

MIMIC-IF: Interpretability and Fairness Evaluation of Deep Learning Models on MIMIC-IV Dataset

Chuiheng Meng

University of Southern California

Loc Trinh (✉ loctrinh@usc.edu)

University of Southern California

Nan Xu

University of Southern California

Yan Liu

University of Southern California

Research Article

Keywords: large-scale healthcare datasets, black-boxed model, MIMIC-IV, IMV-LSTM

Posted Date: April 14th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-402058/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

MIMIC-IF: Interpretability and Fairness Evaluation of Deep Learning Models on MIMIC-IV Dataset

Chuizheng Meng^{1,+,*}, Loc Trinh^{1,+,*}, Nan Xu^{1,+,*}, and Yan Liu^{1,*}

¹Department of Computer Science, University of Southern California, Los Angeles, CA 90089

*chuizhem@usc.edu, loctrinh@usc.edu, nanx@usc.edu, yanliu.cs@usc.edu

+these authors contributed equally to this work

ABSTRACT

The recent release of large-scale healthcare datasets has greatly propelled the research of data-driven deep learning models for healthcare applications. However, due to the nature of such deep black-boxed models, concerns about interpretability, fairness, and biases in healthcare scenarios where human lives are at stake call for a careful and thorough examination of both datasets and models. In this work, we focus on MIMIC-IV (Medical Information Mart for Intensive Care, version IV), the largest publicly available healthcare dataset, and conduct comprehensive analyses of dataset representation bias as well as interpretability and prediction fairness of deep learning models for in-hospital mortality prediction. In terms of interpretability, we observe that (1) the best-performing interpretability method successfully identifies critical features for mortality prediction on various prediction models; (2) demographic features are important for prediction. In terms of fairness, we observe that (1) there exists disparate treatment in prescribing mechanical ventilation among patient groups across ethnicity, gender and age; (2) all of the studied mortality predictors are generally fair while the IMV-LSTM (Interpretable Multi-Variable Long Short-Term Memory) model provides the most accurate and unbiased predictions across all protected groups. We further draw concrete connections between interpretability methods and fairness metrics by showing how feature importance from interpretability methods can be beneficial in quantifying potential disparities in mortality predictors.

1 Introduction

With the release of large scale healthcare datasets, research of data-driven deep learning methods for healthcare applications demonstrates their superior performance over traditional methods on various tasks, including mortality prediction, length-of-stay prediction, phenotyping classification and intervention prediction¹⁻³. However, deep learning models have been treated as black-box universal function approximators, where prediction explanations are no longer available as their traditional counterparts, e.g., Logistic Regression and Random Forests. Lack of interpretability hinders the wide application of deep learning models in critical domains like healthcare. In addition, due to bias in datasets or models, decisions made by machine learning algorithms are prone to be unfair, where an individual or a group is favored compared with the others owing to their inherent traits. As a result, more and more concerns about interpretability, fairness and biases have been raised recently in the healthcare domain where human lives are at stake⁴. These concerns call for careful and thorough analyses of both datasets and algorithms. In this work, we focus on the latest version (version IV⁵) of a widely used large scale healthcare dataset MIMIC⁶, and conduct comprehensive analyses of model interpretability, dataset bias, algorithmic fairness, and the interaction between interpretability and fairness.

Interpretability evaluation. First, we evaluate the performance of common interpretability methods for feature importance estimation on multiple deep learning models trained for the mortality prediction task. Due to the complexity of dynamics in electronic health record data, there is no access to the ground truth of feature importance. Therefore, we utilize ROAR (remove and retrain)⁷ to quantitatively evaluate different feature importance estimations. On all models considered, the ArchDetect⁸ outperforms other interpretation methods in feature importance estimation. Then we qualitatively analyze the feature importance estimation results given by ArchDetect, and verify its effectiveness based on the observations that it successfully identifies critical features for mortality prediction. We also find that demographic features are important for prediction, which leads to our following analyses of dataset bias and algorithmic fairness.

Dataset bias and algorithmic fairness. We adopt the following commonly used demographic features as protected attributes: 1) *ethnicity*, 2) *gender*, 3) *marital status*, 4) *age*, and 5) *insurance type*. For dataset bias, we analyze the average adoption and duration of five types of ventilation treatment on patients from different groups. There exists treatment disparity among patient groups split by different protected attributes, which is most evident across different ethnic groups: Black and Hispanic cohorts are less likely to receive ventilation treatments, as well as shorter treatment duration on average. However, multiple confounders

may lead to the observed disparity in treatment, which adds the difficulty of identifying intentional discriminations. Hence we call for a close look at causal analysis for a better understanding. For algorithmic fairness, we evaluate the performance of state-of-the-art machine learning approaches for mortality prediction in terms of AUC-based fairness metrics. Experiment results indicate a strong correlation between mortality rates and fairness: machine learning approaches tend to obtain lower AUC scores on groups with higher mortality rates. Meanwhile, all of the studied mortality predictors are fair in general while IMV-LSTM⁹ performs the best overall across protected groups.

Interactions between interpretability and fairness. We examine the interaction of interpretability and fairness by drawing connections between feature importance and fairness metrics. Furthermore, we observe substantial disparities in the importance of each demographic feature used for in-mortality prediction across the protected subgroups, which raises a concern about whether these demographic features should be used in mortality prediction.

In summary, our main contributions are:

- We give quantitative evaluations of popular interpretability methods for feature importance estimation on deep learning models in the context of mortality prediction. We observed that the best-performing interpretability method successfully identifies critical clinical and demographic features for mortality prediction on various models.
- For dataset bias, we observe treatment disparity among patient groups split by different protected attributes.
- For algorithmic fairness, we find that all of the studied mortality predictors are fair in general while the IMV-LSTM model performs the best overall across different protected groups.
- We also examine the interaction between interpretability and fairness, and observe disparities of feature importance among demographic subgroups.

2 Related Work

2.1 Interpretability Evaluation

2.1.1 Interpretability of Deep Learning Models

Due to the complexity of deep learning models, interpretability research has developed diversely, and many methods have been used to interpret how a deep learning model works from various aspects, including: (1) *Feature importance estimation*^{10–20}. For a given data sample, these methods estimate the importance of each input feature with respect to a specified output. (2) *Feature interaction attribution*^{8,21–25}. In addition to estimating the importance of individual features, these methods analyze how interactions of feature pairs/groups contribute to predictions. (3) *Neuron/layer attribution*^{20,26–29}. These methods estimate the contribution of specified layers/neurons in the model. (4) *Explanation with high-level concepts*^{30–32}. These methods interpret deep learning models with human-friendly concepts instead of the importance of low-level input features. In this paper, we focus on feature importance estimation due to its importance and the completeness of its evaluation methods.

2.1.2 Evaluation of Feature Importance Interpretation

Since feature importance estimation assigns an importance score for each input feature, the evaluation of results is equivalent to the evaluation of binary classification results when the ground truth of feature importance is available, where the label indicates whether the feature is important for the problem.³³ constructs synthetic datasets with feature importance labels for evaluation.³⁴ obtains feature importance labels from both manually constructed tasks and domain experts.³⁵ derives importance labels from tasks with graph-valued data with computable ground truths. However, these evaluation methods require the accessibility of ground truth labels, which is hard to fulfill and is usually the problem itself we need to solve in domains such as healthcare.

For evaluation without ground truth, A common strategy to evaluate feature importance estimation is to measure the degradation of model performance with the gradual removal of features estimated to be important.³⁶ perturbs features ranked by importance in test samples and calculates the area over the MoRF curve (AOPC): a higher AOPC means the information disappears faster with feature removal and indicates a better importance estimation.⁷ remove features from the entire dataset and retrain the model when obtaining AOPC, which excludes the interference of data distribution shifting.³³ replace features with known feature distributions for evaluation on synthetic tasks to ensure the consistency of data distribution. In this paper, we utilize the evaluation in⁷.

2.2 Fairness Evaluation

2.2.1 Bias and Fairness in Machine Learning

With the open access to large-scale datasets and the development of machine learning algorithms, more decisions in the real world are made by machine learning algorithms with or without human's intervention, e.g., job advertisements promoting³⁷, facial recognition³⁸, treatment recommendation³⁹, etc. Due to bias in datasets or models, decisions made by machine learning algorithms are prone to be unfair, where an individual or a group is favored compared with the others owing to their inherent

traits. One well-known example is the software COMPAS (Correctional Offender Management Profiling for Alternative Sanctions), which was found a bias against African-Americans to assign a higher risk score of recommitting another crime than to Caucasians with the same profile⁴⁰.

Based on the general assumption that the algorithm itself is not coded to be biased, the decision unfairness can be attributed to biases in the data, which is likely to be picked up and amplified by the trained algorithm⁴¹. Three major sources of data biases are⁴¹: 1) *Biased Labels*: the ground-truth labels for the machine learning algorithms to predict are biased; 2) *Imbalanced representation*: imbalanced representation of different demographic groups occurs when some protected groups are underrepresented with fewer observations in the dataset compared with other groups; 3) *Data Quality Disparity*: data from protected groups might be less complete or accurate during data collecting and processing. Mostly widely considered traits, such as gender, age, ethnicity, marital status, are considered as protected or sensitive attributes in literature⁴². Fairness has been defined in various ways considering different contexts or applications, two of them are the most widely leveraged for bias detection and correction: *Equal Opportunity*, where the predictions are required to have equal true positive rate across two demographics, and *Equalized Odds*, where an additional constraint is put on the predictor to have equal false positive rate⁴³. To derive fair decisions with machine learning algorithms, three categories of approaches have been proposed to mitigate biases^{42,44}: 1) *Pre-processing*: the original dataset is transformed so that the underlying discrimination towards some groups is removed⁴⁵; 2) *In-processing*: either by adding a penalization term in the objective function⁴⁶ or imposing a fairness-relevant constraint⁴⁷; 3) *Post-processing*: further recompute the results from predictors to improve fairness⁴⁸.

2.2.2 Bias and Fairness in MIMIC-III

With clinical notes^{49,50} or temporal measurements^{4,51,52} or both⁵³ from MIMIC-III considered, fairness evaluation and bias mitigation have been studied recently for tasks such as mortality prediction^{4,49-53}, phenotyping^{50,53}, readmission⁵¹, length of stay⁵², etc. To evaluate data and prediction fairness for the aforementioned healthcare tasks, attributes like ethnicity^{4,49,50,52,53}, gender^{50,52,53}, insurance^{50,53}, age⁴⁹ and language⁵⁰, are considered most often to split patients into different protected groups.

When making medical decisions based on text data like clinical notes, word embeddings, used as machine learning inputs, have been demonstrated to propagate unwanted relationships with regard to different genders, language speakers, ethnicities, and insurance groups^{50,53}. With respect to gender and insurance type, differences in accuracy and therefore machine bias has been observed for mortality prediction⁵¹. To mitigate biases and improve prediction fairness, Chen *et al.* argued that collecting data with adequate sample sizes and predictive variables measures is an effective approach to reduce discrimination without sacrificing accuracy⁴. Martinez *et al.* proposed an in-processing approach where the fairness problem is characterized as a multi-objective optimization task, where the risk for each protected group is a separate objective⁴⁹. After well-trained machine learning models make predictions, equalized odds post-processing⁵³ and updating predictions according to the weighted sum of utility and fairness⁵² were introduced respectively as effective post-processing approaches.

To continue the dataset bias and algorithmic fairness study on MIMIC-IV, we follow previous fairness study work and adopt the following commonly used demographic features as protected attributes: 1) *ethnicity*, 2) *gender*, 3) *marital status*, 4) *Age*, and 5) *insurance type*. For dataset bias, we analyze the average adoption and duration of five types of ventilation treatment on patients from different groups. For algorithmic fairness, we evaluate the performance of state-of-the-art machine learning approaches for mortality prediction in terms of accuracy and fairness.

2.3 Interactions between Interpretability and Fairness

Besides accuracy, interpretability and fairness are two important aspects that businesses and researchers should take into consideration when designing, deploying, and maintaining machine learning models⁵⁴. It is also well acknowledged that enhancing model interpretability is an important step towards developing fairer ML systems⁵⁵ since interpretations can help detecting and mitigating bias during data collection or labeling⁵⁶⁻⁵⁸. Given evaluation metrics from the two concepts, demonstrations of performance from different predictive models have been shown in literature to further investigate their interactions⁵⁹⁻⁶³. When the model's complexity is determined by the number of features and simpler models are more interpretable, curves showing how model fairness is affected by model complexity were studied besides its influence on accuracy^{59,60}. When the feature importance is leveraged to interpret model predictions, failure of fairness can be identified by detecting whether the feature has a larger effect than it should have^{61,62}. For instance, Adebayo *et al.* showed that gender is of low importance among all studied demographic features in a bank's credit limit model, which indicates that the bank's algorithm is not overly dependent on gender in making credit limit determinations⁶¹. Recently, connections between interpretability and fairness were quantitatively studied by comparing fairness measures and feature importance measure: there is a direct relation between SHAP value difference and equality of opportunity after removing bias with reweighing techniques and measuring feature importance with SHAP on Adult, German, Default and COMPAS datasets⁶³. Given mortality predictions made by state-of-the-art models on MIMIC-IV, we study the connections between feature importance induced by different interpretation approaches and the fairness measures in this paper.

3 MIMIC-IV Dataset

In this section, we describe the following preprocessing steps of the MIMIC-IV dataset: cohort selection, feature selection, and data cleaning. We also report the distributions of demographic, admission and comorbidity variables within the dataset.

3.1 Dataset Description

MIMIC-IV^{5,6} is a publicly available database of patients admitted to the Beth Israel Deaconess Medical Center (BIDMC) in Boston, MA, USA. It contains de-identified data of 383,220 patients admitted to an intensive care unit (ICU) or the emergency department (ED) between 2008 - 2019. Till the day when we finished all experiments, the latest version of MIMIC-IV is v0.4 and only provides public access to the electronic health record data of 50,048 patients admitted to the ICU, which is sourced from the clinical information system MetaVision at the BIDMC. Therefore, we design the following data preprocessing procedures for the ICU data part of MIMIC-IV.

3.2 Preprocessing

3.2.1 Cohort Selection

Following the common practice in^{1,3}, we select ICU stays satisfying the following criteria as the cohort: (1) the patient is at least 15 years old at the time of ICU admission; (2) the ICU stay is the first known ICU stay of the patient; (3) the total duration of ICU stay is between 12 hours and 10 days. After the cohort selection, we collect 45,768 ICU stays as the cohort. According to the cohort selection criterion (2), each ICU stay corresponds to one unique patient and one unique hospital admission.

3.2.2 Data Cleaning & Feature Selection

We follow the same data cleaning procedure in¹ to handle: (1) Inconsistent units. We convert features with multiple units to their major unit. (2) Multiple recordings at the same time. We use the average value for numerical features and the first appearing value for categorical features. (3) Range of feature values. We use the median of the range as the value of the feature.

We select 164 features from the following groups, a detailed list of all selected features is in Table A1 in Appendix:

- Electronic healthcare records (EHR). We modify the feature list used in¹ and extract 122 features after removing features that are no longer available in MIMIC-IV.
- Demographic features. We extract 5 from patients' demographic information.
- Admission features. We extract 4 from admission records.
- Comorbidity features. We extract binary flags of 33 types of comorbidity using patients' ICD codes.

3.2.3 Data Filtering, Truncation, Aggregation and Imputation

Data Filtering After specifying the list of features, we further filter ICU stays from the cohort and only keep those that have records of selected EHR features for at least 24 hours and at most 10 days, starting from the first record within 6 hours prior to ICU admission time. We have 43005 ICU stays after the filtering. Other works³ extract the first 30-hour data and drop data from the last 6 hours to avoid information leakage of positive mortality labels to features measured within 6 hours prior to deathtime. We find that most (96.02%) of the patients with positive in-hospital mortality labels have measurements for over 30 hours prior to their deathtime, thus we omit this processing step. **Truncation** For each ICU stay, we only keep the data of the first 24 hours, starting from the first record within 6 hours prior to its ICU admission time. **Aggregation** For each ICU stay, we aggregate its records hourly by taking the average of multiple records within the same hourly time window. **Imputation** We perform forward and backward imputation to fill missing values. For cases where certain features of some patients are completely missing, we fill with mean values of corresponding features in the training set.

3.3 Dataset Summary

After all preprocessing steps, we obtain features of the shape (N, T, F) , where $N = 43005$ is the number of ICU stays (data samples), $T = 24$ is the number of time steps with 1-hour step size, and $F = 164$ is the total number of features. We also process the data into the tabular form (N, F') by replacing sequential EHR features with the summary over time steps including minimum, maximum, and mean values (for the urinary_output_sum feature we have summation in addition), where $F' = 409$. We show the distribution of demographic, admission, and comorbidity features grouped by patients' in-hospital mortality status in Table A2 in Appendix. We also demonstrate differences between the preprocessed MIMIC-IV data in this work and the preprocessed MIMIC-III data from¹ in Table 1.

Table 1. Differences between preprocessed MIMIC-III in¹ and preprocessed MIMIC-IV.

	MIMIC-III	MIMIC-IV (this work)
# Samples	35627	43005
# Temporal Features	135	122
# Demographic Features	1	5
# Admission Features	1	4
# Comorbidity Features	3	33

Table 2. Classification performance of all considered deep models.

AutoInt		LSTM		TCN		Transformer		IMVLSTM	
AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC
0.508	0.901	0.660	0.938	0.666	0.928	0.686	0.939	0.769	0.955

4 Interpretability Evaluation

In this section, we evaluate the performance of various feature importance interpretability methods on multiple models for the in-hospital mortality prediction task. We describe the task, models, interpretability methods, and the evaluation method in detail and report the evaluation results.

4.1 Task Description

Mortality prediction is one primary outcome of high interest of hospital admissions, and is widely considered in other benchmark works^{1-3,64}. We use the in-hospital mortality prediction task to train different models and evaluate the performance of various interpretability methods. We formulate the in-hospital mortality prediction task as a binary classification task. Given the observed sequence of features $\mathbf{X} \in \mathbb{R}^{T \times F}$ of one patient (or its summary $\mathbf{x} \in \mathbb{R}^F$, depending on the model), the model gives the probability that the patient dies during his/her hospital admission after being admitted to ICU. In MIMIC-IV, a patient has in-hospital mortality if and only if his/her deathtime exists in the `mimic_core.admissions` table. We randomly divide 60% data for training, 20% for validation and 20% for test.

4.2 Models

We consider following models: (1) **AutoInt**⁶⁵. A model that learns feature interaction automatically via self-attentive neural networks. (2) **LSTM**⁶⁶. Long short-term memory recurrent neural network, which is a common baseline for sequence learning tasks. (3) **TCN**⁶⁷. Temporal convolutional networks, which outperform canonical recurrent networks across various tasks and datasets. (4) **Transformer**⁶⁸. A network architecture based solely on attention mechanisms. Here we only adopt its encoder part for the classification task. (5) **IMVLSTM**⁹. An interpretable model that jointly learns network parameters, variable and temporal importance, and gives inherent feature importance interpretation. We use sequence data as input for (2)-(5), and the summary of sequence data as input for (1) since AutoInt only processes tabular data in its original implementation.

We use the area under the precision-recall curve (AUPRC) and the area under the receiver operating characteristic curve (AUROC) as metrics for binary classification. The performance of all models considered in this work is shown in Table 2.

4.3 Interpretability Methods

Interpretation of deep learning models is still a rapidly developing area and contains various aspects. In this work, we focus on the interpretation of feature importance, which estimates the importance of single features for a given model on a specific task. Estimation of feature importance helps improve the model, builds trust in prediction and isolates undesirable behavior⁷. Recent works^{7,33,36} have developed methods for evaluating feature performance estimation without access to the ground truth of feature importance, which fits scenarios in healthcare domains well: ground-truth feature importance for healthcare applications is either the problem we need to solve itself or requires extraction from a huge amount of domain knowledge. Therefore, we choose the interpretation of feature importance as the target aspect for evaluating interpretability methods.

Formally, given a function $M : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R}^{d_{out}}$ and the input (flattened) feature vector $\mathbf{x} \in \mathbb{R}^{d_{in}}$, the interpretation of feature importance gives a non-negative score $\mathbf{s}(\mathbf{x}) \in \mathbb{R}^{d_{in}}$, where $\mathbf{s}(\mathbf{x})_i$ is the importance of x_i to $M(\mathbf{x})$.

We select the following interpretability methods to compare their feature importance estimation results. Notice that some interpretability methods give signed scores (or "attributions"), where signs reflect positive/negative contributions of features to the output, and we use the absolute values of signed scores as importance scores. For methods requiring a baseline input vector, unless otherwise specified, we follow the method in³³ and randomly sample $\mathbf{x}' \in \mathbb{R}^{d_{in}}$, where $x'_i \sim \mathcal{U}[0, 1]$.

(1) Gradient based methods.

- **Saliency**¹⁰. Saliency returns the gradients with respect to inputs as feature importance: $\mathbf{s}(\mathbf{x}) = \frac{\partial M(\mathbf{x})}{\partial \mathbf{x}}$. By taking the first-order Taylor expansion of the neural network at the input, $M(\mathbf{x}) \approx (\frac{\partial M(\mathbf{x})}{\partial \mathbf{x}})^T \mathbf{x} + b$, which is a linear approximation of the network, the gradient $\frac{\partial M(\mathbf{x})}{\partial x_i} = \mathbf{s}(\mathbf{x})_i$ is the coefficient of the i -th feature.
- **IntegratedGradients**¹¹. IntegratedGradients assigns an importance score to each input feature by approximating the integral of gradients of the model's output with respect to the inputs along the path (straight line) from given baselines

$$\text{IntegratedGradients}(\mathbf{x})_i = (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial M(\mathbf{x}' + \alpha(\mathbf{x} - \mathbf{x}'))}{\partial x_i} d\alpha, \quad (1)$$

where \mathbf{x}' is the baseline.

- **DeepLift**^{12,13}. DeepLift decomposes the output prediction of a neural network on a specific input by backpropagating the contributions of all neurons in the network to every feature of the input.

$$\text{DeepLift}(\mathbf{x})_i = (x_i - x'_i) \times \frac{\partial^g M(\mathbf{x})}{\partial x_i}, g(z_t) = \frac{f_t(z_t) - f_t(z'_t)}{z_t - z'_t}, \quad (2)$$

where $\frac{\partial^g M(\mathbf{x})}{\partial x_i} = \sum_{p \in P_{io}} (\prod_{(s,t) \in p} w_{ts} \prod_{(s,t) \in p} g(z_t))$. P_{io} is the set of all paths from the i -th input feature to the output neuron in the network. (s, t) is a pair of connected neurons in path p . Each neuron t contains a linear transformation $z_t = \sum_{q \in Pa(t)} w_{tq} o_q + b_t$ followed by a nonlinear mapping $o_t = f(z_t)$.

- **GradientShap**¹⁴. GradientShap approximates SHAP (SHapley Additive exPlanations) values by computing the expectations of gradients by randomly sampling from the distribution of baselines. It first adds white noise to each input sample and selects a random baseline from a given distribution, then selects a random point along the path between the baseline and the input with noise, and computes the gradient of outputs with respect to the random point. The procedure is repeated for multiple times to approximate the expected values of gradients $E(\frac{\partial M(\mathbf{x})}{\partial \mathbf{x}})$. The final SHAP value for the i -th feature is $E(\frac{\partial M(\mathbf{x})}{\partial x_i}) \times (x_i - x'_i)$.
- **DeepLiftShap**¹⁴. It extends DeepLift algorithm and approximates SHAP values using DeepLift. For each input, it samples baselines from a given distribution and computes the DeepLift score for each input-baseline pair and averages the resulting scores per input example as the output.
- **SaliencyNoiseTunnel**¹⁵. SaliencyNoiseTunnel adds Gaussian noise to the input sample and averages the calculated attributions using Saliency method as the output.

(2) Perturbation based methods.

- **ShapleySampling**^{16,17}. Shapley value gives attribution scores by taking each permutation of the input features and adding them one-by-one to a given value. Since the computation complexity is extremely high for large numbers of features, ShapleySampling takes some random permutations of the input features and averages the marginal contribution of features.
- **FeaturePermutation**¹⁸. FeaturePermutation permutes the input feature values randomly within a batch and computes the difference between original and shuffled outputs as the result.
- **FeatureAblation**¹⁹. FeatureAblation replaces each input feature with a given baseline value and computes the difference in output as the result.
- **Occlusion**²⁰. Occlusion replaces each contiguous rectangular region with a given baseline and computing the difference in output as the result.
- **ArchDetect**⁸. It utilizes the discrete interpretation of partial derivatives. While the original paper considers both single features and feature pairs, we here only apply it to single features, since the evaluation method in this work is designed for single feature importance only. In the single feature case, the importance score of the i -th feature is

$$\text{ArchDetect}(\mathbf{x})_i = \left(\frac{M(\mathbf{x}_{\{i\}} + \mathbf{x}'_{\setminus\{i\}}) - M(\mathbf{x}'_{\{i\}})}{x_i - x'_i} \right)^2, \text{ where } (\mathbf{x}_{\mathcal{S}})_i = \begin{cases} x_i, & \text{if } i \in \mathcal{S}; \\ 0, & \text{otherwise.} \end{cases} \quad (3)$$

Here we select $\mathbf{x}' = \mathbf{0} \in \mathbb{R}^{d_{in}}$.

(3) Glassbox interpretation. If the model's architecture provides feature importance scores directly as a part of the output of the model, such as the attention score of each feature, we call this interpretation as "Glassbox" and regard it as an extra baseline.

(4) Random baseline. As a baseline, we randomly shuffle all features as the feature importance ranking.

For models in Section 4.2, AutoInt maps categorical features to embeddings using learnable dictionaries and has no gradient on categorical features, thus gradient based methods are not applicable. Only IMVLSTM model has Glassbox interpretation.

4.4 Evaluation Method

Since acquiring the ground-truth feature importance is challenging for mortality prediction tasks, we evaluate one feature importance estimation by gradually dropping most important features it gives at certain ratios from the dataset and observe

Table 3. Area under the curve (AUC) of interpretability methods for each model and each classification performance metric evaluated using ROAR. AUC is measured for two prediction metrics (AUPRC and AUROC) respectively. Lower AUC indicates more rapid prediction performance drop and better feature importance interpretation.

Interpreters	AutoInt		LSTM		TCN		Transformer		IMVLSTM	
	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC
Random	0.401	0.842	0.615	0.909	0.605	0.901	0.662	0.918	0.669	0.915
Glassbox	\times	\times	\times	\times	\times	\times	\times	\times	0.533	0.892
Saliency	\times	\times	0.558	0.898	0.587	0.893	0.616	0.909	0.566	0.884
IntegratedGradients	\times	\times	0.586	0.899	0.593	0.899	0.588	0.903	0.465	0.863
DeepLift	\times	\times	0.575	0.900	0.598	0.898	0.594	0.905	0.542	0.883
GradientShap	\times	\times	0.561	0.893	0.592	0.899	0.600	0.904	0.470	0.858
DeepLiftShap	\times	\times	0.569	0.897	0.607	0.901	0.619	0.909	0.554	0.887
SaliencyNoiseTunnel	\times	\times	0.551	0.892	0.581	0.896	0.578	0.899	0.475	0.851
ShapleySampling	0.456	0.866	0.628	0.910	0.613	0.898	0.655	0.916	0.668	0.917
FeaturePermutation	0.454	0.866	0.624	0.910	0.616	0.903	0.655	0.917	0.677	0.918
FeatureAblation	0.279	0.733	0.438	0.811	0.479	0.824	0.425	0.792	0.408	0.830
Occlusion	0.456	0.866	0.617	0.909	0.609	0.898	0.653	0.917	0.684	0.920
ArchDetect	0.251	0.696	0.369	0.774	0.446	0.818	0.379	0.784	0.382	0.805

the degradation of the model’s performance. The larger the degradation is, the better the estimation is, since it identifies the features most helpful for the model on the task.

More specifically, we use **ROAR** (remove and retrain) proposed in⁷ for evaluation. For each interpretability method, we replace the most important features of certain fractions of each data sample with a fixed uninformative value. We conduct this in both training and test sets. Then we retrain the model with the modified training set and evaluate its classification performance on the modified test set. By retraining the model on datasets with features removed, ROAR ensures that train and test data comes from a similar distribution and reduces the impact on the model’s performance of data distribution discrepancy, so that the degradation of performance is caused by the removal of information instead of the shift of data distribution.

For sequence input $\mathbf{X} \in \mathbb{R}^{T \times F}$, we flatten it and give feature importance scores for all $T \times F$ features. For the i -th feature, we use its mean value in the training set as its uninformative value. We evaluate each interpretability method with feature drop ratios 10%, 20%, ..., 100% and plot the curve of model performance with respect to feature drop ratio for each model.

4.5 Results

4.5.1 Evaluation of Interpretability Methods

Figure 1 shows the curves of model performance (measured with AUPRC and AUROC respectively) with respect to the feature drop ratio of different interpretability methods for the top-2 models (Transformer & IMVLSTM), refer to Section A.3 for all curves. Table 3 gives the quantitative results of area under the curve (AUC). A lower value of AUC means that the performance curve drops faster with the increase of feature drop ratio, thus indicates that the interpretability method gives a better ranking of feature importance.

We have the following observations: (1) **ArchDetect gives the best performing feature importance estimation overall.** From Figure 1, we observe that the curve of ArchDetect drops the fastest for all models on both metrics. Quantitative results in Table 3 also show that ArchDetect has the lowest AUC. Therefore, for the in-hospital mortality task, the feature importance ranking given by ArchDetect is the most reasonable one among results of all interpretability methods considered in this work. (2) **Gradient based methods perform well on LSTM, Transformer and IMVLSTM models, but are no better than a random guess on TCN.** AUC of both metrics of gradient based methods is significantly lower than that of random guessing for LSTM, Transformer and IMVLSTM. But for TCN, even the best performing gradient based method SaliencyNoiseTunnel has AUC close to random guessing (0.581 vs 0.605 for AUPRC and 0.896 vs 0.901 for AUROC). (3) **Attention scores are not necessarily the best estimation of feature importance.** In IMVLSTM, the Glassbox baseline utilizes attention scores the model gives as an estimation of feature importance. Although it outperforms the random guessing baseline, it is not among the best interpretation methods and is inferior to methods such as ArchDetect and IntegratedGradients. Similar observations also exist in the natural language processing domain^{69,70}, where attention weights largely do not correlate with feature importance.

4.5.2 Identified Important Features

We further investigate and compare important features given by different prediction models with the best performing interpretability method ArchDetect in Section 4.5.1 for a qualitative evaluation of its effectiveness. Since ArchDetect gives local feature importance for each data sample respectively, we aggregate local results for a global qualitative evaluation with following steps: (1) for each sample, get the rank of importance for each feature; (2) calculate the average of ranks for each

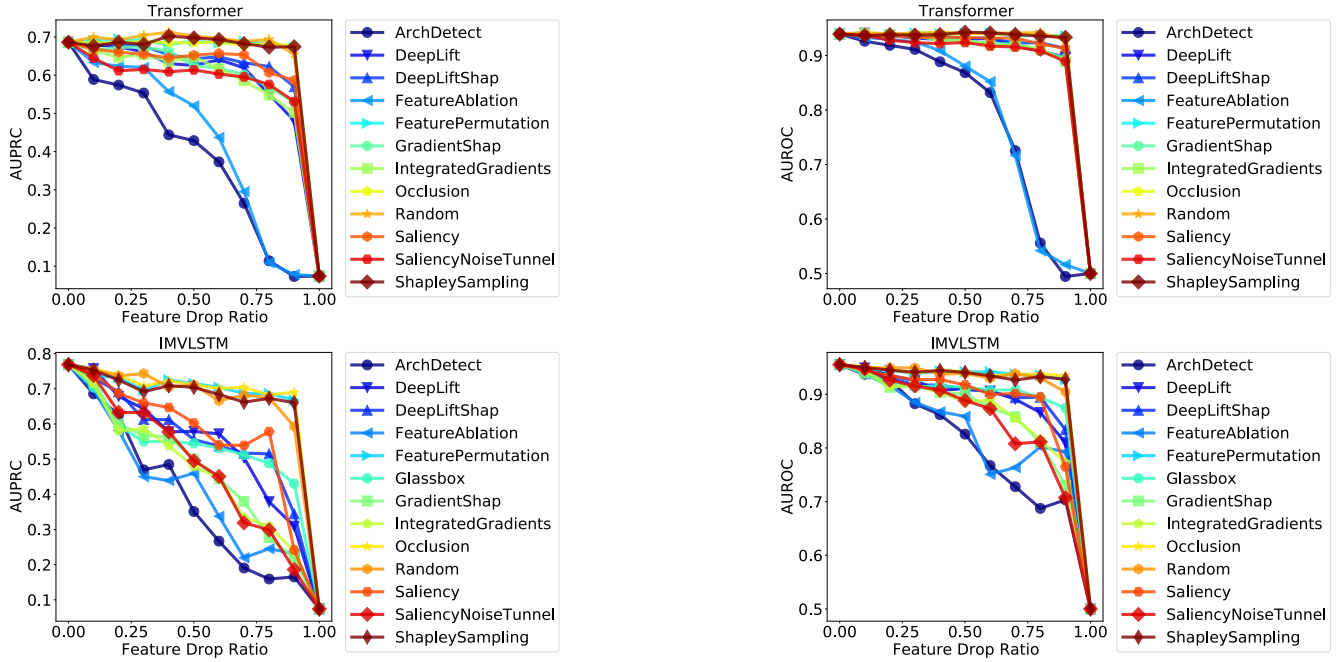


Figure 1. Curves of performance metric w.r.t feature drop ratio.

feature over all data samples; (3) sort the averaged ranks of features from (2) as the global ordering of importance for all features. We then verify the effectiveness of feature importance estimation given by ArchDetect from following aspects:

Similarity of Important Features from Different Models Figure 2 shows the Jaccard similarity of top-50 most important features identified in models. We observe that (1) the Jaccard similarity of top-50 most important features from any pair of two models is above 0.667; (2) each pair of models accepting sequential data (LSTM, TCN, Transformer, and IMVLSTM) has a Jaccard similarity over 0.786. This result demonstrates that ArchDetect identifies similar sets of important features when applied to various models, which is necessary for its correctness since the ground truth set of important features is unique.

Visualization of Global Feature Importance Ranks Figure A2 - Figure A7 visualize the aggregated global feature importance ranks from different models to (1) further demonstrate the similarity of feature importance estimation for different prediction models, and (2) to give an intuitive explanation of what features are important for mortality prediction tasks. We observe that there are several regions of features that all prediction models treat as important features: (1) **labevent features**, including feature 29-46 and feature 108-111, which contains laboratory based measurements of fluid of the patient’s body; (2) **respiratory-related features**, including feature 64 (Respiratory Rate), feature 76 (O2Flow) and feature 103 (respiratory_rate); (3) **SAPS-II features**, including feature 85-98, which are used in the Simplified Acute Physiology Score - II severity classification system; (4) **demographic features**, including feature 122 (age), feature 127 (gender), feature 130 (insurance), feature 132 (marital_status), and feature 133 (ethnicity). We find that the feature importance estimation results successfully identify critical features that are also used in the domain knowledge based SAPS-II system for severity classification, indicating the correctness of results. We also notice that demographic features play important roles in prediction models, which may raise the concern of fairness. We further investigate the fairness of data and models in the following section.

AutoInt	1.000	0.667	0.667	0.724	0.695
LSTM	0.667	1.000	0.786	0.852	0.818
TCN	0.667	0.786	1.000	0.818	0.818
Transformer	0.724	0.852	0.818	1.000	0.852
IMVLSTM	0.695	0.818	0.818	0.852	1.000
	AutoInt	LSTM	TCN	Transformer	IMVLSTM

Figure 2. Jaccard similarity of top-50 most important features identified in all models.

5 Fairness Evaluation

In this section, we first describe the set of demographic features considered as protected attributes. We then investigate the extent of which disparate treatment exists within the MIMIC-IV dataset. Given that the in-hospital mortality predictors can be further utilized in a down-stream decision-making policy, we audit their fairness across various protected attributes.

5.1 Protected Attributes

MIMIC-IV came with a set of demographic features that are helpful for the task of auditing in-hospital mortality predictors for prediction fairness. Protected classes under the Equal Credit Opportunity Act (ECOA) include the following: *age, color,*

Table 4. Protected attributes and subgroups within MIMIC-IV.

Protected Attributes	Groups
Ethnicity	['ASIAN', 'BLACK/AFRICAN AMERICAN', 'HISPANIC/LATINO', 'OTHER', 'WHITE']
Gender	['FEMALE', 'MALE']
Marital Status	['MARRIED', 'SINGLE', 'DIVORCED/WIDOWED']
Age	['<55 YRS', '55-67 YRS', '67-78 YRS', '≥78 YRS']
Insurance	['MEDICAID/MEDICARE', 'PRIVATE']

*marital status, national origin, race, recipient of public assistance, religion, sex*⁷¹. For our task, we consider a subset of such protected classes available within the dataset. To remove uncertainty within our analysis, we further identify and drop examples with unclear attributes, such as 'None', 'Unknown', or 'Unable to obtain'. Table 4 lists the attributes and subgroups used within our analysis. Note that *age* is grouped by quartiles. Refer to Table A2 in the Appendix for more information on each subgroup.

5.2 Fair Treatment Analysis

Disparate treatment is unlawful discrimination in US labor law. Title VII of the United States Civil Rights Act is created to prevent unequal treatment or behavior toward someone because of a protected attribute (e.g. race, gender, or religious beliefs). Although the type and duration of treatment received by patients are determined by multiple factors, analyzing treatment disparities in MIMIC-IV can give us insights in potential biases in treatment received by different groups. Previously, there have been a few works pointing out the racial disparities in end-of-life care between cohorts of black and white patients within MIMIC-III^{72,73}. In a similar spirit, we additionally investigate treatment adoptions and duration across not only ethnicity, but also gender, age, marital status, and insurance type.

5.2.1 Evaluation Method

In MIMIC-IV, 5 categories of mechanical ventilation received by patients have been recorded: HighFlow, InvasiveVent, NonInvasiveVent, Oxygen, and Trach. We first extract the treatment duration and then label the patients with no record as no intervention adoption. If a patient had multiple spans, such as an intubation-extubation-reintubation, then we consider the patient's treatment duration to be the sum of the individual spans.

5.2.2 Results

Figure A9 plots the intervention adoption rate and intervention duration across different protected attributes. We observe that: (1) **There exists disparate treatments, which is most evident across different ethnic groups.** The first column in Figure A9 indicates that on average Black and Hispanic cohorts are less likely to receive ventilation treatments, while also receiving a shorter treatment duration. Similarly, this is also observed across groups split by marital status, where single patients tend to receive shorter and fewer ventilation treatments as opposed to married patients, and similarly with patients with public or private insurances. (2) **There are numerous hidden confounders in analyzing disparate treatment.** The fourth column in Figure A9 indicates more treatments provided to older patients. However, one can imagine that cause of this is medically relevant as the older cohort tends to require more care. Similarly, patients with generous public insurance can more easily afford more treatments. In particular, we note that it is difficult to precisely determine whether the differences in treatment are due to intentional discrimination or differences caused by other confounders. At the current junction, we suspect a close look at causal analysis can help address this problem.

5.3 Fair Prediction Analysis

Fairness in machine learning is a rapidly developing field with numerous definitions and metrics for prediction fairness with respect to two notions: individual and group fairness. For our binary classification task of in-hospital mortality prediction, we consider the group notion where a small number of protected demographic groups G (such as racial groups) is fixed, and we then ask for the classification parity of certain statistics across all of these protected groups.

5.3.1 Fairness Metrics

Most recently, a multitude of statistical measures have been introduced for group fairness, most notable are statistics that ask for the equality of the false positive or negative rates across all groups G (often known as '*equal opportunity*'⁴³) or the equality of classification rates (also known as *statistical parity*). Interestingly, it has been proven that some of the competing definitions and statistics previously proposed are mutually exclusive⁷⁴. Thus, it is impossible to satisfy all of these fairness constraints.

In our case, it is often necessary for mortality assessment algorithms to explicitly consider health-related protected characteristics, especially the age of the patients. Hence, an age-neutral assessment score can systematically overestimate a young person's mortality risk, and can in turn encourage unnecessarily medical interventions. Similarly, enforcing equality of mortality classification rates can likewise lead to discriminatory decision making. Hence, we choose AUC (area under the ROC curve) as our evaluation metrics to audit fairness across subgroups. First, it encompasses both FPR and FNR, which touches on the notion of equalized opportunity and equalized odds. Second, it is robust to class imbalance, which is especially important in

Table 5. Summarized Area under the curve (AUC) performance of the in-hospital mortality predictors evaluated on sets of protected groups. Higher AUC indicates better predictive performance.

Methods	Patient Group	AUC Overall	Minimum AUC over all protected groups	Macro-average AUC over all protected group	AUC for the smallest protected group
AutoInt	All	0.900	0.832	0.897	0.882
LSTM		0.941	0.896	0.939	0.932
TCN		0.937	0.883	0.936	0.948
Transformer		0.941	0.898	0.939	0.953
IMV-LSTM		0.955	0.918	0.954	0.968
AutoInt	HEM, METS	0.795	0.546	0.783	0.546
LSTM		0.842	0.726	0.830	0.777
TCN		0.832	0.696	0.822	0.696
Transformer		0.839	0.778	0.830	0.823
IMV-LSTM		0.884	0.845	0.879	0.862

the task of mortality prediction where mortality rates are $\sim 7\%$. Lastly, AUC is threshold agnostic, which does not necessitate setting a specific threshold for binary prediction that is used across all groups.

5.3.2 Evaluation Method

To evaluate fairness on the MIMIC-IV dataset, we stratify the test set by groups (Table 4), and compute the model’s AUC for each protected group, similarly to⁷⁵. In addition, we also added a stratification for the patient group with the largest common comorbidity, with HEM/METS for patients with lymphoma, leukemia, multiple myeloma, and metastatic cancer. We report (1) AUC(min): minimum AUC over all protected groups, (2) AUC(macro-avg): macro-average over all protected group AUCs and (3) AUC(minority): AUC reported for the smallest protected group in the dataset. Higher AUC is better for all three metrics.

Additionally, as MIMIC-IV is an ongoing data collection effort, we also investigate the relationships between the predictive performance of the mortality predictors and the data distribution with respect to each protected group. It was shown in⁷⁶ that if the risk distributions of protected groups in general differ, such as mortality rates, threshold-based decisions will typically yield error metrics that also differ by group. Hence, we are interested in studying the potential source of the bias/differences in predictive performances from the MIMIC-IV training set.

5.3.3 Results

Figure A10 shows the training data distribution, mortality rates, and testing AUCs across each protected attribute for all patients and patients with HEM/METS, summarized over all five classifiers: AutoInt, LSTM, IMV-LSTM, TCN, and Transformer. Smaller gaps in AUC indicate equality in predictive performances, and larger gaps indicate potential inequalities. Table 5 gives the quantitative results of the area under the curve (AUC). Higher values of AUCs for each of the min, avg, and minority AUC metrics indicate better predictive performance with respect to the protected groups.

We have the following observations: (1) **IMV-LSTM performs the best overall on fairness measure with respect to AUC across different protected groups.** Quantitatively, from Table 5, it is clear that IMV-LSTM has the highest AUC for both overall samples and the subgroups. We see that the minimum AUC for the protected subgroups is highest among the methods considered in this work. This indicates a higher lower bound over all protected attributes. Moreover, the AUC gap for minimum over protected groups is much larger than the next best model, Transformer, for the patient groups with HEM, and METS. (2) **The in-hospital mortality predictors are in general fair, but less so for the subgroup of patients with the comorbidity HEM/METS.** From Figure A10, we observe that the maximum AUC gap across all attributes is at most 0.08, which is smaller than the maximum AUC gap for patients with HEM and METS at 0.11. The difference is more pronounced in the Ethnicity class, but can similarly be observed for other protected classes. In general, we note that all models are quite fair across ethnic groups, with small deviations in gender, and patient’s insurance. Across both sets of patients, we see that all classifiers are in general more accurate for younger patients (<55 years) versus older patients. (3) **There exists a strong correlation between mortality rates and AUCs for each of the protected attributes.** We observe that there is a strong correlation between group mortality rates and group AUC, with Pearson’s $r=-0.922$ and a p -value $< .00001$. This shows that groups with higher mortality rates indicate lower AUC scores. From Figure A10, we also observe that data with imbalanced representation between each subgroup does not impact predictive performance substantively.

6 Interactions between Interpretability and Fairness

Fairness and interpretability are two critical pillars of the recent push for fairness, accountability, and transparency within deep learning. Overall, most interpretability works concern with explaining how the input features impact the final prediction, whether through feature importance or attributions, interactions, and knowledge distillation. Fairness on the other hand

considers fairness metrics, optimization for fairness constraints, and the trade-off between accuracy and fairness. However, to the best of our knowledge, few work attempts to answer the question of how can interpretability help with fairness. What can we learn from our interpretability methods that would indicate either algorithmic bias or representation bias? In this section, we present concrete evidence to establish the initial connection between the two areas, but admittedly leave the fully investigation on the strength of this interaction for future work.

6.1 Feature Importance Correlation with Fairness Metrics

Given mortality predictions made by state-of-the-art models on MIMIC-IV, we study the connections between feature importance induced by different interpretation approaches and the fairness measures in Figure A8. For all the five protected attributes, we compute their respective feature importance by averaging the values produced from interpretability models across time and patients. Taking the feature importance as x axis and the minimum AUC from subgroups split by protected attributes as y axis, we are expecting to see a decreasing trend, where more important features have a higher possibility to lead to performance divergence in the split subgroups. We observe the expected trend consistently among all prediction models, when the interpretability approach *DeepLift* and *DeepLiftShap* are utilized. As shown in Figure A8, age (black dot) is the most important feature compared with other protected attributes and the accuracy difference between young and old is more obvious than other group divisions. Similarly, ethnicity (red dot) and gender (green dot) are the least important features, which leads to much higher minimum AUC than other protected attributes. We plotted but did not observe obvious connections between feature importance from other interpretability approaches and other two fairness evaluation metrics.

6.2 Feature Importance Scores across Protected Attributes

Interpretability often concerns with *global* feature importance for the entire model and *local* feature importance for an individual sample with respect to its prediction. Here, we consider the group feature importance that builds upon *local* feature importance. Ideally, we want to measure how important each feature is across different groups with certain protected attributes. Hence, we define the group feature importance g_i for feature i and protected attribute A , $g_{i,A} = \frac{1}{N_A} \sum_{j=1}^{N_A} \phi_i^j$, where N_A is the size of the group with attribute A , and ϕ_i^j is the *local* feature importance of the feature i for a person j with attribute A . The parity between $g_{i,A}$ would indicate a parity in how each feature is being used for different groups within a certain class of protected attributes. In the MIMIC-IV setting, we are interested in the importance of each of the demographic features used for the in-mortality prediction across the protected subgroups.

Since the scales of the feature importance scores are different for each interpretability method, we calculate the group feature importance for each demographic feature and rank their importance relative to other features within each interpretability method. Since feature importance is provided for {each hour timestep} \times {each feature} within the first 24 hours in the ICU, for all models, we additionally average the feature importance across timesteps. Figure A11 presents the box plot of the feature rankings for each demographic feature for the four models: Transformer, TCN, LSTM, and IMV-LSTM, and each of the 12 interpretability methods: ArchDetect, DeepLiftShap, FeaturePermutation, IntegratedGradients, SaliencyNoiseTunnel, DeepLift, FeatureAblation, GradientShap, Saliency, and ShapleySampling. A lower ranking indicates higher feature importance.

We observe that similar trends exist across different models of varying architectures, where a demographic feature is more important (has lower ranking) for specific groups. Out of 164 features used for each timestep, the feature *ethnicity* has the highest feature importance for the *WHITE* patients, similarly for the *MALE* patients with the feature *gender*, and the age group ≥ 78 YRS with the feature *age*, and so on. The protected attribute age is the most intuitive in this setting, where in-hospital mortality predictors would attribute high importance to elderly patients since that is a strong signal for mortality prediction. A similar case can be made the feature *insurance*, as patients with Medicare are often elderly. However, it is less intuitive for the *ethnicity* feature, as to why one subgroup would use the *ethnicity* feature more strongly than the other subgroups. This stark parity exists for all models, even for different methods of interpretability to obtain feature importance. In summary, we do note that feature importance, especially when viewed as group importance, can concretely reveal how a feature is being used for different groups. However, it is difficult to identify the confounders or features that strongly correlate with the *ethnicity* feature. Therefore we leave further study from causal perspectives for future work.

7 Summary

In this work, we conduct analysis on the MIMIC-IV dataset and several deep learning models in terms of model interpretability, dataset bias, algorithmic fairness, and the interaction between interpretability and fairness. We present quantitative evaluations of interpretability methods on deep learning models for mortality prediction, demonstrate the dataset bias in treatment in MIMIC-IV, verify the fairness of studied mortality prediction models, and reveal the disparities of feature importance among demographic subgroups. We will conduct further analysis from a causal perspective on the relation between the difference of feature importance and the difference of model outcomes among subgroups.

References

1. Purushotham, S., Meng, C., Che, Z. & Liu, Y. Benchmarking deep learning models on large healthcare datasets. *J. Biomed. Informatics* **83**, 112 – 134, DOI: <https://doi.org/10.1016/j.jbi.2018.04.007> (2018).
2. Harutyunyan, H., Khachatrian, H., Kale, D. C., Ver Steeg, G. & Galstyan, A. Multitask learning and benchmarking with clinical time series data. *Sci. data* **6**, 1–18 (2019).
3. Wang, S. *et al.* Mimic-extract: A data extraction, preprocessing, and representation pipeline for mimic-iii. In *Proceedings of the ACM Conference on Health, Inference, and Learning*, 222–235 (2020).
4. Chen, I., Johansson, F. D. & Sontag, D. Why is my classifier discriminatory? In *Advances in Neural Information Processing Systems*, 3539–3550 (2018).
5. Johnson, A. *et al.* Mimic-iv (version 0.4). *PhysioNet* (2020).
6. Goldberger, A. *et al.* Physiobank, physiotoolkit, and physionet: Components of a new research resource for complex physiologic signals. *Circ. [Online]* **101(23)**, e215–e220 (2000).
7. Hooker, S., Erhan, D., Kindermans, P.-J. & Kim, B. A benchmark for interpretability methods in deep neural networks. In *Advances in Neural Information Processing Systems*, 9737–9748 (2019).
8. Tsang, M., Rambhatla, S. & Liu, Y. How does this interaction affect me? interpretable attribution for feature interactions. *Adv. Neural Inf. Process. Syst.* **33** (2020).
9. Guo, T., Lin, T. & Antulov-Fantulin, N. Exploring interpretable lstm neural networks over multi-variable data. In *International Conference on Machine Learning*, 2494–2504 (2019).
10. Simonyan, K., Vedaldi, A. & Zisserman, A. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034* (2013).
11. Sundararajan, M., Taly, A. & Yan, Q. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, 3319–3328 (2017).
12. Shrikumar, A., Greenside, P. & Kundaje, A. Learning important features through propagating activation differences. In *International Conference on Machine Learning*, 3145–3153 (2017).
13. Ancona, M., Ceolini, E., Öztireli, C. & Gross, M. Towards better understanding of gradient-based attribution methods for deep neural networks. In *International Conference on Learning Representations* (2018).
14. Lundberg, S. M. & Lee, S.-I. A unified approach to interpreting model predictions. In *Advances in neural information processing systems*, 4765–4774 (2017).
15. Smilkov, D., Thorat, N., Kim, B., Viégas, F. & Wattenberg, M. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825* (2017).
16. Castro, J., Gómez, D. & Tejada, J. Polynomial calculation of the shapley value based on sampling. *Comput. & Oper. Res.* **36**, 1726–1730 (2009).
17. Strumbelj, E. & Kononenko, I. An efficient explanation of individual classifications using game theory. *The J. Mach. Learn. Res.* **11**, 1–18 (2010).
18. Molnar, C. *Interpretable Machine Learning* (Lulu. com, 2020).
19. Suresh, H. *et al.* Clinical intervention prediction and understanding using deep networks. *arXiv preprint arXiv:1705.08498* (2017).
20. Zeiler, M. D. & Fergus, R. Visualizing and understanding convolutional networks. In *European conference on computer vision*, 818–833 (Springer, 2014).
21. Sundararajan, M., Dhamdhere, K. & Agarwal, A. The shapley taylor interaction index. In *International Conference on Machine Learning*, 9259–9268 (PMLR, 2020).
22. Janizek, J. D., Sturmfels, P. & Lee, S.-I. Explaining explanations: Axiomatic feature interactions for deep networks. *arXiv preprint arXiv:2002.04138* (2020).
23. Sorokina, D., Caruana, R., Riedewald, M. & Fink, D. Detecting statistical interactions with additive groves of trees. In *Proceedings of the 25th international conference on Machine learning*, 1000–1007 (2008).
24. Tsang, M., Cheng, D. & Liu, Y. Detecting statistical interactions from neural network weights. In *International Conference on Learning Representations* (2018).

25. Tsang, M., Liu, H., Purushotham, S., Murali, P. & Liu, Y. Neural interaction transparency (nit): Disentangling learned interactions for improved interpretability. In *Advances in Neural Information Processing Systems*, 5804–5813 (2018).
26. Dhamdhere, K., Sundararajan, M. & Yan, Q. How important is a neuron? *arXiv preprint arXiv:1805.12233* (2018).
27. Shrikumar, A., Su, J. & Kundaje, A. Computationally efficient measures of internal neuron importance. *arXiv preprint arXiv:1807.09946* (2018).
28. Leino, K., Sen, S., Datta, A., Fredrikson, M. & Li, L. Influence-directed explanations for deep convolutional networks. In *2018 IEEE International Test Conference (ITC)*, 1–8 (IEEE, 2018).
29. Springenberg, J., Dosovitskiy, A., Brox, T. & Riedmiller, M. Striving for simplicity: The all convolutional net. In *ICLR (workshop track)* (2015).
30. Kim, B. *et al.* Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International conference on machine learning*, 2668–2677 (PMLR, 2018).
31. Ghorbani, A., Wexler, J., Zou, J. Y. & Kim, B. Towards automatic concept-based explanations. In *Advances in Neural Information Processing Systems*, 9277–9286 (2019).
32. Zhou, B., Sun, Y., Bau, D. & Torralba, A. Interpretable basis decomposition for visual explanation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 119–134 (2018).
33. Ismail, A. A., Gunady, M., Corrada Bravo, H. & Feizi, S. Benchmarking deep learning interpretability in time series predictions. *Adv. Neural Inf. Process. Syst.* **33** (2020).
34. Hardt, M. *et al.* Explaining an increase in predicted risk for clinical alerts. In *Proceedings of the ACM Conference on Health, Inference, and Learning*, 80–89 (2020).
35. Sanchez-Lengeling, B. *et al.* Evaluating attribution for graph neural networks. *Adv. Neural Inf. Process. Syst.* **33** (2020).
36. Samek, W., Binder, A., Montavon, G., Lapuschkin, S. & Müller, K.-R. Evaluating the visualization of what a deep neural network has learned. *IEEE transactions on neural networks learning systems* **28**, 2660–2673 (2016).
37. Lambrecht, A. & Tucker, C. Algorithmic bias? an empirical study of apparent gender-based discrimination in the display of stem career ads. *Manag. Sci.* **65**, 2966–2981 (2019).
38. Raji, I. D. & Buolamwini, J. Actionable auditing: Investigating the impact of publicly naming biased performance results of commercial ai products. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, 429–435 (2019).
39. Schnabel, T., Swaminathan, A., Singh, A., Chandak, N. & Joachims, T. Recommendations as treatments: Debiasing learning and evaluation. *arXiv preprint arXiv:1602.05352* (2016).
40. Dressel, J. & Farid, H. The accuracy, fairness, and limits of predicting recidivism. *Sci. advances* **4**, eaao5580 (2018).
41. Fu, R., Huang, Y. & Singh, P. V. Artificial intelligence and algorithmic bias: Source, detection, mitigation, and implications. In *Pushing the Boundaries: Frontiers in Impactful OR/OM Research*, 39–63 (INFORMS, 2020).
42. Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K. & Galstyan, A. A survey on bias and fairness in machine learning. *arXiv preprint arXiv:1908.09635* (2019).
43. Hardt, M., Price, E. & Srebro, N. Equality of opportunity in supervised learning. In *Advances in neural information processing systems*, 3315–3323 (2016).
44. Bellamy, R. K. *et al.* Ai fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic bias. *arXiv preprint arXiv:1810.01943* (2018).
45. Kamiran, F. & Calders, T. Data preprocessing techniques for classification without discrimination. *Knowl. Inf. Syst.* **33**, 1–33 (2012).
46. Moyer, D., Gao, S., Brekelmans, R., Galstyan, A. & Ver Steeg, G. Invariant representations without adversarial training. *Adv. Neural Inf. Process. Syst.* **31**, 9084–9093 (2018).
47. Singh, H., Singh, R., Mhasawade, V. & Chunara, R. Fair predictors under distribution shift. *arXiv preprint arXiv:1911.00677* (2019).
48. Barda, N. *et al.* Addressing bias in prediction models by improving subpopulation calibration. *J. Am. Med. Informatics Assoc.* (2020).
49. Martinez, N., Bertran, M. & Sapiro, G. Minimax pareto fairness: A multi objective perspective. In *International Conference on Machine Learning*, 6755–6764 (PMLR, 2020).

50. Zhang, H., Lu, A. X., Abdalla, M., McDermott, M. & Ghassemi, M. Hurtful words: quantifying biases in clinical contextual word embeddings. In *Proceedings of the ACM Conference on Health, Inference, and Learning*, 110–120 (2020).
51. Chen, I. Y., Szolovits, P. & Ghassemi, M. Can ai help reduce disparities in general medical and mental health care? *AMA journal ethics* **21**, 167–179 (2019).
52. Cui, S., Pan, W., Zhang, C. & Wang, F. xorder: A model agnostic post-processing framework for achieving ranking fairness while maintaining algorithm utility. *arXiv preprint arXiv:2006.08267* (2020).
53. Chen, J., Berlot-Atwell, I., Hossain, S., Wang, X. & Rudzicz, F. Exploring text specific and blackbox fairness algorithms in multimodal clinical nlp. *arXiv preprint arXiv:2011.09625* (2020).
54. Sharma, S., Henderson, J. & Ghosh, J. Certifai: Counterfactual explanations for robustness, transparency, interpretability, and fairness of artificial intelligence models. *arXiv preprint arXiv:1905.07857* (2019).
55. Chu, E., Gillani, N. & Priscilla Makini, S. Games for fairness and interpretability. In *Companion Proceedings of the Web Conference 2020*, 520–524 (2020).
56. Doshi-Velez, F. & Kim, B. A roadmap for a rigorous science of interpretability. *arXiv preprint arXiv:1702.08608* **2** (2017).
57. Lipton, Z. C. The mythos of model interpretability. *Queue* **16**, 31–57 (2018).
58. Du, M., Yang, F., Zou, N. & Hu, X. Fairness in deep learning: A computational perspective. *IEEE Intell. Syst.* (2020).
59. Kleinberg, J. & Mullainathan, S. Simplicity creates inequity: implications for fairness, stereotypes, and interpretability. In *Proceedings of the 2019 ACM Conference on Economics and Computation*, 807–808 (2019).
60. Jabbari, S., Ou, H.-C., Lakkaraju, H. & Tambe, M. An empirical study of the trade-offs between interpretability and fairness. *ICML 2020 Work. on Hum. Interpret. Mach. Learn.* (2020).
61. Adebayo, J. & Kagal, L. Iterative orthogonal feature projection for diagnosing bias in black-box models. *arXiv preprint arXiv:1611.04967* (2016).
62. Wadsworth, C., Vera, F. & Piech, C. Achieving fairness through adversarial learning: an application to recidivism prediction. *arXiv preprint arXiv:1807.00199* (2018).
63. Cesaro, J. & Cozman, F. G. Measuring unfairness through game-theoretic interpretability. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 253–264 (Springer, 2019).
64. Sjoding, M. *et al.* Democratizing ehr analyses a comprehensive pipeline for learning from clinical data. *Mach. Learn. For Healthc. (Clinical Abstr. Track)* (2019).
65. Song, W. *et al.* Autoint: Automatic feature interaction learning via self-attentive neural networks. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, 1161–1170 (2019).
66. Hochreiter, S. & Schmidhuber, J. Long short-term memory. *Neural computation* **9**, 1735–1780 (1997).
67. Bai, S., Kolter, J. Z. & Koltun, V. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271* (2018).
68. Vaswani, A. *et al.* Attention is all you need. *Adv. neural information processing systems* **30**, 5998–6008 (2017).
69. Jain, S. & Wallace, B. C. Attention is not explanation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 3543–3556 (2019).
70. Grimsley, C., Mayfield, E. & R.S. Bursten, J. Why attention is not explanation: Surgical intervention and causal reasoning about neural models. In *Proceedings of the 12th Language Resources and Evaluation Conference*, 1780–1790 (European Language Resources Association, Marseille, France, 2020).
71. Chen, J., Kallus, N., Mao, X., Svacha, G. & Udell, M. Fairness under unawareness: Assessing disparity when protected class is unobserved. In *Proceedings of the conference on fairness, accountability, and transparency*, 339–348 (2019).
72. Yarnell, C. J. *et al.* Association between immigrant status and end-of-life care in ontario, canada. *JAMA* **318**, 1479–1488 (2017).
73. Lee, J. J., Long, A. C., Curtis, J. R. & Engelberg, R. A. The influence of race/ethnicity and education on family ratings of the quality of dying in the icu. *J. Pain Symptom Manag.* **51**, 9 – 16, DOI: <https://doi.org/10.1016/j.jpainsymman.2015.08.008> (2016).
74. Kleinberg, J., Mullainathan, S. & Raghavan, M. Inherent trade-offs in the fair determination of risk scores. *arXiv preprint arXiv:1609.05807* (2016).

75. Lahoti, P. *et al.* Fairness without demographics through adversarially reweighted learning. In *Advances in neural information processing systems* (2021).
76. Corbett-Davies, S. & Goel, S. The measure and mismeasure of fairness: A critical review of fair machine learning. *arXiv preprint arXiv:1808.00023* (2018).

Author contributions statement

C.M. conceived and conducted the interpretability experiment(s), L.T. and N.X. conducted the fairness experiment(s), C.M., L.T, N.X., and Y.L analysed the results. All authors reviewed the manuscript.

Competing interests

The author(s) declare no competing interests.

Figures

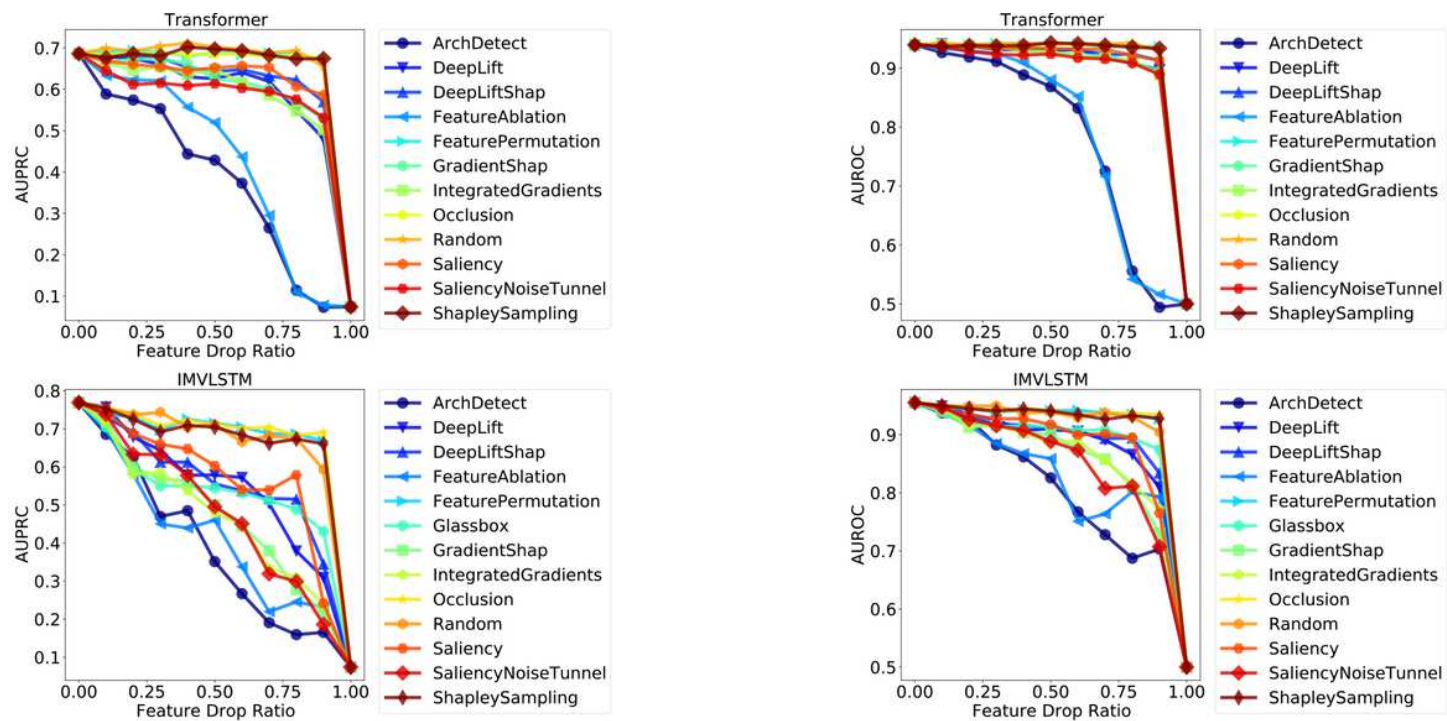


Figure 1

Curves of performance metric w.r.t feature drop ratio.

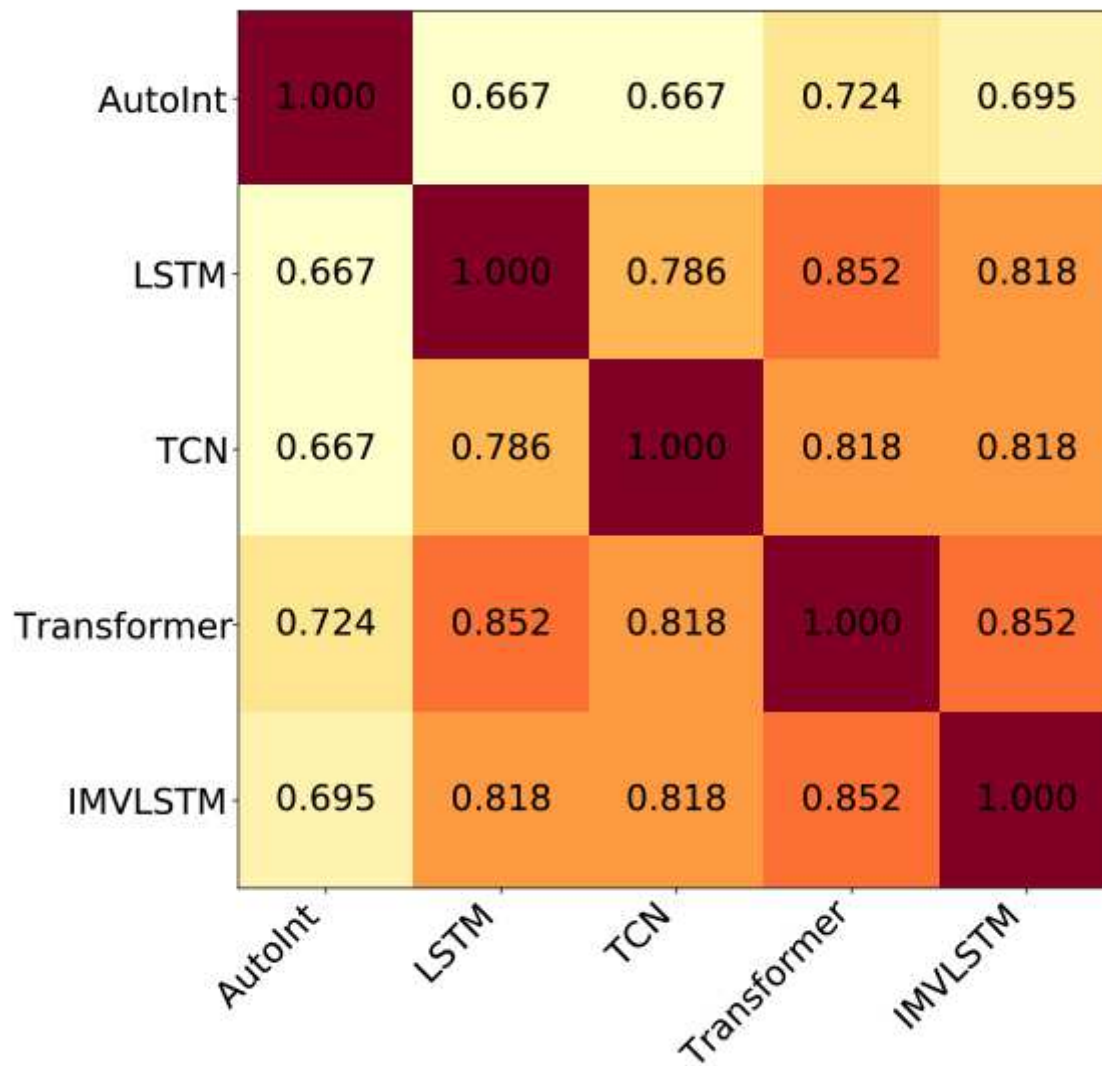


Figure 2

Jaccard similarity of top-50 most important features identified in all models.