

Machine Learning for Survival Analysis: A New Approach

David Dooling^{1,✉}, Angela Kim^{1,‡}, Jennifer Webster^{1,✉}

1 Innovative Oncology Business Solutions, Albuquerque, NM, USA

✉These authors contributed equally to this work.

‡These authors also contributed equally to this work.

* ddooling@innovativeobs.com

Abstract

We have applied a little-known data transformation to subsets of the Surveillance, Epidemiology, and End Results (SEER) publically available data of the National Cancer Institute (NCI) to make it suitable input to standard machine learning classifiers. This transformation properly treats the right-censored data in the SEER data and the resulting Random Forest and Multi-Layer Perceptron models predict full survival curves. Treating the 6, 12, and 60 months points of the resulting survival curves as 3 binary classifiers, the 18 resulting classifiers have AUC values ranging from .765 to .885. Further evidence that the models have generalized well from the training data is provided by the extremely high levels of agreement between the random forest and neural network models predictions on the 6, 12, and 60 month binary classifiers.

Author Summary

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Curabitur eget porta erat. Morbi consectetur est vel gravida pretium. Suspendisse ut dui eu ante cursus gravida non sed sem. Nullam sapien tellus, commodo id velit id, eleifend volutpat quam. Phasellus mauris velit, dapibus finibus elementum vel, pulvinar non tellus. Nunc pellentesque pretium diam, quis maximus dolor faucibus id. Nunc convallis sodales ante, ut ullamcorper est egestas vitae. Nam sit amet enim ultrices, ultrices elit pulvinar, volutpat risus.

Introduction

Opportunities are emerging in many industries today to develop and deploy services that cater to individual needs and preferences. Music afficianados can create their own radio stations tailored to their individual tastes from Pandora¹, bibliophiles can receive highly trustworthy book recommendations from goodreads.com², and Google will provide directions between any two points, giving options such as mode of transportation and as well as warnings of delays in realtime.³ These individualized services share many

¹Pandora Internet Radio - Listen to Free Music You'll Love, <http://www.pandora.com/> (accessed 27 Jan 2016)

²Share Book Recommendations With Your Friends, Join Book Clubs, Answer Trivia, <https://www.goodreads.com/> (accessed 27 Jan 2016)

³Google Maps, <https://goo.gl/1D7Jwf> (accessed 27 Jan 2016)

common features. In particular, they leverage large databases of aggregated information to learn and extract information relevant to individuals. Extracting actionable information from data is changing the fabric of modern business. A class of techniques that transforms data into actionable information goes by the name of Machine Learning [1]. Machine Learning has recently become a popular method to answer questions and solve problems that are too complex to solve via traditional methods.

The primary objective of this study is to show how machine learning methods can be trained with data in cancer registries to produce personalized survival prognosis curves, but the methods presented below can be applied to any type of survival data. Traditionally, cancer survival curves have been estimated using Kaplan-Meier methods [2]. Kaplan-Meier methodology also uses large datasets to make predictions, but the resulting information is not personal; the resulting curves are summaries for a population and not necessarily relevant or particularly accurate for any given individual. This property of Kaplan-Meier methods is exacerbated when dealing with heterogeneous populations. The methods described below also take full advantage of all relevant aggregate information, but are able to provide personalized survival curves relevant to individual subjects. This objective is in keeping with the recent movement in medicine known as Predictive, Preventive and Personalized Medicine (PPPM), which aims to leverage increasing amounts of health related data to maximize quality of care and to intelligently eliminate inefficient and unnecessary use of resources [3]. This capability of providing individualized survival curve prognosis is a direct result of the recent advances in computing power and machine learning algorithms, and similar methodology is becoming commonplace in many industries. These techniques are now infiltrating the healthcare industry, in spite of some of the data aggregation challenges posed by the Health Insurance Portability and Accountability Act (HIPPA) of 1996. This study makes use of a freely available data source that circumvents the restrictions imposed by HIPPA.

The Surveillance, Epidemiology, and End Results (SEER) Program of the National Cancer Institute (NCI) has been collecting data because intuitively researchers feel confident that this data will eventually allow researchers to detect information crucial to patients and providers including the relationships between the types of data collected (demographic as well as staging information, treatment and disease characteristics) and the survival outcomes. Though these relationships evade capture by traditional methods, it is possible to surface them with two machine learning techniques known as *Random Forests* and *Neural Networks*. As will be demonstrated in section , these two methods produce very similar results when applied to the SEER dataset, and are based on almost diametrically opposed learning philosophies, which lends confidence in the validity of the results.

The Surveillance, Epidemiology, and End Results (SEER) Program of the National Cancer Institute (NCI) is the most recognized authoritative source of information on cancer incidence and survival in the United States. SEER currently collects and publishes cancer incidence and survival data from population-based cancer registries covering approximately 28 percent of the US population.

Quoting directly from the SEER website [4]:

The SEER program registries routinely collect data on patient demographics, primary tumor site, tumor morphology and stage at diagnosis, first course of treatment, and follow-up for vital status. This program is the only comprehensive source of population-based information in the United States that includes stage of cancer at the time of diagnosis and patient survival data. The mortality data reported by SEER are provided by the National Center for Health Statistics. The population data used in calculating cancer rates is obtained periodically from the Census Bureau.

Updated annually and provided as a public service in print and electronic formats, SEER data are used by thousands of researchers, clinicians, public health officials, legislators, policymakers, community groups, and the public.

One characteristic of the SEER data that is shared by many datasets in the medical field goes by the name of "censored data." Observations are labeled censored when the survival time information is incomplete. The SEER data contains the number of months each patient survived, as well as an indicator variable showing whether or not the patient is still alive at the end of the data collection period. Methods to deal effectively with this kind of "right-censored data" include Kaplan-Meier curves and Cox Proportional Hazard models [2]. The Kaplan-Meier techniques only give estimates for cohorts of patients and are not applicable for predicting the survival curve for a single patient, and the Cox Proportional Hazard models require a fairly restrictive set of assumptions to be satisfied in order to yield reliable results.

Previous work applying machine learning methods to subsets of the SEER data include creative attempts to deal with the problems presented by "right-censored data." Shin et al. [5] use semi-supervised learning techniques to predict 5 year survival, essentially imputing values for SEER records where the survival months information is censored at a value less than 5 years. Zolbanin et al. [6] investigate the effects of comorbidities; i.e., patients with two different cancer diagnoses, but their treatment of the censored data underestimates the survival probabilities. All records representing patients who survived at least 60 months as well as all those who died earlier than 60 months were considered, but patients alive prior to 60 months but censored out of the study before 60 months were not included. This treatment biases the data and the predictions, leading to overly pessimistic survival probabilities predicted by the models.

Previous work applying machine learning methods based on decision trees to survival data in general have a long history, starting with Gordon et al. [7]. A summary of more recent developments concerning *survival trees* is provided by Bou-Hamad et al. [8]. These methods focus on altering the splitting criteria used in decision tree growth to account for the censoring, and use 1958 Kaplan-Meier methods at the resulting nodes for prediction purposes. These methods do not generalize to non-tree-based machine learning algorithms, though Ishwaran et al. have extended the methodology to *random survival forests*, ensembles of *survival trees* [9].

IOBS has applied a little-known technique to transform the SEER data to make it amenable to more powerful machine learning methods. Instead of modifying existing learning algorithms in drastic ways, we focus attention on the input data. This approach allows for different machine learning algorithms to use the same data with no modification. The essential idea is to recast the problem to an appropriate discrete classification problem instead of a regression problem (predicting survival months). Treating months after diagnosis as just another discrete feature, the SEER data (or any other right-censored data) can be transformed to make predictions for the hazard function (probability of dying in the next month, given that the patient has not yet died). The full survival function can then be derived from the hazard function.

This paper is organized as follows. We introduce the subsets of the SEER data used for this study, and present survival curves computed from traditional methods based on this data for the three cancer types *lung*, *breast*, and *colon*. We then present the essential methodology of this work, the data transformation that allows censored survival data to be used as input to existing machine learning classifiers. Then we present the details of the trained models, including some subtleties arising from the data transformation pertaining to the partition into training and test datasets. The method of deriving binary classifiers from the models' predictions for the survival curves is presented. In this paper, we have constructed binary classifiers corresponding to 6, 12, and 60 months, as these are standard metrics in cancer survival prognosis. Then follows

a discussion of the evaluation of the trained models. The performance metrics are the 18 AUC curves associated with the 6, 12, and 60 month survival binary classifiers for the two models associated with each cancer type. We also present additional evidence supporting validity of the predictions by computing the levels of agreement between the random forest and neural network models for each of the 18 binary classifiers and find striking agreement. Next we provide urls for 6 web applications that use the trained models to predict individual cancer survival prognosis curves. These apps are hosted on the popular Heroku website, and allow for exploration of the nonlinear relationships between the input features and resulting survival prognosis. It is exactly these kinds of tools that are the goal of Predictive, Preventive and Personalized Medicine. Finally, we present avenues for future research.

Materials and Methods

For this study we use the publically available 1973-2012 SEER incidence data files corresponding to colon, breast and lung cancer contained in the list below. SEER requires that researchers submit a request for the data, which includes an agreement form. Detailed documentation explaining the contents of both the incidence data files used in this study as well as a data dictionary for the 1973-2012 SEER incidence data files are available without the need to register or submit a data request [10].

- incidence\yr1973.2012.seer9\COLRECT.txt
- incidence\yr1973.2012.seer9\BREAST.txt
- incidence\yr1973.2012.seer9\RESPIR.txt
- incidence\yr1992.2012.sj_la_rg_ak\COLRECT.txt
- incidence\yr1992.2012.sj_la_rg_ak\BREAST.txt
- incidence\yr1992.2012.sj_la_rg_ak\RESPIR.txt
- incidence\yr2000.2012.ca_ky_lo_nj_ga\COLRECT.txt
- incidence\yr2000.2012.ca_ky_lo_nj_ga\BREAST.txt
- incidence\yr2000.2012.ca_ky_lo_nj_ga\RESPIR.txt
- incidence\yr2005.lo_2nd_half\COLRECT.txt
- incidence\yr2005.lo_2nd_half\BREAST.txt
- incidence\yr2005.lo_2nd_half\RESPIR.txt

Data preparation and preprocessing

A great deal of data munging is necessary before using these SEER incidence files as input into machine learning algorithms. A preprocessing step common to each of the three cancer types studied involves the SEER STATE-COUNTY RECODE variable. The STATE-COUNTY RECODE field is a state-county combination where the first two characters represent the state FIPS code and the last three digits represent the FIPS county code. The FIPS code is a five-digit Federal Information Processing Standard (FIPS) code which uniquely identifies counties and county equivalents in the United States, certain U.S. possessions, and certain freely associated states. This particular field illustrates an important characteristic of machine learning, that is, the difference between *categorical features* and *numeric features*. All input into a machine learning algorithm must be numeric, but real numbers carry with them the usually extremely useful property known as the well-ordering property. Machine learning algorithms use the well-ordering property of the real numbers to learn. But if one is tasked with encoding a categorical feature into suitable numeric format for machine learning, it is necessary to do so in a way that removes the well-ordering property [11].

Etiam eget sapien nibh.

Nulla mi mi, Fig. 1 venenatis sed ipsum varius, volutpat euismod diam. Proin rutrum vel massa non gravida. Quisque tempor sem et dignissim rutrum. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi at justo vitae nulla elementum commodo eu id massa. In vitae diam ac augue semper tincidunt eu ut eros. Fusce fringilla erat porttitor lectus cursus, S1 Video vel sagittis arcu lobortis. Aliquam in enim semper, aliquam massa id, cursus neque. Praesent faucibus semper libero.

Figure 1. Figure Title first bold sentence Nulla mi mi, venenatis sed ipsum varius, volutpat euismod diam. Figure Caption Proin rutrum vel massa non gravida. Quisque tempor sem et dignissim rutrum. A: Lorem ipsum dolor sit amet. B: Consectetur adipiscing elit.

1. react
2. diffuse free particles
3. increment time by dt and go to 1

Results

Nulla mi mi, venenatis sed ipsum varius, Table 1 volutpat euismod diam. Proin rutrum vel massa non gravida. Quisque tempor sem et dignissim rutrum. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi at justo vitae nulla elementum commodo eu id massa. In vitae diam ac augue semper tincidunt eu ut eros. Fusce fringilla erat porttitor lectus cursus, vel sagittis arcu lobortis. Aliquam in enim semper, aliquam massa id, cursus neque. Praesent faucibus semper libero.

Table 1. Table caption Nulla mi mi, venenatis sed ipsum varius, volutpat euismod diam.

Heading1				Heading2			
cell1row1	cell2 row 1	cell3 row 1	cell4 row 1	cell5 row 1	cell6 row 1	cell7 row 1	cell8 row 1
cell1row2	cell2 row 2	cell3 row 2	cell4 row 2	cell5 row 2	cell6 row 2	cell7 row 2	cell8 row 2
cell1row3	cell2 row 3	cell3 row 3	cell4 row 3	cell5 row 3	cell6 row 3	cell7 row 3	cell8 row 3

Table notes Phasellus venenatis, tortor nec vestibulum mattis, massa tortor interdum felis, nec pellentesque metus tortor nec nisl. Ut ornare mauris tellus, vel dapibus arcu suscipit sed.

LOREM and IPSUM Nunc blandit a tortor.

Maecenas convallis mauris sit amet sem ultrices gravida. Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. Curabitur fringilla pulvinar lectus consectetur pellentesque. Quisque augue sem, tincidunt sit amet feugiat eget, ullamcorper sed velit. Sed non aliquet felis. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Mauris commodo justo ac dui pretium imperdiet. Sed suscipit iaculis mi at feugiat.

Sed ac quam id nisi malesuada congue.

Nulla mi mi, venenatis sed ipsum varius, volutpat euismod diam. Proin rutrum vel massa non gravida. Quisque tempor sem et dignissim rutrum. Lorem ipsum dolor sit

amet, consectetur adipiscing elit. Morbi at justo vitae nulla elementum commodo eu id massa. In vitae diam ac augue semper tincidunt eu ut eros. Fusce fringilla erat porttitor lectus cursus, vel sagittis arcu lobortis. Aliquam in enim semper, aliquam massa id, cursus neque. Praesent faucibus semper libero.

Subsection 1

Nulla mi mi, venenatis sed ipsum varius, volutpat euismod diam. Proin rutrum vel massa non gravida. Quisque tempor sem et dignissim rutrum. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi at justo vitae nulla elementum commodo eu id massa. In vitae diam ac augue semper tincidunt eu ut eros. Fusce fringilla erat porttitor lectus cursus, vel sagittis arcu lobortis. Aliquam in enim semper, aliquam massa id, cursus neque. Praesent faucibus semper libero.

Subsection 2

3rd Level Heading. Nulla mi mi, venenatis sed ipsum varius, volutpat euismod diam. Proin rutrum vel massa non gravida. Quisque tempor sem et dignissim rutrum. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi at justo vitae nulla elementum commodo eu id massa. In vitae diam ac augue semper tincidunt eu ut eros. Fusce fringilla erat porttitor lectus cursus, vel sagittis arcu lobortis. Aliquam in enim semper, aliquam massa id, cursus neque. Praesent faucibus semper libero.

Discussion

Nulla mi mi, venenatis sed ipsum varius, Table 1 volutpat euismod diam. Proin rutrum vel massa non gravida. Quisque tempor sem et dignissim rutrum. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi at justo vitae nulla elementum commodo eu id massa. In vitae diam ac augue semper tincidunt eu ut eros. Fusce fringilla erat porttitor lectus cursus, vel sagittis arcu lobortis. Aliquam in enim semper, aliquam massa id, cursus neque. Praesent faucibus semper libero.

LOREM and IPSUM Nunc blandit a tortor.

CO₂ Maecenas convallis mauris sit amet sem ultrices gravida. Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. Curabitur fringilla pulvinar lectus consectetur pellentesque. Quisque augue sem, tincidunt sit amet feugiat eget, ullamcorper sed velit.

Sed non aliquet felis. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Mauris commodo justo ac dui pretium imperdiet. Sed suscipit iaculis mi at feugiat. Ut neque ipsum, luctus id lacus ut, laoreet scelerisque urna. Phasellus venenatis, tortor nec vestibulum mattis, massa tortor interdum felis, nec pellentesque metus tortor nec nisl. Ut ornare mauris tellus, vel dapibus arcu suscipit sed. Nam condimentum sem eget mollis euismod. Nullam dui urna, gravida venenatis dui et, tincidunt sodales ex. Nunc est dui, sodales sed mauris nec, auctor sagittis leo. Aliquam tincidunt, ex in facilisis elementum, libero lectus luctus est, non vulputate nisl augue at dolor. For more information, see S1 Text.

Supporting Information 224

S1 Video 225

Bold the first sentence. Maecenas convallis mauris sit amet sem ultrices gravida. 226
Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. 227
Curabitur fringilla pulvinar lectus consectetur pellentesque. 228

S1 Text 229

Lorem Ipsum. Maecenas convallis mauris sit amet sem ultrices gravida. Etiam eget 230
sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. Curabitur 231
fringilla pulvinar lectus consectetur pellentesque. 232

S1 Fig 233

Lorem Ipsum. Maecenas convallis mauris sit amet sem ultrices gravida. Etiam eget 234
sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. Curabitur 235
fringilla pulvinar lectus consectetur pellentesque. 236

S2 Fig 237

Lorem Ipsum. Maecenas convallis mauris sit amet sem ultrices gravida. Etiam eget 238
sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. Curabitur 239
fringilla pulvinar lectus consectetur pellentesque. 240

S1 Table 241

Lorem Ipsum. Maecenas convallis mauris sit amet sem ultrices gravida. Etiam eget 242
sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. Curabitur 243
fringilla pulvinar lectus consectetur pellentesque. 244

Acknowledgments 245

Cras egestas velit mauris, eu mollis turpis pellentesque sit amet. Interdum et malesuada 246
fames ac ante ipsum primis in faucibus. Nam id pretium nisi. Sed ac quam id nisi 247
malesuada congue. Sed interdum aliquet augue, at pellentesque quam rhoncus vitae. 248

References

1. Sebastian Raschka. Python Machine Learning Essentials. Packt Publishing; 2015.
2. Cam Davidson-Pilon. Quickstart – lifelines 0.8.0.1 documentation; 2016 (accessed 14 Jan 2016).
<http://lifelines.readthedocs.org/en/latest/Quickstart.html>.
3. Van Poucke S, Zhang Z, Schmitz M, Vukicevic M, Laenen MV, Celi LA, et al. Scalable predictive analysis in critically ill patients using a visual open data analysis platform. PLoS ONE. 2016;11(1). Cited By 0. Available from:
<http://www.scopus.com/inward/record.url?eid=2-s2.0-84953931466&partnerID=40&md5=7a0cad7137c03146e4b75f3295f84cc6>.

4. National Cancer Institute, the Surveillance, Epidemiology, and End Results Program. About the SEER Program - SEER; 2016 (accessed 14 Jan 2016). <http://seer.cancer.gov/about>.
5. Shin, Hyunjung and Nam, Yonghyun; ISCB Asia. A coupling approach of a predictor and a descriptor for breast cancer prognosis [Article; Proceedings Paper]. BMC MEDICAL GENOMICS. 2014 MAY 8;7(1). 3rd Annual Translational Bioinformatics Conference (TBC) / ISCB-Asia, Seoul, SOUTH KOREA, OCT 02-04, 2013.
6. Zolbanin, Hamed Majidi and Delen, Dursun and Zadeh, Amir Hassan. Predicting overall survivability in comorbidity of cancers: A data mining approach [Article]. DECISION SUPPORT SYSTEMS. 2015 JUN;74:150–161.
7. Gordon L, Olshen RA. Tree-structured survival analysis. Cancer Treatment Reports. 1985;69(10):1065–1068. Cited By 97. Available from: <http://www.scopus.com/inward/record.url?eid=2-s2.0-0021875130&partnerID=40&md5=9e112ed840960f801b6260b23bf6811d>.
8. Bou-Hamad I, Larocque D, Ben-Ameur H. A review of survival trees. Statistics Surveys. 2011;5:44–71. Cited By 15. Available from: <http://www.scopus.com/inward/record.url?eid=2-s2.0-84857308440&partnerID=40&md5=f8af82017ade68e335fd258c6857bf49>.
9. Ishwaran H, Kogalur UB. Consistency of random survival forests. Statistics and Probability Letters. 2010;80(13-14):1056–1064. Cited By 26. Available from: <http://www.scopus.com/inward/record.url?eid=2-s2.0-77953020220&partnerID=40&md5=1e4478c51150f0159fdc6c1cb631968b>.
10. National Cancer Institute, the Surveillance, Epidemiology, and End Results Program. Documentation for ASCII Text Data Files - SEER Datasets; 2016 (accessed 15 Jan 2016). <http://seer.cancer.gov/data/documentation.html>.
11. Michael Bowles. Machine Learning in Python: Essential Techniques for Predictive Analysis. Wiley; 2015.