

Generalizability Challenges of Mortality Risk Prediction Models

Generalizability Challenges of Mortality Risk Prediction Models: A Retrospective Analysis
on a Multi-center Database

Short Title: Generalizability Challenges of Mortality Risk Prediction Models

Harvineet Singh^{1*}, Vishwali Mhasawade², Rumi Chunara^{2,3}

¹ New York University, Center for Data Science

² New York University, Tandon School of Engineering

³ New York University, School of Global Public Health

*Corresponding author

E-mail: hs3673@nyu.edu (HS)

Abstract

Modern predictive models require large amounts of data for training and evaluation, absence of which may result in models that are specific to certain locations, populations in them and clinical practices. Yet, best practices for clinical risk prediction models have not yet considered such challenges to generalizability. Here we ask whether population- and group-level performance of mortality prediction models vary significantly when applied to hospitals or geographies different from the ones in which they are developed. Further, what characteristics of the datasets explain the performance variation? In this multi-center cross-sectional study, we analyzed electronic health records from 179 hospitals across the US with 70,126 hospitalizations from 2014 to 2015. Generalization gap, defined as difference between model performance metrics across hospitals, is computed for area under the receiver operating characteristic curve (AUC) and calibration slope. To assess model performance by the race variable, we report differences in false negative rates across groups. Data were also analyzed using a causal discovery algorithm “Fast Causal Inference” that infers paths of causal influence while identifying potential influences associated with unmeasured variables. When transferring models across hospitals, AUC at the test hospital ranged from 0.777 to 0.832 (1st-3rd quartile or IQR; median 0.801); calibration slope from 0.725 to 0.983 (IQR; median 0.853); and disparity in false negative rates from 0.046 to 0.168 (IQR; median 0.092). Distribution of all variable types (demography, vitals, and labs) differed significantly across hospitals and regions. The race variable also mediated differences in the relationship between clinical variables and mortality, by hospital/region. In conclusion, group-level performance should be assessed during generalizability checks to identify potential harms to the groups. Moreover, for developing methods to improve model performance in new environments, a better understanding and documentation of provenance of data and health processes are needed to identify and mitigate sources of variation.

Author Summary

With the growing use of predictive models in clinical care, it is imperative to assess failure modes of predictive models across regions and different populations. In this retrospective cross-sectional study based on a multi-center critical care database, we find that mortality risk prediction models developed in one hospital or geographic region exhibited lack of generalizability to different hospitals or regions. Moreover, distribution of clinical (vitals, labs and surgery) variables significantly varied across hospitals and regions. Based on a causal discovery analysis, we postulate that lack of generalizability results from dataset shifts in race and clinical variables across hospitals or regions. Further, we find that the race variable commonly mediated changes in clinical variable shifts. Findings demonstrate evidence that predictive models can exhibit disparities in performance across racial groups even while performing well in terms of average population-wide metrics. Therefore, assessment of subgroup-level performance should be recommended as part of model evaluation guidelines. Beyond algorithmic fairness metrics, an understanding of data generating processes for subgroups is needed to identify and mitigate sources of variation, and to decide whether to use a risk prediction model in new environments.

Introduction

Validation of predictive models on intended populations is a critical prerequisite to their application in making individual-level care decisions since a miscalibrated or inaccurate model may lead to patient harm or waste limited care resources (1). Models can be validated either on the same population as used in the development cohort, named internal validity, or on a different yet related population, named external validity or generalizability (or sometimes transportability) (2). The TRIPOD (Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis) Statement strongly recommends

assessing external validity of published predictive models in multiple ways including testing on data from a different geography, demography, time period, or practice setting (3). However, the guidelines do not specify appropriate external validity parameters on any of the above factors. At the same time, recent studies in the computer science and biomedical informatics literature have indicated that sub-group performance of clinical risk prediction models by race or sex can vary dramatically (4,5), and clinical behavior can guide predictive performance (6). Within the statistics literature are several methods for computing minimum sample size and other best practices for assessing external validity of clinical risk prediction models (7–9). However, the assessment of sub-group-level performance and data-shifts are not explicitly considered in such guidance (10). Moreover, recent analyses have shown that clinical prediction models are largely being developed in a limited set of geographies, bringing significant concern regarding generalizability of models to broader patient populations (11). Amidst such rising challenges, an understanding of how differences among population and clinical data impact external generalizability of clinical risk prediction models is imperative.

When a model fails to generalize for specific patient groups (such as racial or gender identities), using it to guide clinical decisions can lead to disparate impact on such groups. This raises questions of equity and fairness in the use of clinical risk prediction models, for which performance on diverse groups has been repeatedly lacking (12–14). Predictive discrimination quantifies how well a model can separate individuals with and without the outcome of interest (we study mortality prediction here). Calibration quantifies how well the predicted probabilities match with the observed outcomes. These measures can be used to check for aggregate model performance across a study sample or within groups, but do not illuminate variation across groups. Hence, in assessment of generalizability, we add another set of measures to our analyses which we refer to as “fairness” metrics, following the

algorithmic fairness literature (15–17). Such performance checks are important, especially given the evidence on racial bias in medical decision-making tools (4,13,14).

The primary objective of this study is to evaluate the external validity of predictive models for clinical decision making across hospitals and geographies in terms of the metrics – predictive discrimination (area under the receiver operating characteristic curve), calibration (calibration slope) (18), and algorithmic fairness (disparity in false negative rates and disparity in calibration slopes). The secondary objective is to examine the possible reasons for performance changes via shifts in the distributions of different types of variables and their interactions. We focus on risk prediction models for in-hospital mortality in ICUs. Our choice of evaluation metrics are guided by the use of such models for making patient-level care decisions. We note that similar models (e.g., SAPS and APACHE scores) (19,20) are widely-used for other applications as well such as assessing quality-of-care, resource utilization, or risk-adjustment for estimating healthcare costs (19–21), which are not the focus of this study. Recently, prediction models for in-hospital mortality have been prospectively validated for potential use (22), or in case of sepsis, have even been integrated into the clinical workflow (23). With access to large datasets through electronic health records, new risk prediction models leveraging machine learning approaches have been proposed, which provide considerable accuracy gains (24). Being flexible, such approaches might overfit to the patterns in a particular dataset, thus, raising concerns for their generalization to newer environments (25). We use the eICU dataset (26) as a test bed for our analyses. Past studies have employed the dataset for evaluating mortality prediction models (27,28). As the dataset was collected from multiple hospitals across the US, it allows us, in a limited way, to test external validity across hospitals, diverse geographies, and populations.

Materials and Methods

Analyses are based on data obtained from the publicly-available eICU Collaborative Research Database (26), designed to aid remote care of critically-ill patients in a telehealth ICU program. The database is composed of a stratified random sample of ICU stays from hospitals in the telehealth program where the sample is selected such that the distribution of the number of unique patient-stays across hospitals is maintained (26). Data on 200,859 distinct ICU stays of 139,367 patients with multiple visits across 208 hospitals in the US between 2014 and 2015 are included. We followed the Strengthening the Reporting of Observational Studies in Epidemiology (STROBE) reporting guideline (29).

Data Preprocessing

We follow the feature extraction and exclusion procedures including the exclusion criteria used by Johnson et al (27) which removes patient stays conforming to APACHE IV exclusion criteria (20) and removes all non-ICU stays. APACHE IV criteria excludes patients admitted for burns, in-hospital readmissions, patients without a recorded diagnosis after 24 hours of ICU admission, and some transplant patients (26). Only patients 16 years or above are included. Age of patients above 89 years (which is obfuscated to adhere to HIPAA provisions) is coded as 90. After pre-processing, the dataset consists of 70,126 stays from 179 hospitals. For analyses, data is grouped at two levels – by individual hospitals and by U.S. geographic regions (Northeast, South, Midwest, West) (30). Hospital-level analyses are restricted to the top 10 hospitals with the most stays, all of which have at least 1631 stays, to ensure enough examples for model training and evaluation. Data is split into ten separate datasets using a hospital identifier for hospital-specific analyses and into four separate datasets using a region identifier for region-specific analyses. The outcome label is in-hospital

mortality (binary). Mortality rates differed in the range of 3.9%-9.3% (1st-3rd quartile) across hospitals. Summary statistics by hospital and region are included in Table S1 and Table S2.

Mortality Prediction Model

Features from the SAPS II risk scoring model (19) from the first 24 hours of the patient-stay starting from ICU admission were extracted and are summarized in Table S3. These include 12 physiological measurements (vitals and labs), age, and an indicator for whether the stay was for an elective surgery. For features with multiple measurements, their worst values determined using the SAPS II scoring sheet (Table 3 in Le Gall et al (19)) are extracted. For example, for Glasgow Coma Score, we take the minimum value among the measurements. As previously employed for mortality prediction (27), we use logistic regression with ℓ_2 regularization using the implementation in scikit-learn v0.22.2 package with default hyperparameters (31). Missing values in features are imputed with mean values computed across the corresponding columns of the full dataset. We experimented with other imputation methods as well, specifically imputation with mean or median across the train datasets and single imputation with a decision tree (77), however, the conclusions did not change. Features are then standardized to zero mean and unit variance using statistics from the train datasets. As sample size used in training models can affect generalizability, we control for this factor by fixing the number of samples used for training and testing. We use 1631 (or 5000) samples from each hospital (or region) while training and testing models. For model development, each dataset (for a hospital or region) is randomly split with the training set comprising 90% of samples and the validation set comprising the remaining 10%. The test set comprises all samples from the hospital (or region) different from the one included in the training set.

Statistical Analysis

Performance Metrics: Discrimination ability of the models is assessed using area under the receiver operating characteristic curve (AUC). For binary outcomes, calibration slope (CS) is computed as the slope of the regression fit between true outcomes and logits of the predicted mortality, with a logit link function. A perfectly calibrated model has a CS of 1. A value lower than 1 indicates that the risk estimates are extreme, i.e. overestimation for high risk patients and underestimation for low risk patients, and suggests overfitting of the model (32). Hence, a value close to 1 is desirable. In addition to reporting AUC and CS computed on the test sets, we also report how much the metrics differ from their values computed on the validation set. This difference, known as the generalization gap, provides a quantitative measure of generalization performance (e.g. Jiang et al (33)). If this difference is high, i.e. the test metrics are worse than the train metrics, the model is said to lack generalizability. The allowable difference between test and train depends on the application context. Studies typically report confidence intervals around the difference and/or the percentage change relative to the train performance (34). To measure fairness of model predictions, we use the racial/ethnic attributes to form two groups – African American, Hispanic, and Asian as one group and the rest as another. These will be referred to as *minority* and *majority* groups. Note that the racial/ethnic attributes are used only for the purposes of fairness analysis, these are not part of the model building process. We acknowledge that aggregating multiple groups does not represent an exposition of which groups are advantaged in the models. Based on the available dataset sizes, this approach serves to illustrate differences that would ideally be unpacked in detail in future work addressing such issues. Disparity in false negative rates (DisparityFNR), and disparity in calibration slope (DisparityCS) are computed as the difference between the respective metric's value for the minority and the majority group. Differences in these two metrics have been employed in recent studies for bias analysis (16,17). FNR quantifies the rate at which patients with the observed outcome of death were

misclassified. Thus, a high FNR for the score may lead to an increase in undertreatment, and high DisparityFNR (in absolute value) highlights large differences in such undertreatment across groups. For the prediction threshold for FNR we use the mortality rate at the *test* hospital (assuming it is known beforehand). This threshold can be chosen in a more principled way, for example, based on decision-curve analysis (18), which will depend on the application context. We further acknowledge that there are myriad ways to define fairness that will depend on the context of the risk prediction's use and inputs from stakeholders (35).

Dataset Differences: To address our secondary objective of studying external validity-specific performance changes, we test for dataset shifts across hospitals and geographies (i.e. whether the distributions of two datasets differ), and use causal graph discovery to explore the reasons for these differences. Dataset shifts are measured using squared maximum mean discrepancy (MMD^2) (36). We perform the two-sample tests under the null hypothesis that the distributions are the same and threshold the resulting p-values at the significance level of 0.05. Details of the MMD^2 metric and the hypothesis test are included in Method S1. To explain the shifts we leverage the recently introduced framework of Joint Causal Inference (37) which allows constructing a single graphical representation of how variables relate to each other, in the form of a causal graph. We use the Fast Causal Inference (FCI) algorithm (38) for constructing the causal graph as it is methodologically well-developed and requires fewer assumptions on the data generating process as it allows for the presence of unobserved variables affecting the observed variables in the data (i.e. unobserved confounders). We also include race as an indicator in the datasets, before running FCI, to study the causes of unfairness with respect to the race variable. Details of the causal modeling are included in Method S2. Note that we use causal graphs as a compact representation of (conditional) independencies in the datasets. These are not meant to make statements about the causal

effect of treatments on physiological variables, for which randomized controlled trials and other methods may be used to gather better evidence.

The metrics of interest – AUC, CS, DisparityFNR, DisparityCS, p-value, and MMD^2 – are averaged over 100 random subsamples of the datasets (i.e. resampling without replacement). While aggregating across hospitals (or regions), we report the median, 1st, and 3rd quartiles across all train-test set pairs which includes the 100 random subsamples in each pair.

[Figure 1 here]

Results

Figure 1 demonstrates the highly varied external validity of models across hospitals based on AUC, CS, generalization gap in AUC and in CS, DisparityFNR, and DisparityCS. Across all train-test hospital pairs, median AUC is 0.801 with 1st-3rd quartile range (IQR) as 0.778 to 0.832, and CS is 0.853 (IQR 0.725 to 0.983). AUCs are lower than the typical values for mortality risk prediction models of around 0.86 (19,22), although AUCs in the same range (around 0.8) have been observed in other studies (albeit in different populations) and were considered acceptable (39–41). CS of around 0.8, as observed in our case, is considered to indicate overfitting (7). Transferring a model trained on hospital ID 73, which is the hospital with the most samples, to other hospitals results in a median gap in AUC of -0.087 (IQR -0.134 to -0.046) and a median gap in CS of -0.312 (IQR -0.502 to -0.128). In aggregate, we observe a decline in the performance on the test hospitals relative to that on the train hospitals (Figure 1a). Across all train-test hospital pairs, the median generalization gap in AUC is -0.018 (IQR -0.065 to 0.032) and the median generalization gap in CS is -0.074 (IQR -0.279 to 0.121). Figure 1b shows that the majority of models have CS of less than 1, indicating consistent miscalibration of mortality risk at test hospitals. This conforms with the typical

observation of good discriminative power but poor calibration of SAPS II models (39,41–43). The median values of AUC and CS are negative, indicating that both of them decrease in majority of the cases upon transfer. For comparison, the generalization gap in AUC for the SAPS II score in the original study by Le Gall et al (19) was -0.02 (AUC decreased to 0.86 in validation from 0.88 in training data), which is the same as the median gap here. Thus, for more than half of the hospital pairs the AUC drop is worse than the acceptable amount found in the original SAPS II study. Percentage changes in AUC and CS from train to test set, reported in Table S4 also indicate substantial drop in performance (in the range of -2.5% to -31.5% in AUC and -15.9% to -45.4% in CS). In some cases, for example for hospital ID 252, we observe an improvement in AUC (fourth row from bottom, Figure 1d). With regard to fairness metrics, DisparityFNR (absolute value) has median 0.093 (IQR 0.046 to 0.168), i.e. false negative rates across the racial groups differ by 4.6% to 16.8%. DisparityCS, i.e. the absolute value of difference in calibration across racial groups, is large as well (median 0.159; IQR 0.076 to 0.293). Considering that the ideal DisparityCS is 0, when CS is 1 for both the groups, the observed DisparityCS of 0.159 is large. There are both positive and negative values in the disparity metrics (Figures 1c, 1f), i.e. models are unfair to the minority groups for some pairs and vice versa for others. Note that disparity metrics for hospital ID 338 are considerably different from others (third column in Figures 1c, 1f) due to the skewed race distribution with only 74 (3.2%) samples from the minority groups (Table S1). We observe that the variation across hospitals in fairness metrics (DisparityFNR, and DisparityCS) is not captured by the variation in discrimination and calibration metrics (AUC and CS). Thus, fairness properties of the models are not elucidated by the standard metrics and should be audited separately.

[Figure 2 here]

Next, given concerns about the development of machine learning models in a limited set of geographies (11), we pool hospitals by geographic region, and validate models trained in one

region and tested on another (Figure 2). Performance in terms of AUC and CS across regions improves as a result of pooling hospital data. Overall, AUC varies in a small range (median 0.804; IQR 0.795 to 0.813) as does CS (median 0.968; IQR 0.904 to 1.018). The same can be observed through generalization gaps in AUC and CS which are smaller – median generalization gap in AUC is -0.001 (IQR -0.017 to 0.016) and median generalization gap in CS is -0.008 (IQR -0.081 to 0.075). However, such pooling does not alleviate fairness metric disparities. DisparityFNR (absolute value) has a median value of 0.040 (IQR 0.018 to 0.074). This translates to, for example, a disparity between minority and majority groups of 6.36% (95% CI -7.66% to 17.42%) and 8.06% (95% CI -3.81% to 19.74%) when transferring models from Midwest to West and Northeast to West respectively (though CIs are large, both values are greater than 0%; one-sided one-sample t-test, $p=10^{-5}$). DisparityCS (absolute value) is still high with a median value of 0.104 (IQR 0.050 to 0.167). For example, transferring models from South to Northeast (the region with the least minority population size) has a high DisparityCS (median 0.108; IQR -0.015 to 0.216). Percentage change in the test set metrics relative to the train set is reported in Table S5 which shows significant changes in DisparityFNR (ranging from -33% to 65%) and DisparityCS (ranging from -52% to 67%). Apart from geography, differences across hospitals can also be due to differences in patient load and available resources. In Figure S2 we include results for more fine-grained pooling of hospitals based on their number of beds, teaching status, and region where we again find consistent lack of generalizability in fairness metrics.

To investigate reasons for these performance differences across hospitals and regions, we first consider whether the corresponding datasets differ systematically. Figure 3 shows results from statistical tests for dataset shifts across hospitals and regions. Shifts across all pairs of hospitals are significant. Some hospitals are considerably different from others like hospital ID 73 in the first column of Figure 3d, which has significantly lower mortality rate than the other hospitals (Table S1).

[Figure 3 and 4 here]

Finally, we study explanations for the observed shifts in Figure 3a via individual features. The discovered causal graph is included in Figure S3 which shows the estimated causal relationships among clinical variables, and which of these variables shift in distribution based on hospital, geography and other factors (i.e. which variables have a direct arrow from the indicators like hospital or region). Figure 4 summarizes the shifts from the causal graph. From Figure 4, we note that the distribution of all fourteen features and the outcome are affected either directly or indirectly by the hospital indicator. However, restricting to direct effects of hospital (first row in Figure 4), we observe that shifts are explained by few of the features – demography (age and race), vitals (3 out of 6), and labs (3 out of 6). Not all vitals and labs change directly as a result of a change in hospitals; changes in 3 of the vitals and 3 of the labs are mediated through changes in other features. We attempt to further understand these changes across hospitals by including hospital-level contextual information, namely, their region, size (number of beds), and teaching status. We observe that there exists common features that explain shifts across the three attributes (different rows in Figure 4). Thus, some of the variation across hospitals is explained through its region, size, and teaching status. But, notably, the three attributes do not explain all variation among hospitals and more contextual information is required. For the fairness analysis, we observe the direct effects of the race variable (last row in Figure 4). There are direct effects from the race variable to most vitals (4 out of 6), labs (4 out of 6), and indicator for elective surgery. These direct effects support past observations made on racial disparities, e.g. in access to specialized care (race → elective surgery) (44) and in blood pressure measurements (race → sysbp or systolic blood pressure) (45). Out of the 9 features that are directly affected by the race variable, 4 also vary across the hospitals. This suggests that the distribution of clinical variables for the racial groups differs as we go from one hospital to another. As a result, the feature-outcome relationships learnt in

one hospital will not be suitable in another, leading to the observed differences in model performance (as quantified by the fairness metrics) upon model transfer.

Discussion

We retrospectively evaluate generalizability of mortality risk prediction models using data from 179 hospitals in the eICU dataset. In addition to commonly-used metrics for predictive accuracy and calibration, we assessed generalization in terms of algorithmic fairness metrics. To interpret results, we investigated shifts in the distribution of variables of different types across hospitals and geographic regions, leading to changes in the metrics. Findings highlight that recommended measures for checking generalizability are needed, and evaluation guidelines should explicitly call for the assessment of model performance by sub-groups.

Generalizability of a risk prediction model is an important criteria to establish reliable and safe use of the model even under care settings different from the development cohort. A number of studies have reported lack of generalizability of risk prediction models across settings such as different countries (46–49), hospitals (27,50–53), or time periods (54–56). For instance, Austin et al (57) study the validity of mortality risk prediction models across geographies in terms of discrimination and calibration measures. They find moderate generalizability, however, the hospitals considered belonged to a single province in Canada. More importantly, these works did not investigate generalizability for different groups in the patient cohorts. Results in Figure 1 show that the fairness characteristics of the models can vary substantially across hospitals. Prior work investigating algorithmic fairness metrics in a clinical readmission task (4) did not investigate changes in the metric when models are transferred across care settings. A recent study of a mortality prediction model showed good performance across 3 hospitals (1 academic and 2 community-based) as well as good

performance for subgroups *within* a hospital (22). However, the change in performance for the subgroups *across* hospitals was not explored.

One strategy to tackle the lack of generalizability is to pool multiple hospital databases to potentially increase diversity of the data used in modeling (58). However, available databases may not faithfully represent the intended populations for the models even after pooling. A recent study (11) found that US-based patient cohorts used to train machine learning models for image diagnosis were concentrated in only three states. However, the effect of using geographically-similar data on generalizability and fairness had not been studied previously. Results in Figure 2 suggest that pooling data from similar geographies may not help mitigate differences in model performance when transferred to other geographies. This finding adds more weight to the concerns raised about possible performance drop when transferring models from data-rich settings to low-resource settings (59).

Prior work has postulated multiple reasons for lack of generalizability (76) including population differences and ICU admission policy changes (46). In Davis et al (54), reasons including case mix, event rate, and outcome-feature association are assessed. However, specific variables which shift across clinical datasets have not been examined. Through an analysis of the underlying causal graph summarized in Figure 4, we identify specific features that explain the changes in data distributions across hospitals. Demographics (age and race) differ across hospitals which is aligned with population difference being the common reason cited for lack of generalizability (46,54). We find that vitals and labs differ as well but only a few change as a direct consequence of changes in the hospital setting. Often, the changes are mediated via a small number of specific vitals and labs. For understanding the root causes of the shifts, causal graph analysis helps to narrow down candidate features to analyse further. Significant differences found across ICUs in feature distributions and model performance calls for a systematic approach to transferring models. Understanding the reasons for the lack of generalizability is the first step to deciding whether to transfer a model, re-training it for

better transfer, or developing methods which can improve transfer of models across environments (60). Factors affecting generalizability can potentially be due to variations in care practices or may represent spurious correlations learned by the model. These findings reinforce calls to better catalog clinical measurement practices (61) and continuously monitor models for possible generalizability challenges (76).

Beyond the descriptive analyses of dataset shifts with causal graphs, in Figure S1, we examine whether the shifts can *predict* lack of generalizability. We plot the generalization gap in AUC and CS against the amount of shift, as measured in MMD^2 , and report the correlation coefficient. Although we find only low correlation, this suggests a need to develop better methods of quantifying dataset shifts that can predict future model performance. One such metric derived from a model trained to discriminate between training and test samples was found to explain the test performance well (focusing on environments that primarily differ by case-mix, without attention to specific clinical or demographic shifts) (62). We hope that current work motivates development and evaluation of such metrics on larger and more diverse populations and datasets. Recent work has also proposed methodological advances to ensure transferability of models across settings, for example, by pre-training on large datasets from related machine learning tasks (63,64) and with the help of causal knowledge about shifts (65–67). Better metrics for dataset shift can help practitioners decide whether to transfer a model to a new setting based on how large the shift is between hospitals, for example.

Importantly, findings here showed that the race variable often mediated shifts in clinical variables. Reasons for this must be disentangled. As race is often a proxy variable for structural social processes such as racism which can manifest both through different health risk factors as well as different care (differences in health care received by patients' racial group are well-documented) (68–71), shifts across hospitals cannot be mitigated simply by population stratification or algorithmic fairness metrics alone. Indeed, better provenance of the process by which data is generated will be critical in order to disentangle the source of

dataset differences (for example, if clinical practices or environmental and social factors are giving rise to different healthcare measures and outcomes). Following guidelines developed for documenting datasets (72) and models (73) in the machine learning community, similar guidelines should be established for models in healthcare as well (10). An example is the proposal for reporting subgroup-level performances in MI-CLAIM checklist (74). In sum, our findings demonstrate that data provenance, as described above, is needed in addition to applying algorithmic fairness metrics alone, to understand the source of differences in healthcare metrics and outcomes (e.g. clinical practice versus other health determinants), and assess potential generalizability of models.

Limitations

The main limitation of this study is that the results are reported for a collection of ICUs within the same electronic ICU program by a single provider. This collection captures only a part of the diversity in care environments that a mortality prediction model might be deployed in. Further, our investigation is limited to models constructed using the SAPS II feature set containing 14 hand-crafted features for the mortality prediction task. Though we employ widely-used methods, our analysis is limited by the specific methods used for computing dataset differences, building predictive models and causal graphs. For analyzing reasons for dataset shifts, we could only investigate explanations based on the 14 features along with limited hospital characteristics (such as geographic region, number of beds, and teaching status). Multiple factors are left unrecorded in the eICU database, such as patient load, budget constraints, and socioeconomic environment of the hospital's target population, that may affect the care practices and outcomes recorded in the dataset. We do not investigate dataset shifts in time, which are common (75), as the eICU database includes patient records only for a year.

Conclusion

External evaluation of predictive models is important to ensure their responsible deployment in different care settings. Recommended metrics for performing such evaluation focus primarily on assessing predictive performance of the models while ignoring their potential impact on health equity. Using a large, publicly-available dataset of ICU stays from multiple hospital centers across the US, we show that models vary considerably in terms of their discriminative accuracy and calibration when validated across hospitals. Fairness of models, quantified using their differential performance on racial groups, is found to be lacking as well. Furthermore, fairness metrics continue to be poor when validating models across US geographies and hospital types. Importantly, the pattern of out-of-sample variation in the fairness metrics is not the same as that in the accuracy and calibration metrics. Thus, the standard checks do not give a comprehensive view of model performance on external datasets. This motivates the need to include fairness checks during external evaluation. While examining reasons for the lack of generalizability, we find that population demographics and clinical variables differ in their distribution across hospitals, and the race variable mediates some variation in clinical variables. Documentation of how data is generated within a hospital where a model is developed specific to sub-groups, along with development of metrics for dataset shift will be critical to anticipate where prediction models can be transferred in a trustworthy manner.

Acknowledgments

Author Contributions: All authors were involved in conceptualizing and designing the study, and in preparing the manuscript. HS and VM performed the statistical analyses. RC supervised the study and obtained funding.

Conflict of Interest Disclosures: Authors report no conflict of interest.

Funding/Support: This study was funded by the National Science Foundation grant number 1845487.

Role of the Funder/Sponsor: The funder had no role in design, conduct of the study, or in preparing the manuscript.

Data/Code Availability: All data used in the study is publicly accessible from PhysioNet website <https://physionet.org/content/eicu-crd/2.0/>. Scripts for data pre-processing are the same as that used by Johnson et al (26) and are available at <https://github.com/alistairewj/icu-model-transfer>. Scripts for model training and validation are available at <https://github.com/ChunaraLab/medshifts>.

References

1. Van Calster B, Vickers AJ. Calibration of risk prediction models: impact on decision-analytic performance. *Med Decis Making*. 2015 Feb;35(2):162–9.
2. Justice AC, Covinsky KE, Berlin JA. Assessing the generalizability of prognostic information. *Ann Intern Med*. 1999 Mar 16;130(6):515–24.
3. Moons KGM, Altman DG, Reitsma JB, Ioannidis JPA, Macaskill P, Steyerberg EW, et al. Transparent Reporting of a multivariable prediction model for Individual Prognosis or Diagnosis (TRIPOD): explanation and elaboration. *Ann Intern Med*. 2015 Jan 6;162(1):W1-73.
4. Chen IY, Szolovits P, Ghassemi M. Can AI help reduce disparities in general medical and mental health care? *AMA J Ethics*. 2019 Feb 1;21(2):E167-179.
5. Pfohl SR, Foryciarz A, Shah NH. An empirical characterization of fair machine learning for clinical risk prediction. *J Biomed Inform*. 2021 Jan;113(103621):103621.
6. Beaulieu-Jones BK, Yuan W, Brat GA, Beam AL, Weber G, Ruffin M, et al. Machine learning for patient risk stratification: standing on, or looking over, the shoulders of clinicians? *NPJ Digit Med*. 2021 Mar 30;4(1):62.

7. Riley RD, Debray TPA, Collins GS, Archer L, Ensor J, van Smeden M, et al. Minimum sample size for external validation of a clinical prediction model with a binary outcome. *Stat Med* [Internet]. 2021 May 24;(sim.9025). Available from: <http://dx.doi.org/10.1002/sim.9025>
8. Pavlou M, Qu C, Omar RZ, Seaman SR, Steyerberg EW, White IR, et al. Estimation of required sample size for external validation of risk models for binary outcomes. *Stat Methods Med Res*. 2021 Apr 21;9622802211007522.
9. Steyerberg EW, Harrell FE Jr. Prediction models need appropriate internal, internal-external, and external validation. *J Clin Epidemiol*. 2016 Jan;69:245–7.
10. Wawira Gichoya J, McCoy LG, Celi LA, Ghassemi M. Equity in essence: a call for operationalising fairness in machine learning for healthcare. *BMJ Health Care Inform* [Internet]. 2021 Apr;28(1). Available from: <http://dx.doi.org/10.1136/bmjhci-2020-100289>
11. Kaushal A, Altman R, Langlotz C. Geographic distribution of US cohorts used to train deep learning algorithms. *JAMA*. 2020 Sep 22;324(12):1212–3.
12. Yadlowsky S, Hayward RA, Sussman JB, McClelland RL, Min Y-I, Basu S. Clinical implications of revised pooled cohort equations for estimating atherosclerotic cardiovascular disease risk. *Ann Intern Med*. 2018 Jul 3;169(1):20–9.
13. Pierson E, Cutler DM, Leskovec J, Mullainathan S, Obermeyer Z. An algorithmic approach to reducing unexplained pain disparities in underserved populations. *Nat Med*. 2021 Jan;27(1):136–40.

14. Obermeyer Z, Powers B, Vogeli C, Mullainathan S. Dissecting racial bias in an algorithm used to manage the health of populations. *Science* (80-). 2019;366(6464):447–53.
15. Chouldechova A, Roth A. A snapshot of the frontiers of fairness in machine learning. *Commun ACM*. 2020 Apr 20;63(5):82–9.
16. Seyyed-Kalantari L, Liu G, McDermott M, Chen I, Ghassemi M. Medical imaging algorithms exacerbate biases in underdiagnosis [Internet]. Research Square. Research Square; 2021. Available from: <http://dx.doi.org/10.21203/rs.3.rs-151985/v1>
17. Barda N, Yona G, Rothblum GN, Greenland P, Leibowitz M, Balicer R, et al. Addressing bias in prediction models by improving subpopulation calibration. *J Am Med Inform Assoc*. 2021 Mar 1;28(3):549–58.
18. Steyerberg EW. *Clinical prediction models*. 2nd ed. Cham, Switzerland: Springer Nature; 2020. 558 p. (Statistics for Biology and Health).
19. Le Gall JR, Lemeshow S, Saulnier F. A new Simplified Acute Physiology Score (SAPS II) based on a European/North American multicenter study. *JAMA*. 1993;270(24):2957–63.
20. Zimmerman JE, Kramer AA, McNair DS, Malila FM. Acute Physiology and Chronic Health Evaluation (APACHE) IV: hospital mortality assessment for today’s critically ill patients. *Crit Care Med*. 2006 May;34(5):1297–310.
21. Zimmerman JE, Kramer AA, McNair DS, Malila FM, Shaffer VL. Intensive care unit length of stay: Benchmarking based on Acute Physiology and Chronic Health Evaluation (APACHE) IV. *Crit Care Med*. 2006 Oct;34(10):2517–29.

Generalizability Challenges of Mortality Risk Prediction Models

22. Brajer N, Cozzi B, Gao M, Nichols M, Revoir M, Balu S, et al. Prospective and external evaluation of a machine learning model to predict in-hospital mortality of adults at time of admission. *JAMA Netw Open*. 2020 Feb 5;3(2):e1920733.
23. Sendak MP, Ratliff W, Sarro D, Alderton E, Futoma J, Gao M, et al. Real-world integration of a sepsis deep learning technology into routine clinical care: Implementation study. *JMIR Med Inform*. 2020 Jul 15;8(7):e15182.
24. Johnson AEW, Mark RG. Real-time mortality prediction in the Intensive Care Unit. *AMIA Annu Symp Proc*. 2017;2017:994–1003.
25. Beam AL, Manrai AK, Ghassemi M. Challenges to the reproducibility of machine learning models in health care. *JAMA*. 2020 Jan 28;323(4):305–6.
26. Pollard TJ, Johnson AEW, Raffa JD, Celi LA, Mark RG, Badawi O. The eICU Collaborative Research Database, a freely available multi-center database for critical care research. *Sci Data*. 2018 Sep 11;5(1):180178.
27. Johnson AEW, Pollard TJ, Naumann T. Generalizability of predictive models for intensive care unit patients [Internet]. *arXiv [cs.LG]*. 2018. Available from: <http://arxiv.org/abs/1812.02275>
28. Cosgriff CV, Celi LA, Ko S, Sundaresan T, Armengol de la Hoz MÁ, Kaufman AR, et al. Developing well-calibrated illness severity scores for decision support in the critically ill. *NPJ Digit Med*. 2019 Aug 15;2(1):76.
29. von Elm E, Altman DG, Egger M, Pocock SJ, Gøtzsche PC, Vandenbroucke JP, et al. The Strengthening the Reporting of Observational Studies in Epidemiology (STROBE) statement: guidelines for reporting observational studies. *Epidemiology*. 2007 Nov;18(6):800–4.

30. Miles O. 200 400. Census regions and divisions of the United States [Internet]. [cited 2021 Jul 17]. Available from: https://www2.census.gov/geo/pdfs/maps-data/maps/reference/us_regdiv.pdf
31. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine learning in Python. the Journal of machine Learning research. 2011;12:2825–30.
32. Van Calster B, McLernon DJ, van Smeden M, Wynants L, Steyerberg EW, Topic Group ‘Evaluating diagnostic tests and prediction models’ of the STRATOS initiative. Calibration: the Achilles heel of predictive analytics. BMC Med. 2019 Dec 16;17(1):230.
33. Jiang Y, Krishnan D, Mobahi H, Bengio S. Predicting the Generalization Gap in Deep Networks with Margin Distributions. In: International Conference on Learning Representations [Internet]. 2019. Available from: <https://openreview.net/forum?id=HJlQfnCqKX>
34. Wessler BS, Ruthazer R, Udelson JE, Gheorghiu M, Zannad F, Maggioni A, et al. Regional validation and recalibration of clinical predictive models for patients with acute heart failure. J Am Heart Assoc [Internet]. 2017 Nov 18;6(11). Available from: <http://dx.doi.org/10.1161/JAHA.117.006121>
35. Paulus JK, Kent DM. Predictably unequal: understanding and addressing concerns that algorithmic clinical prediction may increase health disparities. NPJ Digit Med. 2020 Jul 30;3(1):99.
36. Gretton A, Borgwardt KM, Rasch MJ, Schölkopf B, Smola A. A kernel two-sample test. J Mach Learn Res. 2012;13(1):723–73.

37. Mooij JM, Magliacane S, Claassen T. Joint Causal Inference from Multiple Contexts. *J Mach Learn Res.* 2020;21(99):1–108.
38. Spirtes P, Glymour C, Scheines R. Causation, prediction, and search. 1993rd ed. New York, NY: Springer; 2012. 554 p. (Lecture Notes in Statistics).
39. Apolone G, Bertolini G, D’Amico R, Iapichino G, Cattaneo A, De Salvo G, et al. The performance of SAPS II in a cohort of patients admitted to 99 Italian ICUs: results from GiViTI. Gruppo Italiano per la Valutazione degli interventi in Terapia Intensiva. *Intensive Care Med.* 1996 Dec;22(12):1368–78.
40. Katsounas A, Kamacharova I, Tyczynski B, Eggebrecht H, Erbel R, Canbay A, et al. The predictive performance of the SAPS II and SAPS 3 scoring systems: A retrospective analysis. *J Crit Care.* 2016 Jun;33:180–5.
41. Harrison DA, Brady AR, Parry GJ, Carpenter JR, Rowan K. Recalibration of risk prediction models in a large multicenter cohort of admissions to adult, general critical care units in the United Kingdom. *Crit Care Med.* 2006 May;34(5):1378–88.
42. Beck DH, Smith GB, Pappachan JV, Millar B. External validation of the SAPS II, APACHE II and APACHE III prognostic models in South England: a multicentre study. *Intensive Care Med.* 2003 Feb;29(2):249–56.
43. Minne L, Eslami S, de Keizer N, de Jonge E, de Rooij SE, Abu-Hanna A. Effect of changes over time in the performance of a customized SAPS-II model on the quality of care assessment. *Intensive Care Med.* 2012 Jan;38(1):40–6.
44. Johnson A. Understanding why black patients have worse coronary heart disease outcomes: Does the answer lie in knowing where patients seek care? *J Am Heart Assoc.* Ovid Technologies (Wolters Kluwer Health); 2019 Dec 3;8(23):e014706.

45. Baldo MP, Cunha RS, Ribeiro ALP, Lotufo PA, Chor D, Barreto SM, et al. Racial differences in arterial stiffness are mainly determined by blood pressure levels: Results from the ELSA-brasil study. *J Am Heart Assoc* [Internet]. 2017 Jun 21;6(6). Available from: <http://dx.doi.org/10.1161/JAHA.117.005477>
46. Pappachan JV, Millar B, Bennett ED, Smith GB. Comparison of outcome from intensive care admission after adjustment for case mix by the APACHE III prognostic system. *Chest*. 1999 Mar;115(3):802–10.
47. Wang X, Liang G, Zhang Y, Blanton H, Bessinger Z, Jacobs N. Inconsistent performance of deep learning models on mammogram classification. *J Am Coll Radiol*. 2020 Jun;17(6):796–803.
48. Gola D, Erdmann J, Läll K, Mägi R, Müller-Myhsok B, Schunkert H, et al. Population bias in polygenic risk prediction models for coronary artery disease. *Circ Genom Precis Med*. 2020 Dec;13(6):e002932.
49. Reps JM, Kim C, Williams RD, Markus AF, Yang C, Duarte-Salles T, et al. Implementation of the COVID-19 Vulnerability Index Across an International Network of Health Care Data Sets: Collaborative External Validation Study. *JMIR Med Inf* [Internet]. 2021 Apr;9(4):e21547. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/33661754>
50. Wiens J, Guttat J, Horvitz E. A study in transfer learning: leveraging data from multiple hospitals to enhance hospital-specific predictions. *J Am Med Inform Assoc*. 2014 Jul;21(4):699–706.
51. Gong JJ, Sundt TM, Rawn JD, Guttat JV. Instance weighting for patient-specific risk stratification models. In: *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '15* [Internet]. New York,

New York, USA: ACM Press; 2015. Available from:

<http://dx.doi.org/10.1145/2783258.2783397>

52. Curth A, Thorat P, van den Wildenberg W, Bijlstra P, de Bruin D, Elbers P, et al. Transferring clinical prediction models across hospitals and electronic health record systems. In: Machine Learning and Knowledge Discovery in Databases. Cham: Springer International Publishing; 2020. p. 605–21. (Communications in computer and information science).
53. Desautels T, Calvert J, Hoffman J, Mao Q, Jay M, Fletcher G, et al. Using transfer learning for improved mortality prediction in a data-scarce hospital setting. Biomed Inform Insights. 2017 Jun 12;9:1178222617712994.
54. Davis SE, Lasko TA, Chen G, Siew ED, Matheny ME. Calibration drift in regression and machine learning models for acute kidney injury. J Am Med Inform Assoc. 2017 Nov 1;24(6):1052–61.
55. Granholm A, Møller MH, Krag M, Perner A, Hjortrup PB. Predictive performance of the Simplified Acute Physiology Score (SAPS) II and the initial Sequential Organ Failure Assessment (SOFA) score in acutely ill intensive care patients: Post-hoc analyses of the SUP-ICU inception cohort study. PLoS One. 2016 Dec 22;11(12):e0168948.
56. Nestor B, McDermott MBA, Boag W, Berner G, Naumann T, Hughes MC, et al. Feature Robustness in Non-stationary Health Records: Caveats to Deployable Model Performance in Common Clinical Machine Learning Tasks. In: Doshi-Velez F, Fackler J, Jung K, Kale D, Ranganath R, Wallace B, et al., editors. Proceedings of the 4th Machine Learning for Healthcare Conference [Internet]. PMLR; 2019. p. 381–405. (Proceedings of Machine Learning Research; vol. 106). Available from: <https://proceedings.mlr.press/v106/nestor19a.html>

57. Austin PC, van Klaveren D, Vergouwe Y, Nieboer D, Lee DS, Steyerberg EW. Geographic and temporal validity of prediction models: different approaches were useful to examine model performance. *J Clin Epidemiol*. 2016 Nov;79:76–85.
58. Roth HR, Chang K, Singh P, Neumark N, Li W, Gupta V, et al. Federated learning for breast density classification: A real-world implementation. In: *Domain Adaptation and Representation Transfer, and Distributed and Collaborative Learning*. Cham: Springer International Publishing; 2020. p. 181–91. (Lecture notes in computer science).
59. Nsoesie EO. Evaluating artificial intelligence applications in clinical settings. *JAMA Netw Open*. 2018 Sep 28;1(5):e182658.
60. Futoma J, Simons M, Panch T, Doshi-Velez F, Celi LA. The myth of generalisability in clinical research and machine learning in health care. *Lancet Digit Health*. 2020 Sep;2(9):e489–92.
61. Agniel D, Kohane IS, Weber GM. Biases in electronic health record data due to processes within the healthcare system: retrospective observational study. *BMJ*. 2018 Apr 30;k1479.
62. Debray TPA, Vergouwe Y, Koffijberg H, Nieboer D, Steyerberg EW, Moons KGM. A new framework to enhance the interpretation of external validation studies of clinical prediction models. *J Clin Epidemiol*. 2015 Mar;68(3):279–89.
63. Mustafa B, Loh A, Freyberg J, MacWilliams P, Wilson M, McKinney SM, et al. Supervised transfer learning at scale for medical imaging [Internet]. *arXiv [cs.CV]*. 2021. Available from: <http://arxiv.org/abs/2101.05913>

64. Ke A, Ellsworth W, Banerjee O, Ng AY, Rajpurkar P. CheXtransfer. In: Proceedings of the Conference on Health, Inference, and Learning [Internet]. New York, NY, USA: ACM; 2021. Available from: <http://dx.doi.org/10.1145/3450439.3451867>
65. Subbaswamy A, Saria S. From development to deployment: dataset shift, causality, and shift-stable models in health AI. *Biostatistics*. 2020 Apr 1;21(2):345–52.
66. Subbaswamy A, Saria S. I-SPEC: An end-to-end framework for learning transportable, shift-stable models [Internet]. *arXiv [stat.ML]*. 2020. Available from: <http://arxiv.org/abs/2002.08948>
67. Singh H, Singh R, Mhasawade V, Chunara R. Fairness violations and mitigation under covariate shift. In: Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency [Internet]. New York, NY, USA: ACM; 2021. Available from: <http://dx.doi.org/10.1145/3442188.3445865>
68. Wenneker MB. Racial inequalities in the use of procedures for patients with ischemic heart disease in Massachusetts. *JAMA*. 1989 Jan 13;261(2):253–7.
69. Kjellstrand CM. Age, sex, and race inequality in renal transplantation. *Arch Intern Med*. 1988 Jun 1;148(6):1305.
70. Yergan J, Flood AB, LoGerfo JP, Diehr P. Relationship between patient race and the intensity of hospital services. *Med Care*. 1987 Jul;25(7):592–603.
71. Blendon RJ, Aiken LH, Freeman HE, Corey CR. Access to medical care for black and white Americans. A matter of continuing concern. *JAMA*. 1989 Jan 13;261(2):278–81.
72. Gebru T, Morgenstern J, Vecchione B, Vaughan JW, Wallach H, III HD, et al. Datasheets for Datasets. *Commun ACM [Internet]*. 2021 Nov;64(12):86–92. Available from: <https://doi.org/10.1145/3458723>

73. Mitchell M, Wu S, Zaldivar A, Barnes P, Vasserman L, Hutchinson B, et al. Model cards for model reporting. In: Proceedings of the Conference on Fairness, Accountability, and Transparency [Internet]. New York, NY, USA: ACM; 2019. Available from: <http://dx.doi.org/10.1145/3287560.3287596>
74. Norgeot B, Quer G, Beaulieu-Jones BK, Torkamani A, Dias R, Gianfrancesco M, et al. Minimum information about clinical artificial intelligence modeling: the MI-CLAIM checklist. *Nat Med*. 2020 Sep;26(9):1320–4.
75. Sáez C, Gutiérrez-Sacristán A, Kohane I, García-Gómez JM, Avillach P. EHR temporal Variability: delineating temporal data-set shifts in electronic health records. *Gigascience* [Internet]. 2020 Aug 1;9(8). Available from: <http://dx.doi.org/10.1093/gigascience/giaa079>
76. Finlayson SG, Subbaswamy A, Singh K, Bowers J, Kupke A, Zittrain J, et al. The Clinician and Dataset Shift in Artificial Intelligence. *N Engl J Med*. 2021 Jul 15;385(3):283–6.
77. Little RJA, Rubin DB. Statistical analysis with missing data. Third edit. *Statistical Analysis with Missing Data*. 2014. 1–381 p. (Wiley series in probability and statistics).

Figure Titles

Figure 1. Generalization of performance metrics across individual hospitals.

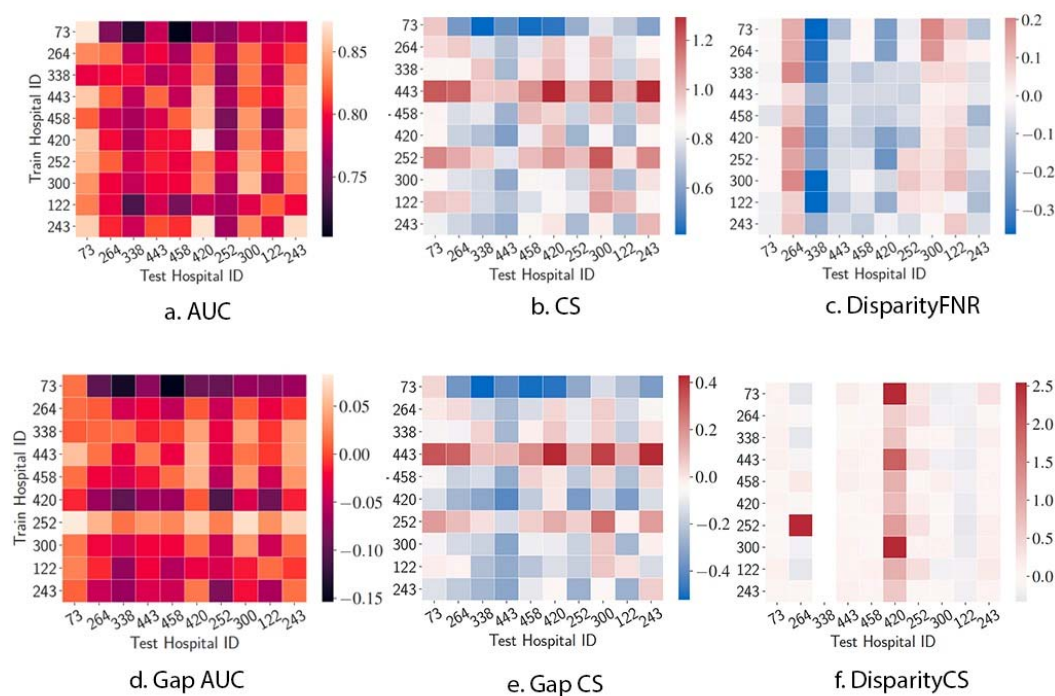
Figure 2. Generalization of performance metrics across US geographic regions.

Figure 3. Statistical tests for dataset shifts.

Figure 4. Shifts in variable distributions due to hospital, region, and other factors based on mortality causal graph.

Figures

Figure 1. Generalization of performance metrics across individual hospitals. Results of transferring models across top 10 hospitals by number of stays. Models are trained and tested on a fixed number of samples (1631, the least in any of the 10 hospitals) from each hospital. Results are averaged over 100 random subsamples for each of the 10×10 train-test hospital pairs. All 6 metrics show large variability when transferring models across hospitals. Abbreviations: AUC, area under ROC curve; CS, calibration slope; FNR, false negative rate.



Generalizability Challenges of Mortality Risk Prediction Models

Figure 2. Generalization of performance metrics across US geographic regions. Results of transferring models after pooling hospitals into 4 regions (northeast, south, midwest, west). Models are trained and tested on 5000 samples from each region. Results are averaged over 100 random subsamples for each of the 4×4 train-test hospital pairs. DisparityFNR and DisparityCS show large variability when transferring models across regions. Abbreviations: AUC, area under ROC curve; CS, calibration slope; FNR, false negative rate.

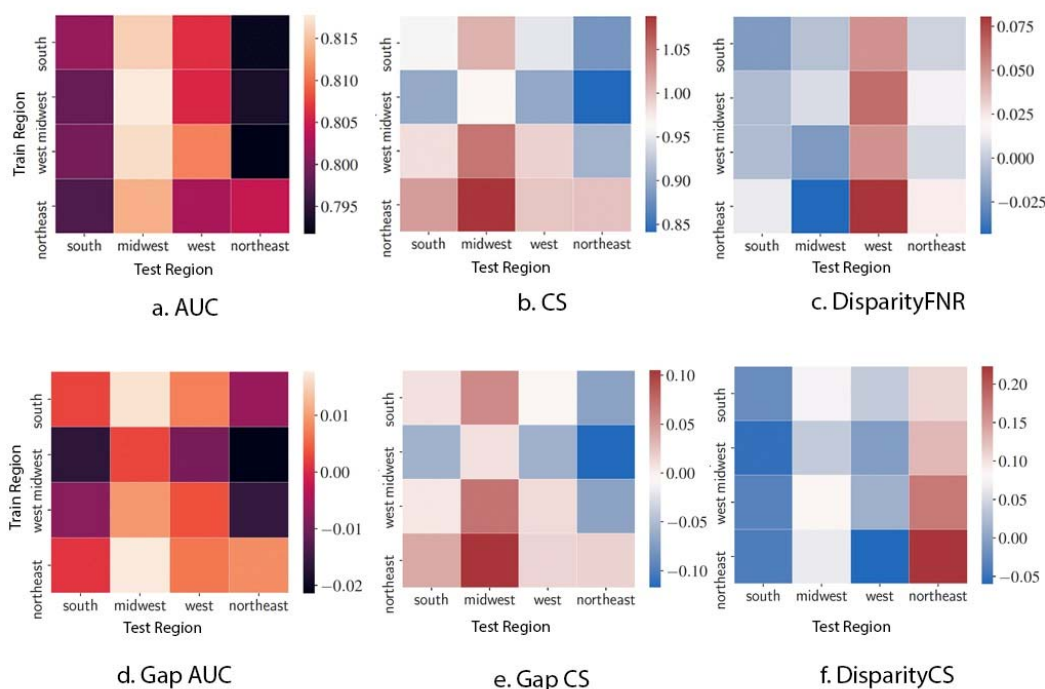


Figure 3. Statistical tests for dataset shifts. Results for two-sample tests with and without pooling of hospitals by region. Test results are plotted in (a,c) and test statistics are plotted in (b,d) to examine the test results in more detail. Since the order of hospitals considered in the two-sample test does not change the test statistic, we plot only the lower halves of the matrices. Results are averaged over 100 random subsamples. Feature distribution changes across all hospital and region pairs.

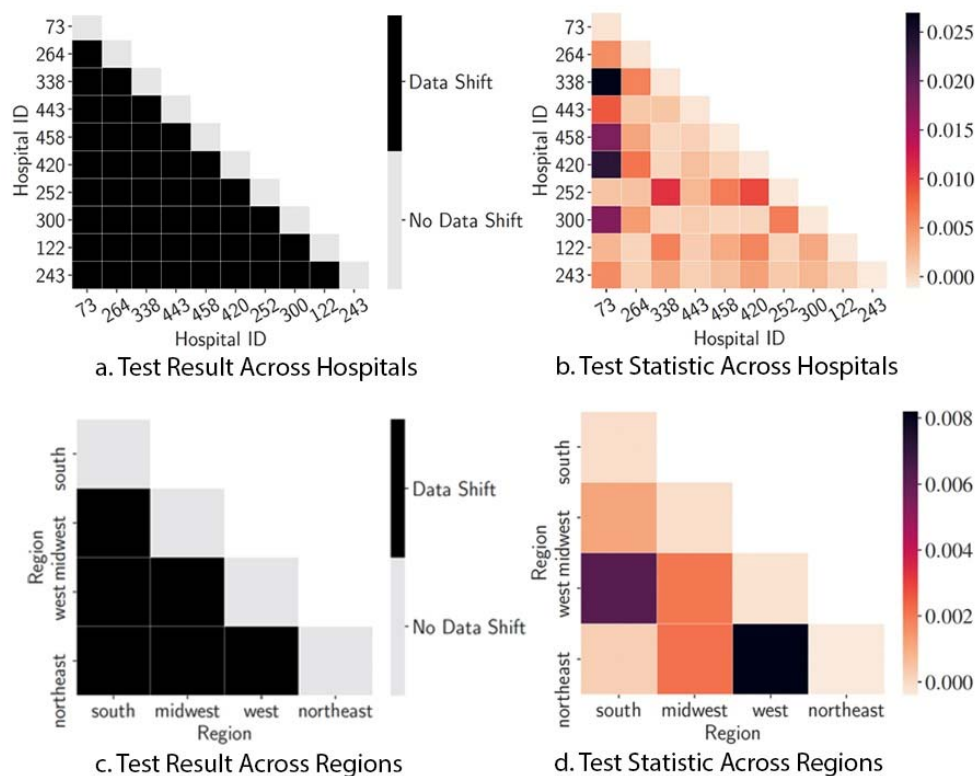


Figure 4. Shifts in variable distributions due to hospital, region, and other factors based on mortality causal graph. Each row represents (in red) the features which explain the shifts across each of the indicators labeling the row, i.e. the features with an edge from the indicator in the causal graph. For instance, shift across the hospital ID indicator (first row) is explained by shifts in the distributions of age, race, temp (or temperature), urine output, and so on. We observe that shifts are explained by changes in a few variables which are common across indicators. Full forms of the abbreviated feature names are added in Table S3.

