Some Lessons Learned from Using Machine Learning Methods to Predict Adverse Health Outcomes

Mark Craven

Department of Biostatistics & Medical Informatics Department of Computer Sciences University of Wisconsin





Predicting adverse health outcomes



Given a patient's clinical history up to a given point in time

Predict whether the patient will have a specific type of adverse event in the future

Learning task

Can we find a model that accurately represents outcome y as a function of feature vector x?



Example: predicting post-hospitalization VTE risk



Given a patient's clinical history up to EHR up to time of hospitalization **Predict** whether the patient is high risk for VTE after release from hospital

Predicting post-hospitalization VTE risk

- 720 subjects
- cases/controls determined by expert review
- two feature sets
 - 119 features based on 78 risk factors for VTE or thrombophilia
 - 3330 features in "unabridged" representation



Emily Kawaler





Peggy Peissig

Steven Yale

Assessing the accuracy of learned VTE models [Kawaler et al. AMIA 2012]



1. We can learn models that are more accurate than conventional risk assessment tools.

2. Learned models can identify novel risk/protective factors.

We compare three types of models for VTE task

curated representation (known risk factors)



unabridged representation (all features)







conventional risk assessment tools



Survival curves for models that stratify patients for post-hospitalization VTE risk



Predicting asthma exacerbations

- 28,101 subjects with prior history of asthma •
- exacerbations phenotyped by a rule: an urgent visit to a healthcare provider for asthma • symptoms followed by treatment with oral corticosteroids
- features represent •
 - demographics •
 - diagnoses, problem-list diagnoses ٠
 - medications ٠
 - vitals •
 - asthma control scores •
 - prior exacerbations ٠



Akshay Sood

Alex Cobian





Theresa Guilbert Lawrence Hanrahan



Assessing the accuracy of learned exacerbation prediction models



3. In some applications, we can learn more accurate models by using patient genetics in addition to clinical variables.

Predicting breast cancer risk

- 738 subjects for which both genetics and mammogram data was available
- cases/controls determined by cancer registry
- excluded cases with known BRCA1 or BRCA2 mutations
- features
 - 49 variables from BI-RADS lexicon
 - 77 high-frequency/low-penetrance SNPs identified in GWAS



Beth Burnside



Ming Yuan



Jun Fan





Shara Feld

Improved breast-cancer risk prediction with structureleveraged methods [Fan et al. *JMLR* 2016]

49 mammography descriptors



learning method optimizes:

$$\min_{(\alpha,\beta)\in R^{d+1}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \log(1 + \exp[-y_i(\alpha + x_i^T \beta)]) + \lambda_1 \sum_{g=1}^{G} \sqrt{d_g} \left\| \beta_g \right\|_2 + \lambda_2 \sum_{g=1}^{G} \sum_{j=1}^{d_g} \left\| \beta_g^j - \beta_g^{j+1} \right\|_p^p \right\}$$
logistic regression
group lasso
(grouping structure)
(dependence structure)

Assessing the accuracy of learned breast-cancer risk prediction models



Assessing the accuracy of learned breast-cancer risk prediction models [Feld et al. *AMIA Informatics Summit* 2018]

the predictive value of genetics depends on the age of subjects



4. Simple, linear models are often competitive with fancier models.

Assessing the accuracy of learned exacerbation prediction models





ARTICLE OPEN Scalable and accurate deep learning with electronic health records

Alvin Rajkomar ^{1,2}, Eyal Oren¹, Kai Chen¹, Andrew M. Dai¹, Nissan Hajaj¹, Michaela Hardt¹, Peter J. Liu¹, Xiaobing Liu¹, Jake Marcus¹, Mimi Sun¹, Patrik Sundberg¹, Hector Yee¹, Kun Zhang¹, Yi Zhang¹, Gerardo Flores¹, Gavin E. Duggan¹, Jamie Irvine¹, Quoc Le¹, Kurt Litsch¹, Alexander Mossin¹, Justin Tansuwan¹, De Wang¹, James Wexler¹, Jimbo Wilson¹, Dana Ludwig², Samuel L. Volchenboum³, Katherine Chou¹, Michael Pearson¹, Srinivasan Madabushi¹, Nigam H. Shah⁴, Atul J. Butte², Michael D. Howell¹, Claire Cui¹, Greg S. Corrado¹ and Jeffrey Dean¹

Predictive modeling with electronic health record (EHR) data is anticipated to drive personalized medicine and improve healthcare quality. Constructing predictive statistical models typically requires extraction of curated predictor variables from normalized EHR data, a labor-intensive process that discards the vast majority of information in each patient's record. We propose a representation of patients' entire raw EHR records based on the Fast Healthcare Interoperability Resources (FHIR) format. We demonstrate that deep learning methods using this representation are capable of accurately predicting multiple medical events from multiple centers without site-specific data harmonization. We validated our approach using de-identified EHR data from two US academic medical centers with 216,221 adult patients hospitalized for at least 24 h. In the sequential format we propose, this volume of EHR data unrolled into a total of 46,864,534,945 data points, including clinical notes. Deep learning models achieved high accuracy for tasks such as predicting: in-hospital mortality (area under the receiver operator curve [AUROC] across sites 0.93–0.94), 30-day unplanned readmission (AUROC 0.75–0.76), prolonged length of stay (AUROC 0.85–0.86), and all of a patient's final discharge diagnoses (frequency-weighted AUROC 0.90). These models outperformed traditional, clinically-used predictive models in all cases. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. In a case study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

npj Digital Medicine (2018)1:18; doi:10.1038/s41746-018-0029-1

ARTICLE OPEN Scalable and accurate deep learning with electronic health records

Supplemental Table 1: Prediction accuracy of each task of deep learning model compared to baselines

| | Hospital A | Hospital B |
|--|-------------------------|-------------------------|
| Inpatient Mortality, AUROC ¹ (95% CI) | | |
| Deep learning 24 hours after admission | 0.95 (0.94-0.96) | 0.93 (0.92-0.94) |
| Full feature enhanced baseline at 24 hours after admission | 0.93(0.92-0.95) | 0.91(0.89-0.92) |
| Full feature simple baseline at 24 hours after admission | 0.93(0.91-0.94) | 0.90(0.88-0.92) |
| Baseline ($aEWS^2$) at 24 hours after admission | 0.85(0.81-0.89) | 0.86(0.83-0.88) |
| 30-day Readmission, AUROC (95% CI) | | |
| Deep learning at discharge | 0.77 (0.75-0.78) | 0.76(0.75-0.77) |
| Full feature enhanced baseline at discharge | 0.75(0.73-0.76) | 0.75(0.74-0.76) |
| Full feature simple baseline at discharge | 0.74(0.73-0.76) | 0.73(0.72 - 0.74) |
| Baseline (mHOSPITAL ³) at discharge | 0.70(0.68-0.72) | 0.68(0.67 - 0.69) |
| Length of Stay at least 7 days AUROC (95% CI) | | |
| Deep learning 24 hours after admission | 0.86 (0.86-0.87) | 0.85(0.85-0.86) |
| Full feature enhanced baseline at 24 hours after admission | 0.85(0.84-0.85) | 0.83(0.83-0.84) |
| Full feature simple baseline at 24 hours after admission | 0.83(0.82-0.84) | 0.81(0.80-0.82) |
| Baseline (mLiu ^{4}) at 24 hours after admission | 0.76(0.75-0.77) | 0.74(0.73-0.75) |

"baseline" models are regularized logistic regression; nearly as good as deep networks for all tasks 5. We can we gain some understanding from complicated learned models (e.g. random forests, neural networks).

Interpreting black-box models [Craven & Shavlik, *NeurIPS* 1996; Lee et al., *AAAI* 2019]



Kyubin Lee



Akshay Sood

Identifying important features via perturbations

 Feature importance can be ascertained by perturbing a feature and measuring its effect on the model loss



foreach test instance $(\mathbf{x}^{(i)}, y^{(i)})$ **do** $\begin{vmatrix} \text{let } \Delta \mathbf{x}_{j}^{(i)} \text{ represent } \mathbf{x}^{(i)} \text{ with feature } j \text{ perturbed in some way} \\ \text{compare loss } L\left[y^{(i)}, f\left(\mathbf{x}^{(i)}\right)\right] \text{ to } L\left[y^{(i)}, f\left(\Delta \mathbf{x}_{j}^{(i)}\right)\right] \end{vmatrix}$

From features to hierarchies over features

• In many domains, there is hierarchical structure over features



False Discovery Rate Control

- There is a **multiple comparisons problem** due to the large number of base features and feature groups under test
- We use hierarchical FDR control [Yekutieli JASA 2008] to identify the finest level of resolution at which we can find features that have a statistically significant effect on model loss

Analysis of learned asthma exacerbation model

- 1. use standard feature hierarchies (e.g. ICD-9)
- 2. do perturbations by erasure
- 3. calculate *p*-values for all nodes in the hierarchy
- 4. apply hierarchical FDR control



Identifying important features for asthma exacerbation prediction model







Conclusions

- learned models are sometimes more accurate than conventional risk assessment tools
- learned models can identify novel risk/protective factors
- genetics can augment predictive value of clinical data in some cases
- simple models often work well
- but we have tools to gain insight from complex black-box models

NIH U54 AI117924 NIH T32 HG002760 NIH T15 LM007359 NIH R01 EY023292 NSF IIS 1218880





Learning associations between HSV-1 genotype and host phenotype [Lee et al. J. Virology 2016; Kolb et al. PLoS Pathogens 2016]



- 69 viral genomes
- features represent 547 haplotype blocks
- we have learned regression models using Lasso, ridge, and random forest approaches



Analysis of learned HSV-1 genotype-to-phenotype models



- 40 "outer" features/feature groups
- 34 base features

Important variables are localized in the viral genome







Using machine-learning to predict adverse outcomes

- 1. select cohort of subjects
- 2. retrospectively phenotype outcomes (e.g. determine who is a case/control)
- 3. define feature representation
- 4. learn models
- 5. evaluate models

Defining feature representations



An LSTM neural network for asthma exacerbations



diagnosesinterventionsvitals, prior(meds + procedures)exacerbations, etc.

An LSTM neural network model



Identifying important features via perturbation

• We can use **hypothesis testing** to characterize the effect of perturbing a feature

| Instance | Perturbed Instance | Loss | Perturbed Loss |
|-------------------------------|---------------------------|--|----------------|
| $(\mathbf{x}^{(1)}, y^{(1)})$ | | $L\left[y^{(1)}, f\left(\mathbf{x}^{(1)}\right) ight]$ | |
| $(\mathbf{x}^{(2)},y^{(2)})$ | | $L\left[y^{(2)}, f\left(\mathbf{x}^{(2)}\right) ight]$ | |
| ••• | | • • • | |

Null hypothesis: Median change in loss when perturbing feature j is 0