found in the online supplementary appendix.

RESULTS

The cohorts for the two analyses are summarised in [table 1](#). The comparison between the algorithmic predictions and the clinician predictions are summarised in [table 2](#). The model had comparable performance to that of the clinicians, with the performance of the machine-learning model marginally better (AUC ranging from 0.76 to 0.87) than that of the clinicians (AUC ranging from 0.75 to 0.79) for all three prediction periods. This similarity in accuracy between algorithmic predictions and the clinician predictions was observed across different cancer types. Consider the predictions for 6-month survival. Of 15 breast cancer cases, the

clinicians made 15 correct predictions and the algorithm made 14; of 18 lung cancer cases, the clinicians made 13 correct predictions and the algorithm made 14; of 7 haematological cases, the clinicians and the algorithm made all predictions correctly. Similar results were observed on 12-month and 24-month survival predictions for different cancers.

Prediction of 6-month survival using the three models is shown in [table 3](#). There were no deaths from breast cancer during this period. Comparing the ECO model with the EAR model, AUCs were comparable for colorectal, genitourinary, haematological, head and neck, and skin tumours. The EAR model was significantly better ($p < 0.05$) for rare tumours, central nervous system (CNS), upper gastrointestinal and unassigned primary source tumours. For each tumour type, the model using ECO and EAR data yielded similar or better performance than the models using information from only one of the two databases. AUCs for the combined model ranged from 0.76 to 1.0. The combined data model showed particularly improved performance over ECO data ($p < 0.05$) for all tumour streams except breast and CNS tumours.

Data for 12-month survival prediction is shown in [table 4](#). Cancer-specific ECO data yielded better prediction than EAR data ($p < 0.05$) for gynaecological, haematological, lung, skin and unknown primary cancers. Otherwise, ECO and EAR models yielded generally similar results. The model using combined data performed better than EAR ($p < 0.05$) for all tumour streams other than CNS, head and neck and upper gastrointestinal tumours. The model using combined data was better than ($p < 0.05$) ECO for all cancers except breast, CNS, gynaecological and haematological cancers.

[Table 5](#) shows data for 24-month survival prediction by the three models. The ECO model yielded superior prediction ($p < 0.05$) to the EAR model for breast,

Table 3 Prediction performance of machine-learning algorithms: 6-month survival

Cancer type	Area under ROC curve (95% CI)		
	EAR only	ECO only	EAR+ECO
Genitourinary	0.81 (0.77 to 0.85)	0.82 (0.78 to 0.86)	0.88 (0.85 to 0.91)*, †
Colorectal	0.84 (0.80 to 0.88)	0.85 (0.81 to 0.89)	0.88 (0.84 to 0.91)*, †
Lung	0.71 (0.67 to 0.76)	0.73 (0.69 to 0.77)*	0.77 (0.73 to 0.82)*, †
Breast	no deaths in the period		
Haematological	0.73 (0.68 to 0.79)	0.74 (0.69 to 0.79)	0.76 (0.71 to 0.81)
Upper gastrointestinal	0.74 (0.69 to 0.78)	0.64 (0.60 to 0.69)	0.84 (0.80 to 0.87)†
Skin	0.84 (0.77 to 0.90)	0.85 (0.79 to 0.91)	0.91 (0.86 to 0.96)*, †
Head and neck	0.66 (0.61 to 0.71)	0.70 (0.64 to 0.75)	0.77 (0.72 to 0.82)*, †
Gynaecological	0.97 (0.94 to 0.99)	0.99 (0.98 to 1)*	1 (0.99 to 1)*
CNS	0.89 (0.85 to 0.94)	0.84 (0.78 to 0.90)	0.82 (0.77 to 0.88)
Unknown primary	0.92 (0.89 to 0.95)	0.79 (0.75 to 0.84)	0.90 (0.87 to 0.93)*, †

*Significantly greater than *EAR only*.

†Significantly greater than *ECO only*.

CNS, central nervous system; EAR, electronic administrative records; ECO, Evaluation of Cancer Outcomes; ROC, receiver operating characteristic.

Table 4 Prediction performance of machine-learning algorithms: 12-month survival

Cancer type	Area under ROC curve (95% CI)		
	EAR only	ECO only	EAR+ECO
Genitourinary	0.79 (0.75 to 0.83)	0.79 (0.75 to 0.83)	0.84 (0.80 to 0.87)*,†
Colorectal	0.82 (0.78 to 0.86)	0.83 (0.79 to 0.86)	0.87 (0.83 to 0.90)*,†
Lung	0.73 (0.69 to 0.77)	0.78 (0.73 to 0.82)*	0.82 (0.78 to 0.86)*,†
Breast	0.71 (0.65 to 0.78)	0.90 (0.86 to 0.94)	0.92 (0.89 to 0.96)*
Haematological	0.63 (0.59 to 0.68)	0.70 (0.66 to 0.75)*	0.69 (0.64 to 0.74)*
Upper gastrointestinal	0.62 (0.57 to 0.66)	0.70 (0.65 to 0.74)*	0.72 (0.68 to 0.76)*
Skin	0.76 (0.71 to 0.88)	0.89 (0.85 to 0.93)*	0.93 (0.90 to 0.96)*
Head and neck	0.77 (0.73 to 0.88)	0.68 (0.63 to 0.73)	0.79 (0.75 to 0.84)†
Gynaecological	0.95 (0.92 to 0.97)	1 (1 to 1)*	0.99 (0.98 to 1)*
CNS	0.66 (0.58 to 0.73)	0.68 (0.61 to 0.76)	0.69 (0.63 to 0.76)
Unknown primary	0.87 (0.84 to 0.91)	0.81 (0.77 to 0.85)	0.88 (0.84 to 0.91)

*Significantly greater than EAR only.

†Significantly greater than ECO only.

CNS, central nervous system; EAR, electronic administrative records; ECO, Evaluation of Cancer Outcomes; ROC, receiver operating characteristic.

genitourinary, gynaecological, lung, skin and unknown primary cancers, while the EAR model was superior to the ECO model for haematological and head and neck tumours. Once more, the model that performed the best was that derived from ECO and EAR data with AUCs ranging from 0.71 to 0.97 across the range of cancers and particularly enhanced performance for all cancers except breast, colorectal, gynaecological and unknown primary tumours compared with the ECO. In summary, over all time periods, the performance of the combined model was better than ECO ($p < 0.05$) for genitourinary, head and neck, lung, skin, and upper gastrointestinal tumours.

One of the key advantages of using MLT is that it can combine the large number of non-clinical factors with the few clinical risk factors. In this study, the model selected most of the known clinical risk factors including *patient age, cancer staging, performance status and tumour size*. In addition, it also found some useful non-clinical risk factors, including the type of the last hospital

admission (emergency vs elective), the frequency of ED visits within the previous 3 and 6 months (related to cancer and other medical conditions).

DISCUSSION

In this study, using cancer outcome prediction as a model, we wished to test the hypothesis that routinely collected digital health data, if analysed by MLT, could be used to assist conventional tools in predicting clinical outcomes.

Applying machine learning to data from the EAR alone predicted clinical outcomes with reasonable accuracy. Using the purpose-built ECO data set, the predictive tool also performed well across a broad range of cancer types, and in both cases the predictive accuracies were at least as good as that of a panel of five expert clinicians. Importantly, a predictive tool derived from the purpose-built clinical registry and administrative data had even greater predictive ability.

Table 5 Prediction performance of machine-learning algorithms: 24-month survival

Cancer type	Area under the ROC curve (AUC)		
	EAR only	ECO only	EAR+ECO
Genitourinary	0.73 (0.69 to 0.78)	0.84 (0.81 to 0.88)*	0.86 (0.82 to 0.89)*,†
Colorectal	0.76 (0.72 to 0.80)	0.76 (0.72 to 0.80)	0.76 (0.72 to 0.80)
Lung	0.74 (0.69 to 0.78)	0.78 (0.73 to 0.82)*	0.82 (0.79 to 0.86)*,†
Breast	0.67 (0.61 to 0.73)	0.86 (0.82 to 0.90)*	0.88 (0.84 to 0.92)*
Haematological	0.73 (0.68 to 0.77)	0.70 (0.66 to 0.75)	0.80 (0.76 to 0.84)*,†
Upper gastrointestinal	0.81 (0.77 to 0.85)	0.77 (0.72 to 0.81)	0.87 (0.83 to 0.9)*,†
Skin	0.71 (0.65 to 0.76)	0.85 (0.8 to 0.89)*	0.94 (0.92 to 0.97)*,†
Head and neck	0.74 (0.7 to 0.78)	0.66 (0.51 to 0.61)	0.71 (0.67 to 0.76)†
Gynaecological	0.96 (0.94 to 0.99)	0.99 (0.98 to 1)*	0.97 (0.95 to 0.99)
CNS	0.83 (0.78 to 0.89)	0.87 (0.82 to 0.93)	0.96 (0.93 to 0.99)*,†
Unknown primary	0.74 (0.7 to 0.79)	0.78 (0.74 to 0.82)*	0.8 (0.76 to 0.84)*

*Significantly greater than EAR only.

†Significantly greater than ECO only.

CNS, central nervous system; EAR, electronic administrative records; ECO, Evaluation of Cancer Outcomes; ROC, receiver operating characteristic.



The wealth of administrative data contained in the EAR includes information on comorbid conditions and previous clinic and hospital attendances as well as a drug history. There is considerable potential to use this data to improve clinical care across a spectrum of diseases.^{5 6}

Most patients in the study were followed up for 3 years, which may not be adequate to capture all oncological outcomes, especially for those cancers with low mortality rate. We have designed this study as retrospective and in a single centre; it will be of major interest to observe how it performs in a variety of settings. The number of cases used to assess performance of the models is relatively small. The strengths include the comparison of machine-learning tools with expert clinical opinion and the fact that very detailed and well-validated data was available both directly related to the cancer and that contained in the EAR. The generic nature of this approach makes it unnecessary to generate separate predictive models for different types of cancer. This was a particular advantage for rarer forms of cancer where predications using more conventional methods are very challenging.

Predictive tools derived from clinical data items have considerable potential to improve clinical care, but must be suitably optimised and shown to perform equally well in diverse clinical settings.^{26 27} Clinical databases have become more widely available and increasingly complex in recent years. The extent and complexity of data available to clinicians means that novel approaches to managing data and supporting clinical decisions are needed. Machine-learning approaches can not only cope with complex data sets, but also adapt in real time and across different clinical settings.

The approach used in this study offers superior performance to previous machine-learning approaches in predicting cancer survival.^{13–17} Previous models have been derived for single cancer types, or for a limited range of cancers. The model described here performed well across a wide range of cancers. One advantage of this generic approach may be the ability to predict outcomes in less common cancers where limited data might preclude development of specific models. The fact that our model derived from administrative and cancer-related data performed slightly better than a panel of expert clinicians validates the potential utility of the model and suggests that it may be useful in assessing quality of care and also in settings where specialist care is not available. An alternative approach to borrow information across different cancer types is called multitask learning. We are currently exploring this approach as well.

Clinical outcomes in any illness are determined by specific factors related to the illness itself and also by the patient's general state of health and by the presence of other chronic medical conditions often coded in an EAR if the individual traffics the health service.^{7–10} As well, a particularly novel and important aspect of the use

of historical data from the EAR in machine learning is that it effectively captures the healthcare institution's current and previous performance. These data can be applied to any individual entering the system with a newly diagnosed cancer, as we have modelled here. As well, they could also be used for quality and performance monitoring.

In conclusion, machine learning applied to information from a disease-specific (cancer) database and the EAR can be used to predict outcomes. Improved prediction of outcome has the potential to help clinicians make more meaningful decisions about treatment and to assist with planning of future social and care needs. Most importantly, the approach described makes use of digital data that is already routinely collected but under-exploited by clinical health systems.

Author affiliations

¹Centre for Pattern Recognition and Data Analytics, Deakin University, Geelong, Victoria, Australia

²Department of Computing, Curtin University, Perth, Western Australia, Australia

³School of Medicine, Deakin University, Geelong, Victoria, Australia

⁴Andrew Love Cancer Centre, Barwon Health, Geelong, Victoria, Australia

⁵Barwon Southwest Integrated Cancer Service, Geelong, Victoria, Australia

Contributors All authors of this research paper have directly participated in the planning, execution or analysis of the study. All authors of this paper have read and approved the final version submitted.

Funding This research received no specific grant from any funding agency in the public, commercial or not-for-profit sectors.

Competing interests None.

Ethics approval Ethics approval was obtained from the Hospital and Research Ethics Committee at Barwon Health (number 12/83). Deakin University has reciprocal ethics authorisation with Barwon Health.

Provenance and peer review Not commissioned; externally peer reviewed.

Data sharing statement No additional data are available.

Open Access This is an Open Access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 3.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/3.0/>

REFERENCES

1. Zhao X, Rodland EA, Sorlie T, *et al.* Combining gene signatures improves prediction of breast cancer survival. *PLoS ONE* 2011;6: e17845.
2. Chang CM, Su YC, Lai NS, *et al.* The combined effect of individual and neighborhood socioeconomic status on cancer survival rates. *PLoS ONE* 2012;7:e44325.
3. Li C, Zhang S, Zhang H, *et al.* Using the k-nearest neighbor algorithm for the classification of lymph node metastasis in gastric cancer. *Comput Math Methods Med* 2012;2012:876545.
4. Huang ML, Hung YH, Lee WM, *et al.* Usage of case-based reasoning, neural network and adaptive neuro-fuzzy inference system classification techniques in breast cancer dataset classification diagnosis. *J Med Syst* 2012;36:407–14.
5. Appari A, Eric Johnson M, Anthony DL. Meaningful use of electronic health record systems and process quality of care: evidence from a panel data analysis of U.S. acute-care hospitals. *Health Serv Res* 2013;48:354–75.
6. Fitzhenry F, Murff HJ, Matheny ME, *et al.* Exploring the frontier of electronic health record surveillance: the case of postoperative complications. *Med Care* 2013;51:509–16.

7. Lund L, Borre M, Jacobsen J, *et al.* Impact of comorbidity on survival of Danish prostate cancer patients, 1995–2006: a population-based cohort study. *Urology* 2008;72:1258–62.
8. Tetsche MS, Norgaard M, Jacobsen J, *et al.* Comorbidity and ovarian cancer survival in Denmark, 1995–2005: a population-based cohort study. *Int J Gynecol Cancer* 2008;18:421–7.
9. Lieffers JR, Baracos VE, Winget M, *et al.* A comparison of Charlson and Elixhauser comorbidity measures to predict colorectal cancer survival using administrative health data. *Cancer* 2011;117:1957–65.
10. Braithwaite D, Moore DH, Satariano WA, *et al.* Prognostic impact of comorbidity among long-term breast cancer survivors: results from the lace study. *Cancer Epidemiol Biomarkers Prev* 2012;21:1115–25.
11. Jones LE, Doebbeling CC. Beyond the traditional prognostic indicators: the impact of primary care utilization on cancer survival. *J Clin Oncol* 2007;25:5793–9.
12. Sant M, Minicozzi P, Allemani C, *et al.* Regional inequalities in cancer care persist in Italy and can influence survival. *Cancer Epidemiol* 2012;36:541–7.
13. Burke HB, Goodman PH, Rosen DB, *et al.* Artificial neural networks improve the accuracy of cancer survival prediction. *Cancer* 1997;79:857–62.
14. Lundin M, Lundin J, Burke HB, *et al.* Artificial neural networks applied to survival prediction in breast cancer. *Oncology* 1999;57:281–6.
15. Manilich EA, Kiran RP, Radivoyevitch T, *et al.* A novel data-driven prognostic model for staging of colorectal cancer. *J Am Coll Surg* 2011;213:579–88, 588.e1–2.
16. Gao P, Zhou X, Wang ZN, *et al.* Which is a more accurate predictor in colorectal survival analysis? Nine data mining algorithms vs. the TNM staging system. *PLoS ONE* 2012;7:e42015.
17. Kim W, Kim KS, Lee JE, *et al.* Development of novel breast cancer recurrence prediction model using support vector machine. *J Breast Cancer* 2012;15:230–8.
18. Johnson CJ, Weir HK, Fink AK, *et al.* Accuracy of Cancer Mortality Study Group. The impact of National Death Index linkages on population-based cancer survival rates in the United States. *Cancer Epidemiol* 2013;37:20–8.
19. Khoury MJ, Lam TK, Ioannidis JP, *et al.* Transforming epidemiology for 21st century medicine and public health. *Cancer Epidemiol Biomarkers Prev* 2013;22:508–16.
20. Cox DR, Oakes D. *Analysis of survival data*. CRC Press, 1984.
21. Cortes C, Vapnik V. Support vector machine. *Mach Learn* 1995;20:273–97.
22. Lin H-T, Lin C-J, Weng RC. A note on Platt's probabilistic outputs for support vector machines. *Mach Learn* 2007;68:267–76.
23. Politis D, Romano J, Wolf M. *Subsampling*. New York: Springer-Verlag, 1999.
24. Friedman J, Hastie T, Tibshirani R. Regularization paths for generalized linear models via coordinate descent. *J Stat Softw* 2010;33:1–22.
25. Hastie T, Tibshirani R, Friedman J, *et al.* The elements of statistical learning: data mining, inference and prediction. *Math Intelligencer* 2005;27:83–5.
26. Chen HC, Kodell RL, Cheng KF, *et al.* Assessment of performance of survival prediction models for cancer prognosis. *BMC Med Res Methodol* 2012;12:102.
27. Chen HC, Chen JJ. Assessment of reproducibility of cancer survival risk predictions across medical centers. *BMC Med Res Methodol* 2013;13:25.

Appendix

In this section we describe the procedure used to build our machine learning model.

Derivation of the machine learning model

We used an ensemble of classifiers to achieve a low variance model. From the derivation cohort, data is randomly split to extract 80% for training (derivation train set) and 20% for testing (derivation test set). This is done by subsampling without replacement. This procedure is repeated 400 times to generate 400 random subsamples (or training/test pairs). The training sets were used to estimate an ensemble of classifiers while the test sets were used to assess the performance of these classifiers (mean Area under ROC curve and 95% CI).

For each training set subsample, a classification model was estimated using the derivation train set. Estimation of the classifier contains two phases: feature selection and classifier design. In *feature selection*, we used an established statistical technique - a generalized linear model with l_1 -norm and l_2 -norm penalty (alpha parameter set to 0.1 and lambda parameter selected using 5-fold internal cross-validation) [1]. Features with nonzero coefficients were selected. Next, using this feature set, the parameters of a *linear Support Vector machine* [2] classifier were estimated. For SVM implementation, we used the open source package LIBSVM [3].

The above procedure generates an ensemble of 400 classifiers to be tested against on the held-out validation cohort. Three such classifier-ensembles were built, one for each survival prediction tasks (i.e. prediction at 6, 12 and 24 months periods).

Predictors for the machine learning models

Table 1 EMR-based predictors

demographics

- gender
- age
- spoken language
- country of origin
- religion
- occupation
- marital status
- insurance type

cancer specific diagnoses

- primary site
- tumor stream (e.g., breast)
- tumor
- morphology code
- topology code

patient history (in the previous 1 month, 3 months, and 6 months)

- number of inpatient admissions
- number of ED visits
- number of admissions from ED
- longest length of hospital stay
- average length of hospital stay
- number of operations
- number of oncology visits
- number of histology tests
- discharge diagnoses in ICD-10
- diagnosis-related groups codes
- procedure codes

Table 2 ECO-based predictors.

patient demographics

Gender

Age

tumour characteristics

primary site (in ICD-10 code)

tumour stream

morphology (in ICD-O-3 code)

histologic grade

metastatic sites

most valid basis of diagnosis

performance status diagnosis

stage basis (pathological or clinical)

stage (TNM)

tumour size

nodes taken

positive nodes

breast cancer related variables

oestrogen receptor

progesterone receptor

human epidermal growth factor receptor 2 (HER2)

References

1. Tibshirani R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* 1996;**58**(1):267-88
2. Cortes C, Vapnik V. Support vector machine. *Machine learning* 1995;**20**(3):273-97
3. Chang C-C, Lin C-J. LIBSVM: a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology (TIST)* 2011;**2**(3):27