
Research and Applications

Modeling physician variability to prioritize relevant medical record information

Mohammadamin Tajgardoon ¹, Gregory F. Cooper,^{1,2} Andrew J. King,³ Gilles Clermont,³ Harry Hochheiser ^{1,2}, Milos Hauskrecht,^{1,4} Dean F. Sittig ⁵, and Shyam Visweswaran^{1,2}

¹Intelligent Systems Program, School of Computing and Information, University of Pittsburgh, Pittsburgh, Pennsylvania, USA, ²Department of Biomedical Informatics, University of Pittsburgh, Pittsburgh, Pennsylvania, USA, ³Department of Critical Care Medicine, University of Pittsburgh, Pittsburgh, Pennsylvania, USA, ⁴Department of Computer Science, University of Pittsburgh, Pittsburgh, Pennsylvania, USA and ⁵Department of Biomedical Informatics, University of Texas Health Science Center at Houston, Houston, Texas, USA

Corresponding Author: Shyam Visweswaran, MD, PhD, The Offices at Baum, 5607 Baum Blvd., Suite 523, Pittsburgh, PA 15206, USA (shv3@pitt.edu)

Received 14 August 2020; Revised 5 October 2020; Accepted 2 November 2020

ABSTRACT

Objective: Patient information can be retrieved more efficiently in electronic medical record (EMR) systems by using machine learning models that predict which information a physician will seek in a clinical context. However, information-seeking behavior varies across EMR users. To explicitly account for this variability, we derived hierarchical models and compared their performance to nonhierarchical models in identifying relevant patient information in intensive care unit (ICU) cases.

Materials and methods: Critical care physicians reviewed ICU patient cases and selected data items relevant for presenting at morning rounds. Using patient EMR data as predictors, we derived hierarchical logistic regression (HLR) and standard logistic regression (LR) models to predict their relevance.

Results: In 73 pairs of HLR and LR models, the HLR models achieved an area under the receiver operating characteristic curve of 0.81, 95% confidence interval (CI) [0.80–0.82], which was statistically significantly higher than that of LR models (0.75, 95% CI [0.74–0.76]). Further, the HLR models achieved statistically significantly lower expected calibration error (0.07, 95% CI [0.06–0.08]) than LR models (0.16, 95% CI [0.14–0.17]).

Discussion: The physician reviewers demonstrated variability in selecting relevant data. Our results show that HLR models perform significantly better than LR models with respect to both discrimination and calibration. This is likely due to explicitly modeling physician-related variability.

Conclusion: Hierarchical models can yield better performance when there is physician-related variability as in the case of identifying relevant information in the EMR.

Key words: electronic medical records, information-seeking behavior, machine learning, physician variability, hierarchical modeling

LAY SUMMARY

In current electronic medical record (EMR) systems, large amounts of information are frequently collected in complex clinical settings such as the intensive care unit, which makes it challenging for physicians to identify relevant patient information. To make EMR systems more useful, one approach is to augment them with the capability to intelligently highlight patient information that is relevant to the clinical task. In an effort to address this challenge, we developed a learning EMR system that uses machine learning to leverage physicians' information-seeking behavior to identify and highlight the relevant information. One challenge for machine learning is the variability in information-seeking behavior across physicians that typically used models may not be able to capture. In this article, we show that hierarchical models such as hierarchical logistic regression perform better at identifying relevant patient information compared to logistic regression. The hierarchical model learns a distinct model for each physician while sharing useful knowledge learned from one physician with the rest of the physicians in order to improve performance.

INTRODUCTION

A key source of frustration with electronic medical record (EMR) systems stems from the inability to retrieve relevant patient information efficiently.¹⁻⁴ Current EMR systems do not possess sophisticated search capability or do they prioritize patient information relative to the clinical task at hand.^{5,6} The inability to identify relevant patient information can lead to poor care and medical errors.⁷⁻⁹ Further, in complex clinical environments, such as the intensive care unit (ICU), large quantities of data per patient accumulate rapidly,¹⁰ which can exacerbate information retrieval challenges. EMR systems that prioritize the display of relevant patient information are therefore needed to minimize the time and effort that physicians spend in identifying relevant information.

Various solutions have been proposed for effective prioritization and display of patient information in EMR systems,¹¹⁻¹⁴ most of which are based on rules that have been developed to customize and organize the display of patient information. In contrast to rule-based approaches, we developed and evaluated a data-driven approach called the learning EMR (LEMUR) system in a prior study.^{15,16} The LEMUR system tracks physician information-seeking behavior and uses it to learn machine learning models that predict which information is relevant in a given clinical context. Those predictions are used to highlight the relevant data in the EMR system to draw a physician's attention.

However, information-seeking behavior has been shown to vary across individual physicians as well as across EMR system user types such as physicians, nurses, and pharmacists.^{1,5} In this study, we use hierarchical models to explicitly model this variability because such models have been shown to be useful when the data are collected from subjects with different behaviors.¹⁷ In particular, we compare the performance of hierarchical logistic regression (HLR) models and standard logistic regression (LR) models in predicting relevant patient information in a LEMUR system.

The remainder of this article is organized as follows. In the Background section, we review the LEMUR system, briefly describe hierarchical models, and describe prior work on physician-related variability. In the Methods section, we describe the data collection and preparation, the experimental details, and the evaluation measures. We present the results of the experiments in the Results section and close with the Discussion and Conclusion sections.

BACKGROUND

In this section, we provide brief descriptions of the LEMUR system, hierarchical models, and past studies that have examined physician-related variability.

The LEMUR system

The LEMUR system uses a data-driven approach to prioritize patient information that is relevant in the context of a clinical task.^{15,16} The system uses machine learning to automatically identify and highlight relevant patient information for a specified task, for example, the task of summarizing a patient's clinical status at morning rounds in the ICU. In ICU morning rounds, the clinical team reviews pertinent information and the status of each patient; for each patient, one team member reviews information in the EMR system and orally presents a summary of the patient's clinical status to the team. Reviewing and identifying relevant patient information, called *prerounding*, is time-consuming and laborious. The goal of the LEMUR system is to use machine learning to automatically identify and highlight the relevant information required for a given clinical task such as *prerounding*. The predictive models of the LEMUR system are derived using the information-seeking behavior of physicians when they search for relevant information in the EMR in the context of the clinical task. In particular, 11 critical care physicians reviewed the EMRs of ICU patients and marked the information that was relevant to *prerounding*, and predictive models were developed from this data.

Hierarchical models

Hierarchical models, also known as *multilevel* models, are useful in modeling hierarchically structured data because they can capture variability at different levels of the hierarchy.¹⁷ For example, consider predicting the mortality rate in a hospital with several units, such as critical care, general medical care, and emergency care. The data have a two-level hierarchical structure with the hospital at the first level and the units at the second level of the hierarchy. The overall mortality rate at the hospital level is obtained by combining the unit-level mortality rates in some fashion. A hierarchical model explicitly estimates the variability of the mortality rates across the units and uses those estimates to derive the hospital-level mortality rate, which can result in a better estimate of the overall mortality rate compared to using nonhierarchical models.

In a similar fashion, the information-seeking data used to develop the LEMUR models has a two-level hierarchical structure, where the top level corresponds to data that denote the entire *population* of physician reviewers and the bottom level corresponds to data that denote individual physicians. For specific patient information such as serum creatinine, its relevance is expected to differ across physician reviewers. A hierarchical model of the LEMUR data explicitly captures this variability that is likely to be useful in deriving more accurate predictive models.

Physician-related variability

Physician-related variability in healthcare outcomes has been of interest for decades, going back to the 1970s with studies reporting the effects of geographic location on clinical outcomes such as mortality and length of stay.¹⁸ In particular, variation in individual physician characteristics and practice styles has been recognized as a source of variability in clinical outcomes after adjusting for the health status of patients and the quality of healthcare services.^{5,19–25} For example, variability in cesarean section rates has been attributed to physician practice style after controlling for patient characteristics and risk factors, the status of the medical facility, and physician years of experience.²⁵ A study concluded that variability across individual physicians may impact the quality of preference-sensitive critical care delivery.²⁰ A recent study analyzed physician search patterns in the EMR and uncovered considerable variation in information-seeking behavior.⁵ In general, hierarchical modeling has been applied in various clinical settings to account for physician-related variability where the data have a hierarchical structure and can be grouped by a variety of factors such as country, state, or hospital site.^{26–32}

METHODS

In this section, we first describe the data set and the data preparation steps. Then, we describe the experimental methods including the development and evaluation of predictive models.

Data set

A set of ICU patient cases with a diagnosis of either acute kidney failure (AKF; ICD-9 584.9 or 584.5; 93 cases) or acute respiratory failure (ARF; ICD-9 518.81; 85 cases) were selected randomly from patients who were admitted between June 2010 and May 2012 to an ICU at the University of Pittsburgh Medical Center. Eleven critical care medicine physicians reviewed the patient cases in the LEMR system and for each patient indicated which patient information was relevant to the task of prerounding in the ICU. The recruited reviewers included ICU fellows and attending clinicians from the Department of Critical Care Medicine at the University of Pittsburgh. Each physician reviewer was instructed to review up to 30 patient cases. The first four cases were the same for all reviewers and were used as burn-in cases; these cases were not included in the data set. The remaining cases were different for each reviewer and each physician reviewed and annotated as many cases as they could during one to two sessions that lasted a total of 4–6 hours. Because the cases had some variation in the amount of patient information they contained and the physicians varied in the speed of reviewing, all physicians did not review the same number of cases.

The data set consists of two sets of variables including the *predictor* variables (or *predictors*) and *target* variables (or *targets*) that we now describe in detail. Predictor variables include demographics, admitting diagnosis, vital signs, ventilator settings, input and output measurements, laboratory test results, and medication administration data. A few variables such as demographics and admitting diagnosis are *static*, that is, their values do not change during the ICU stay, while the remaining variables, which constitute the majority of the predictors, are *temporal* and have multiple values during the ICU stay. For example, *age (in years)* is a static predictor variable while *blood urea nitrogen (BUN)* is a temporal predictor variable as it is usually measured multiple times during an ICU stay.

Target variables include any data in the EMR, such as vital signs, ventilator settings, input and output measurements, laboratory test results, and medication administration data that a physician may annotate as relevant for the task of prerounding. A target variable can take either *relevant* or *not relevant* values. As an example, for a patient with AKF, BUN=*relevant* denotes that BUN was measured for the patient and was sought, found, and annotated by a physician as relevant. If BUN was measured for the patient but was not sought by a physician, then the target is denoted as BUN = *not relevant*. A target variable may be missing too; for example, when BUN is not measured for the patient, it would not be available for a physician to seek and find. We developed a predictive model for each target variable such as BUN that predicts whether it is relevant in a particular patient. To develop a BUN model, we used all predictor variables described in the previous paragraph and used only data in which the BUN target was not missing.

The difference between predictor and target variables is in the values they take; that is, a target variable takes values of either *relevant* or *not relevant*, whereas a predictor variable's values are the measured values that are recorded in the EMR. For example, when BUN is a predictor, it takes numeric values in milligrams per deciliter (mg/dL) unit (Note that BUN may have been measured multiple times for a patient case and therefore, take several numeric values. We summarize these values as a fixed-length vector as described in the Data preparation section.), whereas as a target variable, it takes a value of either relevant or not relevant. Consequently, a model for predicting whether or not BUN is relevant may contain numeric values for BUN as a predictor variable.

Data preparation

We transformed the data set into a representation that is amenable to the application of machine learning methods. In particular, for each temporal predictor variable, we generated between 4 and 36 features (feature expansion in Figure 1).

The number of features for a temporal predictor was based on (1) the data domain of the predictor variable (eg, medication administration or laboratory result) and (2) the type of the predictor variable (eg, nominal or continuous). For example, for each medication variable, we generated four features including an indicator of whether the drug is currently prescribed, the time elapsed between first administration and the current time, the time elapsed between the most recent administration and the current time, and the dose at the most recent administration. For each laboratory test result, vital sign, and ventilator setting, we generated up to 36 features including an indicator of whether the event or measurement ever occurred, the value of the most recent measurement, the highest value, the lowest value, the slope between the two most recent values, and 30 other features. More details on the feature expansion are given in reference 33.

The data set consisted of 178 patient cases and 1864 raw predictor variables. Feature expansion resulted in a total of 30 770 features. Since the dimensionality of the data was high, we reduced the number of features (feature reduction in Figure 1) by removing those features where the values were missing in every patient case, had the same value for every case (ie, had zero variance), or the values were duplicates of another variable. Feature reduction resulted in a total of 6935 features.

We selected as target variables 73 EMR data items that had been annotated as relevant (positive) in 9 or more patient cases. [Supplementary Table S1](#) contains the list of target variables along with the

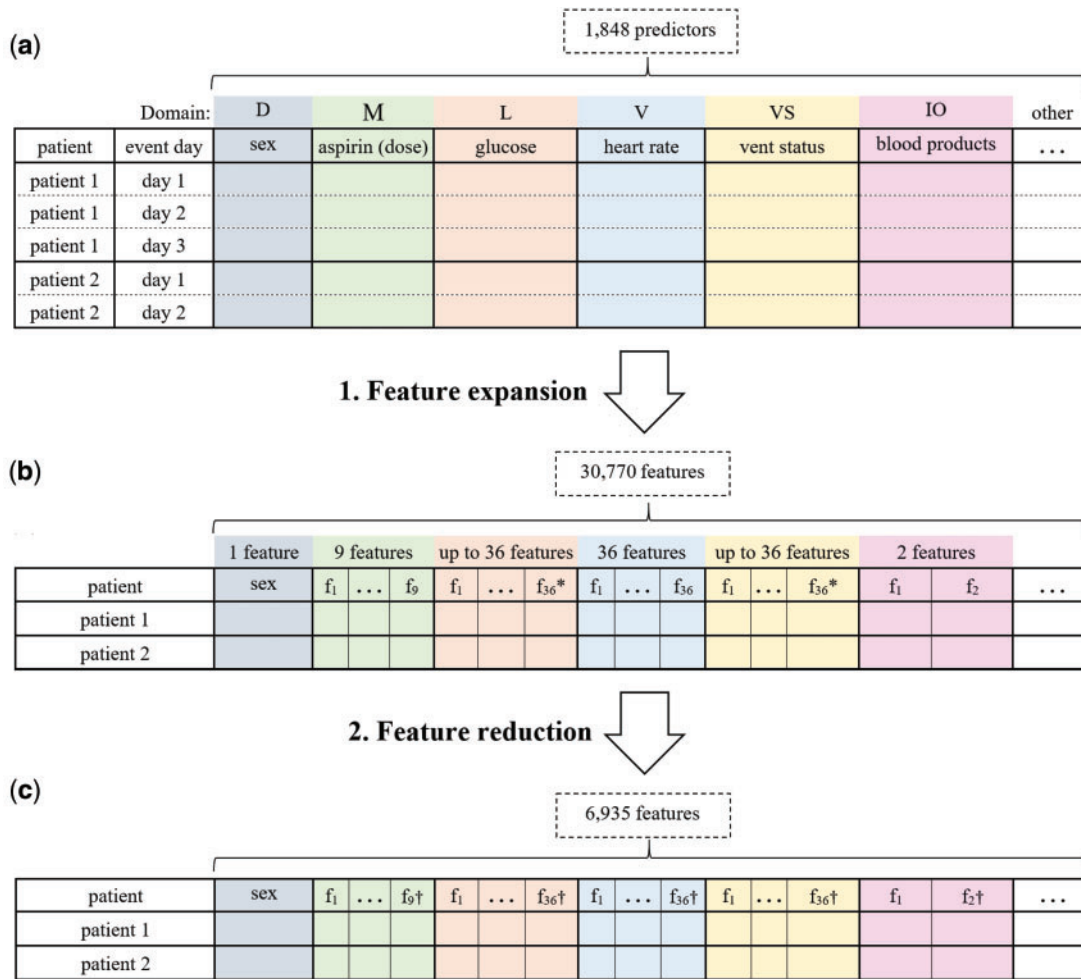


Figure 1. Steps in preparing the predictor variables. (A) The predictor variables for two example patients as measurements with one row per day are presented. The colors represent data domains; D: demographics, M: medication administrations, L: laboratory test results, V: vital signs, VS: ventilator settings, IO: input/output, and other: other domains. (B) The result of expanding the temporal predictor variables (total = 1848) to features (total = 30 770) is shown. This step flattens the data so that a patient that is represented by multiple rows is now represented by a single row. * denotes that the number of expanded predictors differs depending on the predictor value type (eg, nominal or continuous). (C) The features after feature reduction are shown, in which the number of features is reduced to 6935 features. † indicates that the number of features may be different for each variable in the domain.

number of cases in which each target variable was relevant, as well as the number of cases where the target variable was available for selection (ie, the value was not *missing*).

Experimental methods

Predictive models

An HLR model is a generalization of a standard logistic regression model in which the data are clustered into groups and the model intercept and coefficients can vary by group.¹⁷ Figure 2 shows the structure of a 2-level HLR model in which the LEMR data are clustered into groups of patient cases reviewed by each physician. Parameters at the lower level represent the physician-level models for the 11 physician reviewers, and parameters at the upper level represent the model for the entire population of physician reviewers (ie, population-level model). For a more detailed description of HLR models, see [Supplementary Material](#).

We developed HLR predictive models for each of the selected 73 targets. Each predictive model of a target variable is formulated as a binary classification problem where the model learns to identify cases in which the target variable is relevant. To investigate the util-

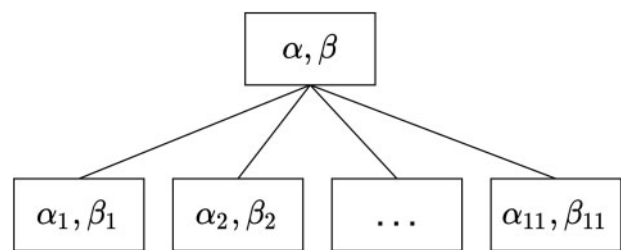


Figure 2. A 2-level HLR model for LEMR data. The lower level represents physician-level intercepts (α_i) and coefficients (β_i) where $i=1, \dots, 11$ denotes the physician identifier. The upper level represents the intercept and coefficients (α, β) for the population-level model.

ity of HLR over nonhierarchical models, we used LR as baseline models in which the physician identifier was included as an indicator variable. We implemented the HLR models using the *brms* package³⁴ in R, which uses No-U-Turn Sampler (NUTS) (as an extension of the Hamiltonian Monte Carlo algorithm) to estimate the posterior distribution of model parameters. In our experiments, we set

the NUTS sampler to use 4 Markov chains; each chain included 400 iterations of sampling where the first 200 were used to calibrate the sampler. A total of $4 \times 200 = 800$ posterior samples for each HLR model parameter were obtained. LR models were implemented using the *glmnet* package in R.³⁵

Cross validation

Each model was trained and evaluated independently in a stratified 10-fold cross-validation setting. At each iteration of the cross validation, the patient cases were randomly split into a training set (9 folds) and a test set (1 fold), while preserving the original distribution of the target variable. Hyperparameter tuning and data preprocessing such as imputing missing values and feature selection were performed during cross validation. More details are described in [Supplementary Materials](#).

Performance measures

We measured the predictive performance of each model with the area under the receiver operating characteristic (ROC) curve (AUROC), area under the precision-recall curve (AUPRC), and expected calibration error (ECE).³⁶ AUROC is a measure of model discrimination and varies from 0.5 and 1, where 0.5 denotes an uninformative model and 1 represents perfect discrimination.

Table 1. Years of ICU experience for each physician and the number of patient cases each physician reviewed

Physician identifier	Years of ICU experience	No. of cases reviewed (no. of ARF, no. of AKF)
1	<1	15 (8, 7)
2	1	15 (10, 5)
3	3	12 (5, 7)
4	<1	17 (8, 9)
5	1	15 (9, 6)
6	1	15 (7, 8)
7	2	22 (10, 12)
8	1	20 (11, 9)
9	1	16 (8, 8)
10	2	16 (8, 8)
11	7	15 (9, 6)

AUPRC summarizes the precision-recall curve where precision (or positive predictive value) and recall (or sensitivity) values at different thresholds are plotted as a curve. The AUPRC varies from 0 to 1 and is commonly used in binary classification problems when the data are imbalanced (ie, when cases with one label are more prevalent than cases with the other label).

ECE is a measure of model calibration. In a perfectly calibrated model, outcomes with predicted probability p correspond to a fraction p of positive cases in the data. ECE is derived from the probability calibration curve³⁷ where the sorted predicted probabilities are partitioned into k bins; in each bin i , calibration error is defined as the absolute difference between the mean of predicted probabilities (P_i) and the fraction of positive outcomes (o_i). ECE is the weighted average of the calibration errors over all k bins:

$$ECE = \sum_{i=1}^k w_i |p_i - o_i|, \quad (1)$$

where w_i denotes the fraction of cases that fall into bin i . Lower ECE denotes a better-calibrated model.

RESULTS

We report the variability across the physician reviewers and then report the results of the predictive performance of LR and HLR models from three perspectives: *overall*, *per-target*, and *per-physician*. [Table 1](#) summarizes the physician characteristics and the number of patients that each physician reviewed within the two diagnostic groups, AKF and ARF.

Table 2. Overall AUROC, AUPRC, and ECE for LR and HLR models over all 73 target variables and across all physicians

Measure	LR	HLR	P-value
AUROC	0.75 (0.74–0.76)	0.81 (0.80–0.82)	<.001
AUPRC	0.665 (0.663–0.667)	0.763 (0.762–0.765)	<.001
ECE	0.16 (0.14–0.17)	0.07 (0.06–0.08)	<.001

Higher AUROC and AUPRC show better discrimination power while lower ECE denotes better probability calibration. The best values for each metric are in boldface.

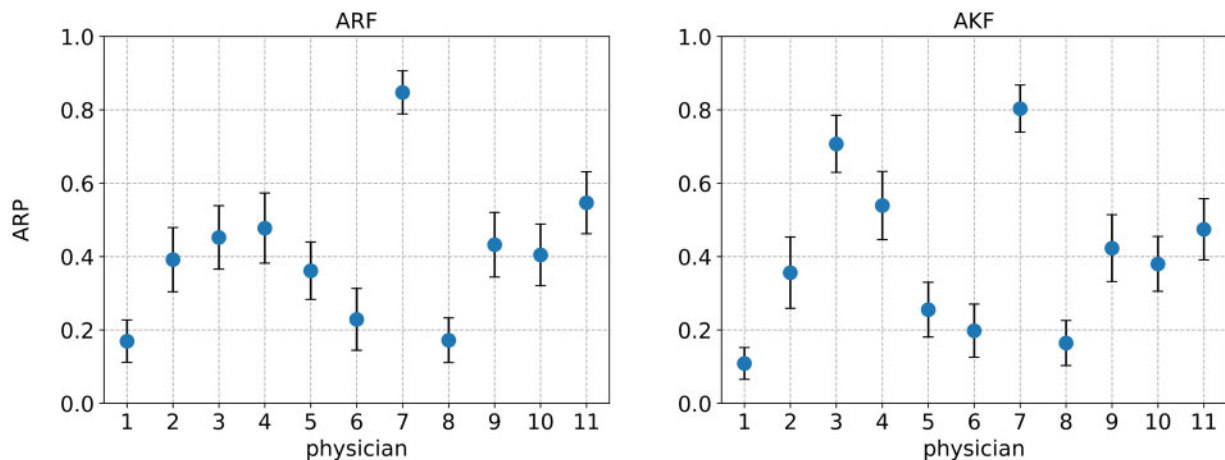


Figure 3. Per-physician ARP values over 73 target variables. A blue circle denotes the ARP value and an error bar denotes a 95% CI. The panel on the left is for ARF cases and the panel on the right is for AKF cases.

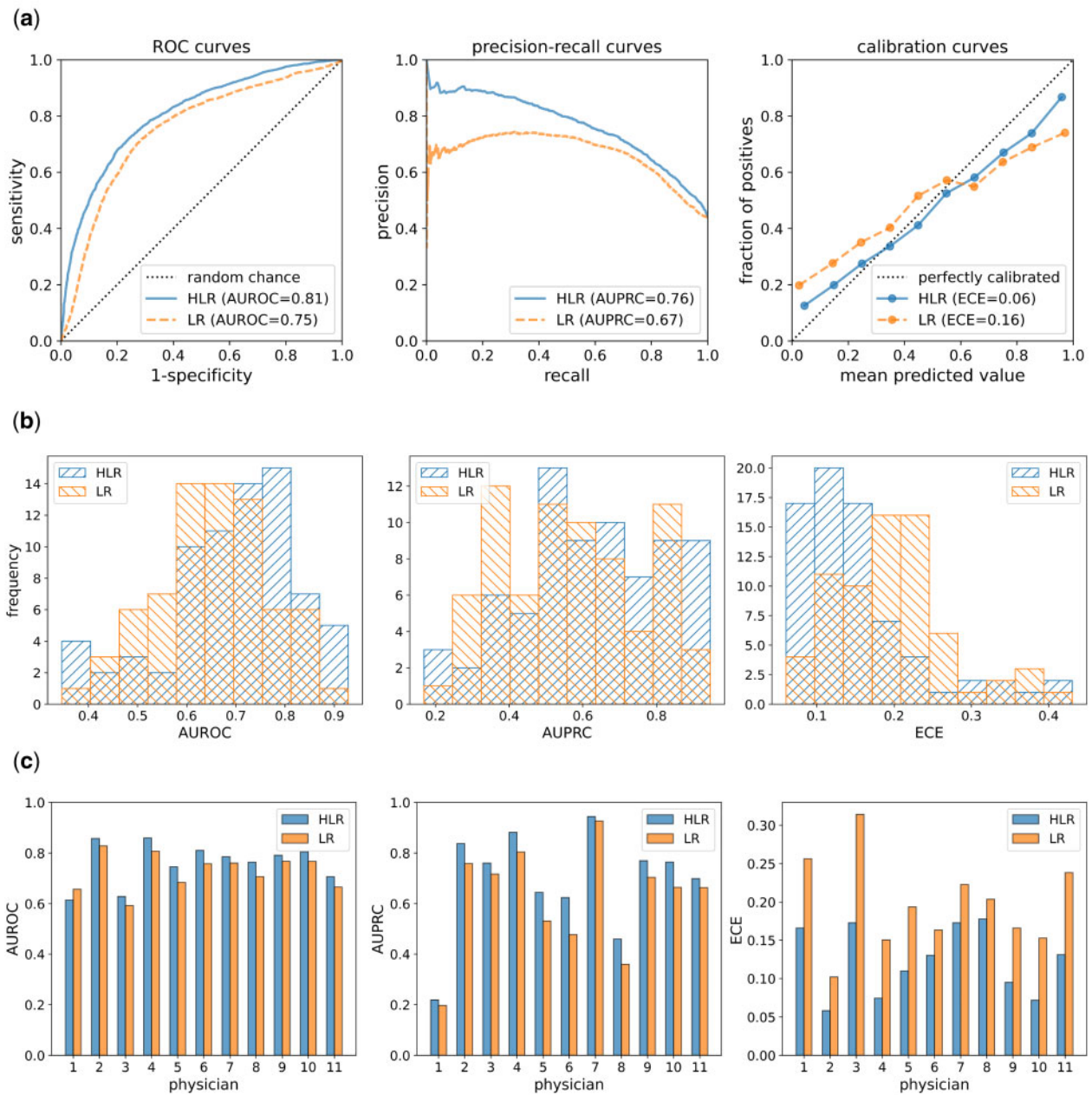


Figure 4. (A) ROC, precision-recall, and calibration curves over all 73 target variables across all physicians. For the calibration curves, the closer a curve is located to the dotted diagonal line, the more calibrated the corresponding approach is. (B) Distribution of AUROC, AUPRC, and ECE values for 73 models. Forward-slash hatches in blue represent the distributions for HLR models and backslash hatches in orange denote the distributions for LR models. The AUROC and AUPRC distributions for HLR models are right-skewed relative to the LR models, which shows that HLR models generally have better discrimination power. The distribution of ECE values of HLR models is left-skewed relative to the LR models, which means that HLR models are generally better calibrated than LR models. (C) AUROC, AUPRC, and ECE values for each physician reviewer over all 73 models. The values for HLR models are shown in blue and the values for LR models are shown in orange. The AUROC and AUPRC values are higher for HLR models than for LR models, except for the AUROC value for physician 1. All the ECE values are lower for HLR models, which means that HLR models are better calibrated than the LR models.

Variability in information-seeking behavior

We define a descriptive statistic called *average relevance proportion* (ARP) to measure the information-seeking behavior of each physician reviewer. An ARP value for a physician is defined as the average proportion of EMR data items that the physician sought as relevant. We calculated the ARP values over the 73 EMR data items that were used as target variables. Figure 3 shows the physician ARP values separately for each of the diagnostic groups. Each circle denotes

the ARP value for the corresponding physician on the x-axis and each error bar represents a 95% confidence interval (CI) for an ARP value. In the ARF diagnosis group, the ARP CIs for physicians 1, 7, and 8 do not overlap with those of the other physicians, which indicates a potential variability in information-seeking behavior between these physicians and the rest. Similar variability is observed in the AKF group, where the ARP CIs of physicians 1, 3, 7, and 8 differ from those of the other physicians.

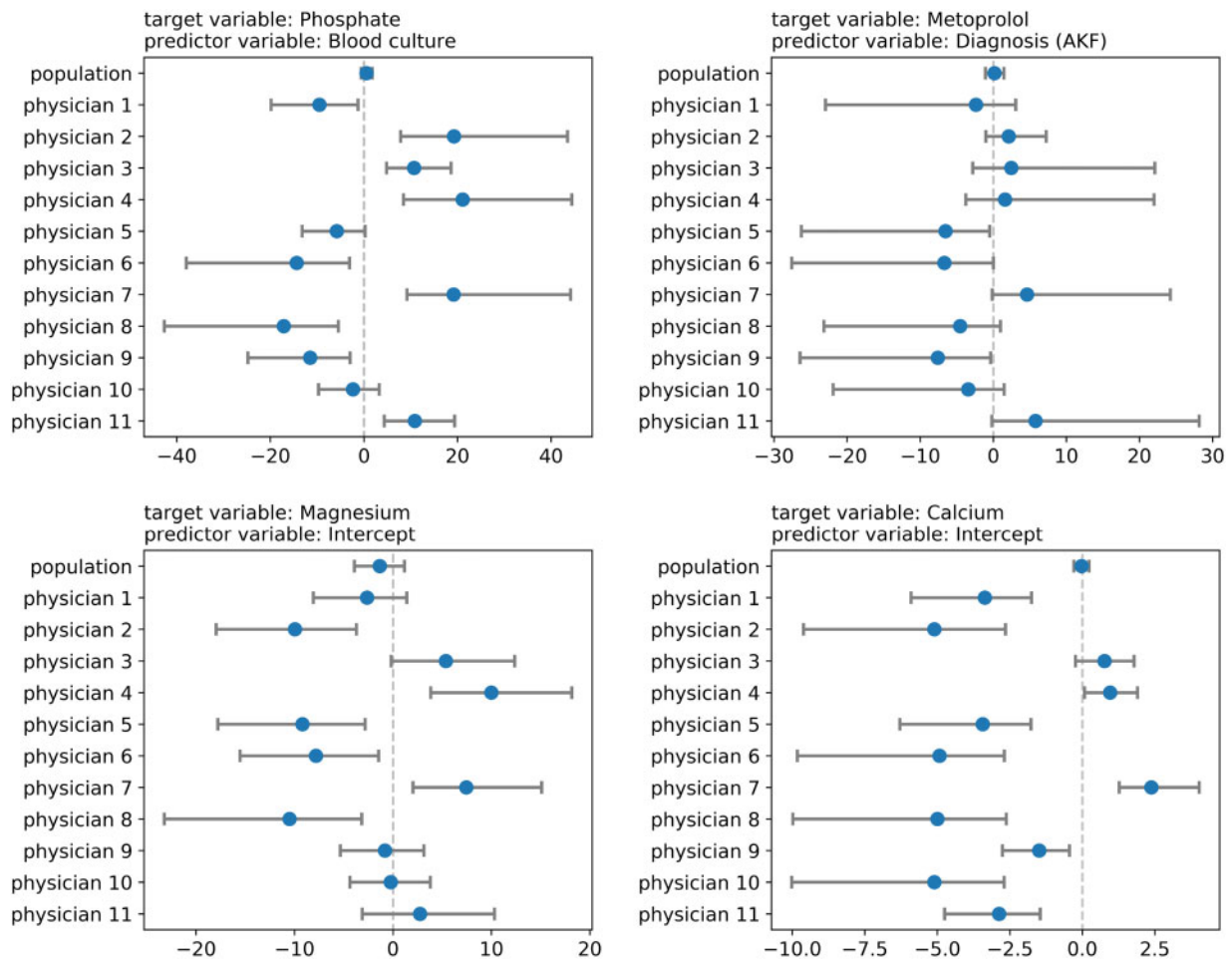


Figure 5. Examples of variation among physicians as seen from the values of the coefficients of a specific predictor variable. Each panel shows estimates of the coefficients of a predictor variable in an HLR model. A circle denotes the median value and the bar denotes the 80% credible interval for the posterior distribution of the model parameter.

Overall performance

The overall performance of each model family (LR and HLR) was calculated by concatenating the predictions for all 73 target variables into a single vector and using that vector to compute the performance metrics. Table 2 reports the AUROC, AUPRC, and ECE for the LR and HLR models across all 73 target variables. For AUROC values, the 95% CI and P -value were calculated using Delong's method.^{38,39} The 95% CI for AUPRC values was derived using the logit intervals method⁴⁰ and the P -value was calculated using the Wald z -test. For ECE values, we set $k = 100$ in Equation 1 and obtained a vector of 100 calibration errors to compute 95% CIs and a t -test P -value. Figure 4A shows the overall ROC and calibration curves for LR and HLR models. Note that for the calibration curves, we set the number of bins to $k = 10$ for better visibility.

Per-target performance

For per-target performance, we computed the predictive performance for each target variable, which resulted in vectors of AUROC, AUPRC, and ECE values each with a length of 73, for each model family (LR and HLR). Distributions of per-target performance measures are shown as histograms in Figure 4B for each model family. Histograms of the two model families are overlaid for better comparison. Additional details are provided in Supplementary

Table S1 where AUROC, AUPRC, and ECE values are reported for each target variable.

Per-physician performance

For per-physician performance, we computed the predictive performance for each physician, which resulted in 11 AUROC, AUPRC, and ECE values for each model family (LR and HLR). Figure 4C presents the per-physician bar plots of the performance measure values; the bars for HLR and LR models are displayed side by side for better comparison. Per-physician calibration curves are presented in Supplementary Figure S1.

DISCUSSION

Our results show that HLR models perform better than LR models when predicting which information a physician will seek in a future patient case. Moreover, the ECE results show that HLR models are generally better calibrated than LR models. In general, the more calibrated the probabilities are that are output by a predictive model, the higher the expected utility of the decisions that will be made using that model; in the case of the LEMR system, those decisions involve which information is worthwhile to highlight in the EMR of a given patient.

Although most physician reviewers had similar years of ICU experience, we observed a considerable degree of variability in information-seeking behavior across physicians in terms of ARP values. Because the study patients were selected to have a similar level of complexity, patient cases are unlikely to be the source of this variability. Controlling for physicians' years of experience in LR models was not as effective in improving predictive performance as estimating individual physician variability using the HLR models. This shows the advantage of HLR models over standard models in the presence of unexplained variability.

The per-physician performance measures in Figure 4C show that HLR models learn physician-specific models that perform better in terms of both discrimination and calibration. Although HLR models fit a separate model for each physician, the inherent regularization in these models prevents overfitting. In particular, as population and physician-specific parameters are estimated at the same time, a pooling effect occurs that prevents a physician-specific model from overfitting when the sample size is small.

Furthermore, HLR models allow for a detailed investigation at the physician level because each physician model has its own set of parameters. Figure 5 demonstrates a few instances of the detailed information that can be obtained from an HLR model. Each panel in Figure 5 represents the distributions of a model parameter in an HLR model for each physician and for all physicians as a whole. Investigating the physician-specific parameters can lead to a better understanding of factors that influence a physician's information-seeking behavior.

LIMITATIONS

One limitation of this study was the relatively modest amount of annotated data. Having experts review and annotate data is an expensive and time-consuming task in many domains, especially in medicine. It takes many hours for a physician to review and annotate a small number of patient cases in the EMR, which makes it challenging to collect large amounts of annotated data in the LEMR system. Due to this limitation, the number of positive samples for most target variables was modest. As a result, we derived models for only 73 target variables out of 865 available target variables. This challenge can be addressed by using scalable data collection methods. For example, a scalable solution based on eye-tracking technology has been proposed to automatically identify information that physicians seek in the EMR.⁴¹

Another limitation of this study is that we did not model the reliability of the annotations across physicians. In practice, physicians will not agree on which information is relevant due to differences in knowledge, level of expertise, and subjective preferences. In a previous study, we observed poor agreement among physicians on which information is relevant for the same case and clearly specified clinical task.¹⁶ Moreover, inexperienced physicians may provide annotations that are judged to be erroneous by expert physicians. By pooling annotations from all physicians, the HLR models are more influenced by physicians who have annotated a larger number of cases, and if disproportionately more cases are erroneously annotated in the training data, the models will provide poor recommendations. Thus, it is important to distinguish between correct and erroneous annotations which may be challenging in complex patient cases where even experts may disagree on which patient information is relevant. One approach to addressing this limitation is to develop a consensus model that physicians collectively converge to and combine such a model with physician-specific models.⁴²

Despite the advantage of HLR models in terms of performance, a drawback of HLR models is the added complexity due to the additional per-level parameters. This complexity creates new challenges in parameter estimation and interpretation. Compared to LR models, training HLR models requires more computing power and there are more hyperparameters to tune, including the choice of prior distributions.

CONCLUSION

Displaying large quantities of patient information in EMR systems with little prioritization can adversely influence the decision-making process of physicians and compromise the safety of patients. A data-driven solution was recently proposed as a learning EMR (LEMUR) system that uses machine learning to identify and prioritize relevant data in the EMR for physicians. The current study improves the performance of LR models by using HLR models.

We trained 2-level HLR models that simultaneously learn physician-specific models at one level and a population model at another level. We evaluated the discrimination and calibration performance of HLR models in identifying relevant data items in the EMR. Our results show that HLR models perform significantly better than LR models. Moreover, we demonstrated that HLR models provide details about the physician-specific models that can be used to investigate physicians' information-seeking behaviors in the EMR system.

SUPPLEMENTARY MATERIAL

Supplementary material is available at *Journal of the American Medical Informatics Association* online.

ACKNOWLEDGMENTS

This research was supported in part by the University of Pittsburgh Center for Research Computing (CRC) through the resources provided. We specifically acknowledge the assistance of the research faculty consultants at CRC. The study was approved by the University of Pittsburgh IRB under protocol PRO14020588.

FUNDING

The research reported in this publication was supported by the National Library of Medicine of the National Institutes of Health under award number R01 LM012095, and a Provost Fellowship in Intelligent Systems at the University of Pittsburgh (awarded to M.T.). The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

CONFLICT OF INTEREST STATEMENT

None declared.

REFERENCES

1. Nolan ME, Siwani R, Helmi H, *et al.* Health IT usability focus section: data use and navigation patterns among medical ICU clinicians during electronic chart review. *Appl Clin Inform* 2017; 08 (04): 1117–26.
2. Roman LC, Ancker JS, Johnson SB, *et al.* Navigation in the electronic health record: a review of the safety and usability literature. *J Biomed Inform* 2017; 67: 69–79.

3. Yang L, Mei Q, Zheng K, *et al.* Query log analysis of an electronic health record search engine. *AMIA Annu Symp Proc* 2011; 2011: 915–24.
4. Natarajan K, Stein D, Jain S, *et al.* An analysis of clinical queries in an electronic health record search utility. *Int J Med Inform* 2010; 79 (7): 515–22.
5. Ruppel H, Bhardwaj A, Manickam RN, *et al.* Assessment of electronic health record search patterns and practices by practitioners in a large integrated health care system. *JAMA Netw Open* 2020; 3 (3): e200512.
6. Mazur LM, Mosaly PR, Moore C, *et al.* Association of the usability of electronic health records with cognitive workload and performance levels among physicians. *JAMA Netw Open* 2019; 2 (4): e191709.
7. Hall A, Walton G. Information overload within the health care system: a literature review. *Health Info Libr J* 2004; 21 (2): 102–8.
8. Ahmed A, Chandra S, Herasevich V, *et al.* The effect of two different electronic health record user interfaces on intensive care provider task load, errors of cognition, and performance. *Crit Care Med* 2011; 39: 1626–34.
9. Pollack AH, Pratt W. Association of health record visualizations with physicians' cognitive load when prioritizing hospitalized patients. *JAMA Netw Open* 2020; 3 (1): e1919301.
10. Manor-Shulman O, Beyene J, Frndova H, *et al.* Quantifying the volume of documented clinical information in critical illness. *J Crit Care* 2008; 23 (2): 245–50.
11. Law AS, Freer Y, Hunter J, *et al.* A comparison of graphical and textual presentations of time series data to support medical decision making in the neonatal intensive care unit. *J Clin Monit Comput* 2005; 19 (3): 183–94.
12. Koch SH, Weir C, Westenskow D, *et al.* Evaluation of the effect of information integration in displays for ICU nurses on situation awareness and task completion time: a prospective randomized controlled study. *Int J Med Inform* 2013; 82 (8): 665–75.
13. Wright MC, Borbolla D, Waller RG, *et al.* Critical care information display approaches and design frameworks: a systematic review and meta-analysis. *J Biomed Informatics X* 2019; 3: 100041.
14. Pickering BW, Dong Y, Ahmed A, *et al.* The implementation of clinician designed, human-centered electronic medical record viewer in the intensive care unit: a pilot step-wedge cluster randomized trial. *Int J Med Inform* 2015; 84 (5): 299–307.
15. King AJ, Cooper GF, Clermont G, *et al.* Using machine learning to selectively highlight patient information. *J Biomed Inform* 2019; 100: 103327.
16. King AJ, Cooper GF, Hochheiser H, *et al.* Using machine learning to predict the information seeking behavior of clinicians using an electronic medical record system. *AMIA Annu Symp Proc* 2018; 2018: 673–82.
17. Gelman A, Hill J. *Data Analysis Using Regression and Multilevel/Hierarchical Models*. Cambridge, England: Cambridge University Press; 2006.
18. Burns LR, Wholey DR. The effects of patient, hospital, and physician characteristics on length of stay and mortality. *Med Care* 1991; 29 (3): 251–71.
19. Wilkinson DJC, Truog RD. The luck of the draw: physician-related variability in end-of-life decision-making in intensive care. *Intensive Care Med* 2013; 39 (6): 1128–32.
20. Yadav KN, Josephs M, Gabler NB, *et al.* What's behind the white coat: potential mechanisms of physician-attributable variation in critical care. *PLoS One* 2019; 14 (5): e0216418.
21. Garland A, Shaman Z, Baron J, *et al.* Physician-attributable differences in intensive care unit costs: a single-center study. *Am J Respir Crit Care Med* 2006; 174 (11): 1206–10.
22. Guterman JJ, Lundberg SR, Scheib GP, *et al.* Wide variability in emergency physician admission rates: a target to reduce costs without compromising quality. *West J Emerg Med* 2016; 17 (5): 561–6.
23. Obermeyer Z, Powers BW, Makar M, *et al.* Physician characteristics strongly predict patient enrollment in hospice. *Health Aff* 2015; 34 (6): 993–1000.
24. Pollack CE, Radhakrishnan A, Parker AM, *et al.* Are physicians social networks linked to breast cancer screening recommendations for older adults? *J Clin Oncol* 2017; 35 (15_suppl): 6550.
25. DeMott RK, Sandmire HF. The Green Bay cesarean section study. *Am J Obstet Gynecol* 1990; 162 (6): 1593–602.
26. Wang Y, Jiang Y, Luo S, *et al.* Multilevel model estimation of age-dependent individual-specific trajectories for left ventricular echocardiographic indexes in an asymptomatic elderly cohort. *J Appl Physiol* 2012; 112 (7): 1128–34.
27. Chung H, Ng E, Ibrahim S, *et al.* Welfare state regimes, gender, and depression: a multilevel analysis of middle and high income countries. *Int J Environ Res Public Health* 2013; 10 (4): 1324–41.
28. Pan H, Gao S, Grant K, *et al.* Analyzing national and state opioid abuse treatment completion with multilevel modeling. In: *Systems and Information Engineering Design Symposium (SIEDS)*; 2018: 123–8; Charlottesville, VA.
29. Berta P, Martini G, Moscone F, *et al.* The association between asymmetric information, hospital competition and quality of healthcare: evidence from Italy. *J R Stat Soc A* 2016; 179 (4): 907–26.
30. Berta P, Vinciotti V. Multilevel logistic cluster-weighted model for outcome evaluation in health care. *Stat Anal Data Min ASA Data Sci J* 2019; 12 (5): 434–43.
31. Towne SD, Fair K, Smith ML, *et al.* Multilevel comparisons of hospital discharge among older adults with a fall-related hospitalization. *Health Serv Res* 2018; 53 (4): 2227–48.
32. Lin X, Cai M, Tao H, *et al.* Insurance status, in-hospital mortality and length of stay in hospitalised patients in Shanxi, China: a cross-sectional study. *BMJ Open* 2017; 7 (7): e015884.
33. King A. *The development and evaluation of a learning electronic medical record system*. [Doctoral dissertation]. Pittsburgh: University of Pittsburgh; 2018.
34. Bürkner PC. brms: an R package for Bayesian multilevel models using Stan. *J Stat Softw* 2017; 80: 1–28.
35. Friedman J, Hastie T, Tibshirani R. Regularization paths for generalized linear models via coordinate descent. *J Stat Softw* 2010; 33 (1): 1–22.
36. Naeini MP, Cooper GF, Hauskrecht M. Obtaining well calibrated probabilities using Bayesian binning. *AAAI Conference on Artificial Intelligence*. 2015: 2901–7.
37. DeGroot MH, Fienberg SE. The comparison and evaluation of forecasters. *J R Stat Soc Ser D (the Stat)* 1983; 32: 12–22.
38. DeLong ER, DeLong DM, Clarke-Pearson DL. Comparing the areas under two or more correlated receiver operating characteristic curves: a non-parametric approach. *Biometrics* 1988; 44 (3): 837–45.
39. Robin X, Turck N, Hainard A, *et al.* pROC: an open-source package for R and S+ to analyze and compare ROC curves. *BMC Bioinformatics* 2011; 12 (1): 77.
40. Boyd K, Eng KH, Page CD. Area under the precision-recall curve: point estimates and confidence intervals In: *Lecture Notes in Computer Science (Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*. Berlin, Heidelberg: Springer; 2013: 451–66.
41. King AJ, Hochheiser H, Visweswaran S, *et al.* Eye-tracking for clinical decision support: a method to capture automatically what physicians are viewing in the EMR. *AMIA Jt Summits Transl Sci Proc* 2017; 2017: 512–521.
42. Valizadegan H, Nguyen Q, Hauskrecht M. Learning classification models from multiple experts. *J Biomed Inform* 2013; 46 (6): 1125–35.