
Measuring the Stability of EHR- and EKG-based Predictive Models

Andrew C. Miller*
Data Science Institute
Columbia University

Ziad Obermeyer
School of Public Health
University of California at Berkeley

Sendhil Mullainathan
Booth School of Business
University of Chicago

Abstract

Databases of electronic health records (EHRs) are increasingly used to inform clinical decisions. Machine learning methods can find patterns in EHRs that are predictive of future adverse outcomes. However, statistical models may be built upon patterns of health-seeking behavior that vary across patient subpopulations, leading to poor predictive performance when training on one patient population and predicting on another. This note proposes two tests to better measure and understand model generalization. We use these tests to compare models derived from two data sources: (i) historical medical records, and (ii) electrocardiogram (EKG) waveforms. In a predictive task, we show that EKG-based models can be more stable than EHR-based models across different patient populations.

1 Introduction

Patient history-based risk scores can inform physician decision-making and alter the course of treatment. The Framingham risk score for long-term coronary heart disease is one such an example [1]. In the past decade, large volumes of observational health care data — administrative claims, digitized physician notes, and laboratory test values among others — are now commonly used to build predictive models for a range of clinical purposes, including tracking disease progression [2], predicting near-term ICU interventions [3] or the onset of sepsis [4], and are incorporated into standardized prediction frameworks [5].

Patterns in observational health data, by construction, reflect health-seeking behavior in the population served. Consequently, predictive accuracy of statistical models can vary significantly if the underlying patterns of patient behavior vary in smaller groups within the population [6]. Further, under-served populations will, by definition, have less historical data on which to base predictive algorithms. Waveform and image data, such as electrocardiograms (EKGs) and echocardiograms, are a complementary source of information that can be used to build models of patient risk. EKGs, for instance, measure a patient’s cardiac function, which may be predictive of future disease. Patterns in a patient’s EKG (which directly measures cardiac activity) that are predictive of heart failure may be more *portable* than patterns in patient medical records across patient populations.

In this note, we study the generalization performance of predictive algorithms trained on one population and tested on another. We propose two statistics for measuring feature portability. As a case study, we construct models that predict the outcome of common lab test used to confirm heart attack — troponin levels — using (i) past diagnoses and medications and (ii) electrocardiogram waveform data. We show that EKG-based predictors are more stable across two sub-populations of patients — those that use the health care system with high frequency and those that use the system with lower frequency. For each prediction algorithm (e.g. data-source and training population), we examine the sources of distribution shift that can lead to poor generalization. We conclude with a discussion of future research directions.

*am5171@columbia.edu, <http://andymiller.github.io/>

2 Prediction and Conditional Stability

Our goal is to construct an accurate predictive model for some outcome $y \in \{0, 1\}$ (e.g. an adverse cardiac event) given some set of patient-specific data $\mathbf{x} \in \mathbb{R}^D$ (e.g. historical medical records or raw EKG waveforms). We observe a set of these data drawn from a population distribution P

$$\mathbf{x}_n, y_n \sim P \quad \text{population distribution} \quad (1)$$

$$\mathcal{D}_{P,N} \triangleq \{\mathbf{x}_n, y_n\}_{n=1}^N \quad P\text{-distributed dataset of size } N. \quad (2)$$

We use these data to train a model that predicts the conditional probability of $y = 1$ given \mathbf{x}

$$m_{\mathcal{D}_{P,N}} \leftarrow \text{train}(\mathcal{D}_{P,N}) \quad (3)$$

where the model $m_{\mathcal{D}_{P,N}}(\mathbf{x}) \triangleq P(y = 1|\mathbf{x})$ approximates the conditional distribution. For example, if \mathbf{x} are medical record features, $m_{\mathcal{D}_{P,N}}$ may be a logistic regression model fit with maximum likelihood. If \mathbf{x} are EKG waveforms, $m_{\mathcal{D}_{P,N}}$ may be a convolutional neural network.

Model Portability How well does $m_{\mathcal{D}_{P,N}}$ perform on another distribution, $\mathbf{x}, y \sim Q$? For example, the Q may be another hospital, health care system, or patient sub-population. For a specific model $m_{\mathcal{D}_P}$ (e.g. logistic regression with EHR data) and a *different* sub-population Q , we can measure cross-generalization with the area under the ROC curve

$$G^{(m)}(\mathcal{D}_P, \mathcal{D}_Q) = \text{AUC}(m_{\mathcal{D}_P}, \mathcal{D}_Q) \quad (4)$$

where an entry $\mathcal{D}_P, \mathcal{D}_Q$ denotes a model trained on population P and tested on population Q . The diagonal entries $G^{(m)}(\mathcal{D}_P, \mathcal{D}'_P)$ measure the standard notion of generalization with no distribution shift — predictive performance on held out data from the same distribution. The off-diagonal entries $G^{(m)}(\mathcal{D}_P, \mathcal{D}_Q)$ measure performance under a new test distribution, a problem addressed by domain adaption techniques [7]. When generalization performance for population P and population Q are different, a natural question is what is driving that gap? What features are portable?

Distribution shift Why do predictive features \mathbf{x} fail to generalize? Consider a predictive model trained on data from population P , which we denote with shorthand $m_P(\mathbf{x})$. The model $m_P(\mathbf{x})$ is just a function of covariates \mathbf{x} . When \mathbf{x} is drawn from P , $m_P(\mathbf{x})$ and y have a joint distribution, induced by $P(\mathbf{x}, y)$. When \mathbf{x} is drawn from a different population Q , $m_P(\mathbf{x})$ and y have a *different* joint distribution, induced by $Q(\mathbf{x}, y)$. For shorthand, define $\hat{y}_P = m_P(\mathbf{x})$. The conditional distributions induced by P , Q , and $m_P(\mathbf{x})$ can tell us what is “stable” and “unstable”

- $P(\hat{y}_P | y)$ vs. $Q(\hat{y}_P | y)$: covariate stability
- $P(y | \hat{y}_P)$ vs. $Q(y | \hat{y}_P)$: predictive stability

If covariates are conditionally unstable, this suggests different patterns of health-seeking behavior (or physiological differences) can be driving the generalization gap. If outcomes y are conditionally unstable, this suggests that similar patterns in \mathbf{x} do not suggest y the same way in population P and population Q — there is something special about the prediction function tuned to population P that does not apply to population Q .

We summarize *covariate stability* by estimating the distance between conditional distributions

$$CS(m_P, y, P, Q) = \text{KS}(P(\hat{y}_P|y); Q(\hat{y}_P|y)) \quad \text{covariate stability} \quad (5)$$

where $\text{KS}(\cdot; \cdot)$ is the Kolmogorov-Smirnov distance between the two distributions — we use a two-sample estimator on held out data for this distance [8]. Intuitively, this will be low when covariates from the same class have the same distribution.

We summarize *predictive stability* with the difference in conditional expectation of y given \hat{y}_P

$$PS(\hat{y}; m_P, y, P, Q) = \mathbb{E}_{y, \hat{y} \sim P} [y|\hat{y}] - \mathbb{E}_{y', \hat{y}' \sim Q} [y'|\hat{y}'] \quad (6)$$

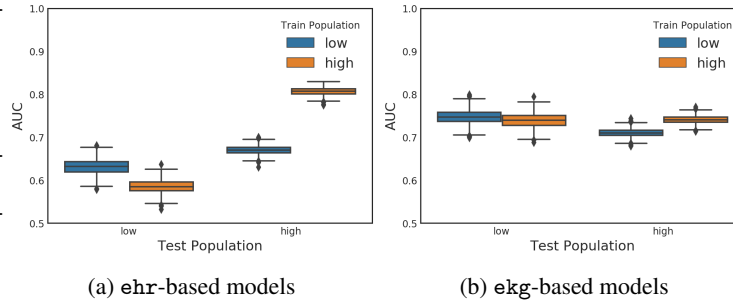
This defines a function of \hat{y} — for each value of \hat{y} the conditional distribution of outcome y may be similar or different between Q and P . We visualize an estimate of this whole function, and compute a one dimensional summary (by averaging over an equal mixture of the P and Q distributions).²

²The notation \hat{y}_P , $P(\hat{y}_P, y)$ and $Q(\hat{y}_P, y)$ can appear confusing — there are really three distributions to keep in mind: the distribution of the model training set (P), and the two distributions being compared (in our examples P again and Q .)

healthcare use	split	trop frac. pos	# pos	# obs	age mean	std	frac. female mean
low	train	0.119	832	7000	58.4	18.4	0.459
	val	0.116	116	1000	58.4	18.2	0.447
	test	0.120	276	2298	58.2	18.4	0.438
high	train	0.177	1236	7000	63.9	15.3	0.583
	val	0.217	217	1000	62.2	17.1	0.533
	test	0.210	941	4491	63.5	15.5	0.555

Table 1: Data summary. We train prediction functions using equally sized populations of low and high usage patients (7,000 train, 1,000 validation). We report test statistics on held-out patients.

Figure 1: Predictive performance of each classifier by AUC. The test population is labeled on the horizontal axis; the train population is labeled by color — blue indicates trained on low-use (Q) patients, and orange indicates trained on high-use (P) patients.



3 Experiments

We build models that predict the outcome of a troponin lab test in a cohort of emergency department patients. A troponin test measures the troponin T or I protein levels, which are released when heart muscle is damaged during a heart attack. Troponin lab tests are used to determine if a patient is currently or about to suffer a heart attack. Based on the guidelines in [9], we define high to be a value of .04 ng/mL or greater. We label a patient “positive” ($y = 1$) if they have received a troponin result greater than or equal to .04 ng/mL within the seven days following their emergency department admission. A patient is labeled “negative” ($y = 0$) if they received a troponin result less than .04 ng/mL within the same time period. An algorithm that reliably predicts troponin outcome in the near term can be useful for assessing a patient’s risk of heart attack.

We study the stability of predictive models derived from two distinct sources of covariate (x) data:

- **ehr**: all diagnoses and medications administered over the past year, derived from ICD-9 codes collected by a network of hospitals. We represent a patient’s one-year history with a sparse binary vector of 2,372 diagnoses and medications — a 1 indicates that the diagnosis/medication was present in the past year. We fit predictive models using ehr data with XGBoost [10], a gradient boosted decision tree algorithm, tuning tree depth using a validation set.
- **ekg**: raw electrocardiogram waveforms from an EKG administered at the beginning of the emergency department visit. For the ekg data, we use the convolutional residual network described in [11] on the three “long leads” (II, V1, V5). We use stochastic gradient descent, keeping the model with the best predictive performance on a validation set.

Following the notation above, we study ehr- and ekg-based predictors over two populations:

- $P \triangleq$ *high-use patients*: patients with more than 8 visits to any department within the network of hospitals in the year before the emergency department encounter.
- $Q \triangleq$ *low-use patients*: patients with 8 or fewer such visits in the same time frame.

Summary statistics of each population are presented in Table 1. We compare models derived from four combinations: trained on $\{P, Q\}$ using features $\{\text{ehr}, \text{ekg}\}$.

Generalization We present the *within-distribution* generalization performance (e.g. $G^{(m)}(P, P)$) and the *between-distribution* performance (e.g. $G^{(m)}(P, Q)$) in Figure 1. The ehr model trained on the low-use (Q) population generalizes slightly better on new patients from Q than the ehr model trained on the high-use (P) population. The ehr model trained on P (high-use) generalizes significantly better on patients from P than the model trained on patients from Q , as summarized in Figure 1a. On the other hand, we see that models trained on ekg data from either P or Q have similar predictive performance on both the P and Q distributions, as seen in Figure 1b.

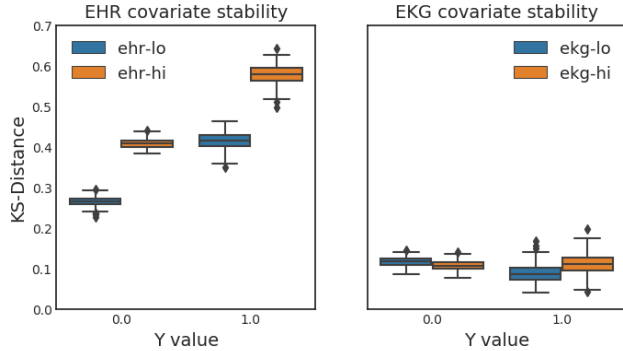
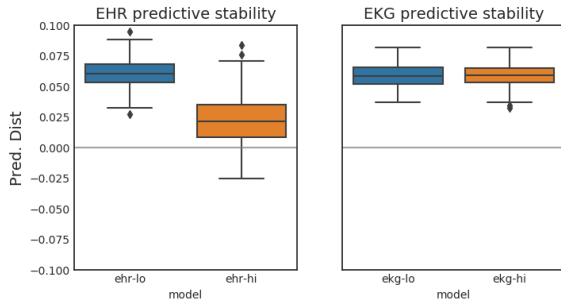


Figure 2: Distribution comparison summary of covariate stability. The left plot compares models derived from EHR features. Each entry estimates the statistical distance described in Equation 5 — in words, for true value y (and accompanying x , how different are the distributions $\hat{y}(x)$ when $x \sim P$ and $x \sim Q$? We see EKG features are far more stable under this measure (lower is better).

Figure 3: Predictive stability. We summarize Equation 6 for ehr- (left) and ekg- (right) based models. The vertical axis measures the difference between $E_Q[y|\hat{y}] - E_P[y|\hat{y}]$ averaged over bins of \hat{y} . Blue corresponds to models trained on low-use patients (Q) and orange to models trained on high-use patients (P). For three models, we see a detectable bias, suggesting the functional relationship between x and y may differ slightly between groups.



Covariate stability Figure 2 shows the covariate stability statistics for the two sets of models, trained on the two populations. We can see that the covariate distributions $P(\hat{y}|y)$ and $Q(\hat{y}|y)$ are significantly closer when \hat{y} is based on ekg features than ehr features. This strongly suggests that behavioral characteristics between high- and low-use patients *with the same outcome* are driving a statistically significant difference in diagnosis and medication covariates x . The ekg features, on the other hand, are significantly more stable between the high- and low-use patients. Intuitively, physiological measurements are not subject to the same biases that human behavior and decision-making introduce. The distributions of \hat{y} given $y = 0$ and $y = 1$ for both the ehr and ekg models are depicted in Appendix Figure 4.

Predictive stability We present our one-dimensional predictive stability summary of Equation 6 in Figure 3, and depict full estimates in Figure . The average deviation between $E_Q[y|\hat{y}]$ and $E_P[y|\hat{y}]$ are similar for the ekg and ehr models and small, though generally non-zero (on the order of .05 to .08). Both sets of models can make modest improvements, but this source of instability does not appear to be driving a large difference in generalization performance. More details are in Figure 5.

4 Discussion

We proposed two tests to better understand the portability of ehr and ekg features between patient populations. Our analysis suggests that EKG features can be stable across low- and high-use patients, whereas diagnoses and medications are not. We see multiple avenues of future work. Given the observed covariate instability, we hope to model patient health-seeking behavior which lead to statistical differences in ehr data. We also hope that models for patterns of missing data can be effectively combined with ekg data to construct a predictor that is both accurate and portable.

Additionally, it is presciently noted in [6] that a strategy for making predictors more equitable is to pay the cost of additional (and focused) data collection. Our analysis suggests a variant of this remedy — find sources of data that are more likely to generalize across populations. Physician notes, for example, are likely to be influenced by physician-specific biases, but are still may be more portable than billing records — a notion we hope to quantify and measure. EKGs and other physiological signals are not subject to the same patterns of presence and missingness in health record data, which is driven by unobserved human behavior and decision-making.

References

- [1] Peter WF Wilson, Ralph B D’Agostino, Daniel Levy, Albert M Belanger, Halit Silbershatz, and William B Kannel. Prediction of coronary heart disease using risk factor categories. *Circulation*, 97(18):1837–1847, 1998.
- [2] Joseph A Cruz and David S Wishart. Applications of machine learning in cancer prediction and prognosis. *Cancer informatics*, 2:117693510600200030, 2006.
- [3] Marzyeh Ghassemi, Mike Wu, Michael C Hughes, Peter Szolovits, and Finale Doshi-Velez. Predicting intervention onset in the icu with switching state space models. *AMIA Summits on Translational Science Proceedings*, 2017:82, 2017.
- [4] Katharine E Henry, David N Hager, Peter J Pronovost, and Suchi Saria. A targeted real-time early warning score (trewscore) for septic shock. *Science translational medicine*, 7(299):299ra122–299ra122, 2015.
- [5] Jenna M Reps, Martijn J Schuemie, Marc A Suchard, Patrick B Ryan, and Peter R Rijnbeek. Design and implementation of a standardized framework to generate and evaluate patient-level prediction models using observational healthcare data. *Journal of the American Medical Informatics Association*.
- [6] Irene Chen, Fredrik D Johansson, and David Sontag. Why is my classifier discriminatory? *arXiv preprint arXiv:1805.12002*, 2018.
- [7] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1-2):151–175, 2010.
- [8] Michael A Stephens. Edf statistics for goodness of fit and some comparisons. *Journal of the American statistical Association*, 69(347):730–737, 1974.
- [9] L Kristin Newby, Robert L Jesse, Joseph D Babb, Robert H Christenson, Thomas M De Fer, George A Diamond, Francis M Fesmire, Stephen A Geraci, Bernard J Gersh, Greg C Larsen, et al. Accf 2012 expert consensus document on practical clinical considerations in the interpretation of troponin elevations: a report of the american college of cardiology foundation task force on clinical expert consensus documents. *Journal of the American College of Cardiology*, 60(23):2427–2463, 2012.
- [10] Tianqi Chen and Carlos Guestrin. XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’16*, pages 785–794, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-4232-2. doi: 10.1145/2939672.2939785. URL <http://doi.acm.org/10.1145/2939672.2939785>.
- [11] Pranav Rajpurkar, Awni Y Hannun, Masoumeh Haghpanahi, Codie Bourn, and Andrew Y Ng. Cardiologist-level arrhythmia detection with convolutional neural networks. *arXiv preprint arXiv:1707.01836*, 2017.

A Additional Figures

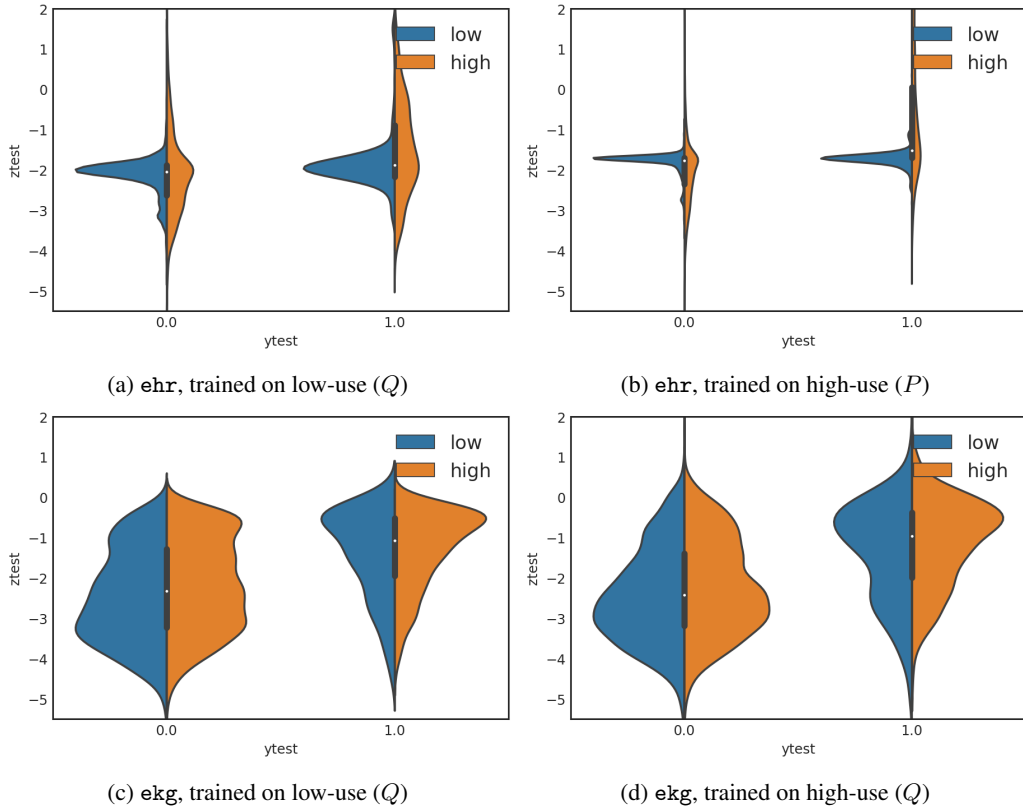


Figure 4: Covariate stability. Above we visualize the conditional distribution of \hat{y} given the true label $y = 0$ or $y = 1$ (denoted along the horizontal axis). Within the $y = 0$ or $y = 1$ subset, the blue distributions denote the predictive \hat{y} distribution among the $Q = \text{low}$ use patients, and the orange denotes the predictive \hat{y} distribution among the $P = \text{high}$ use patients. Intuitively, a stable predictor would have a similar distribution of \hat{y} across P and Q (but within the group true $y = 1$ or $y = 0$). For example (a) shows the conditional distribution of $\hat{y}|y = 0$ for the Q (blue) and P (orange) populations, for \hat{y} trained on ehr features from population Q . (b) depicts the same conditional distributions for ehr features trained on P . (c) and (d) depict the same breakdown for models trained on ekg features. The covariates are “stable” when the conditional distributions match — we see the ekg-based covariates are much more stable than ehr-based covariates. This similarity is summarized by the covariate similarity statistic Figure 2.

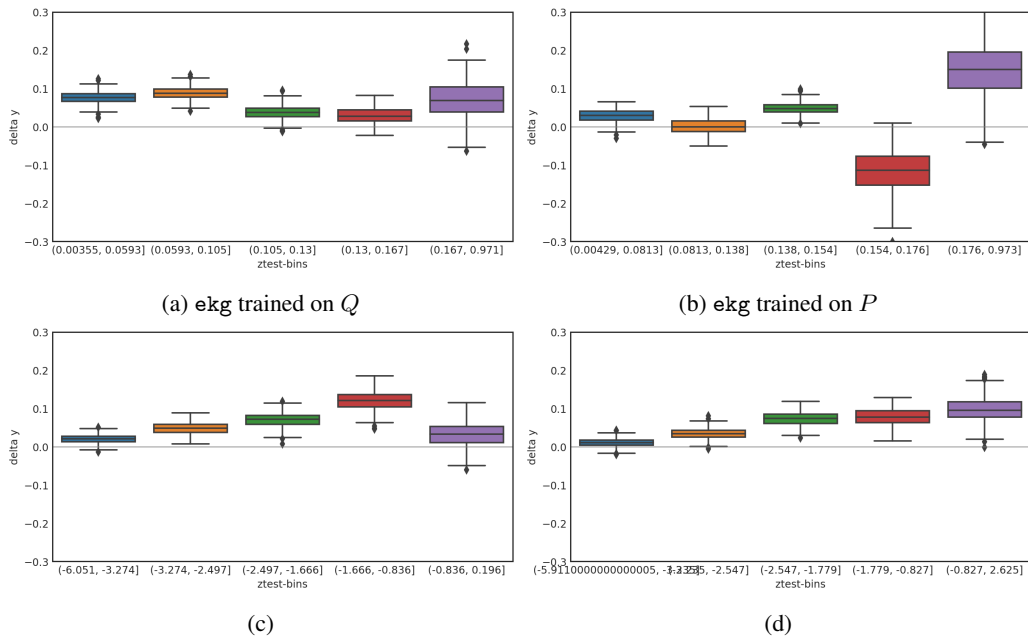


Figure 5: Predictive stability. Here we depict estimates of the predictive stability curve described in Equation 6. Each bin along the horizontal axis reflects a quintile of test data, stratified by \hat{y} . Here we see a more complete picture — the ehr-based predictors tend to be quite predictively stable. The ekg-based predictors appear to be slightly less stable, which may indicate slight differences in the physiology of low-use Q and high-use P patients (for instance, high-use patients are on average a bit older and male in this sample).