

Evaluating Fairness in Healthcare Machine Learning: A Quantitative Approach

K. Roopa

Assistant Professor, Department of Computer Applications, Madanapalle Institute of Technology & Science, Madanpalle

Abstract: *As machine learning models become increasingly integral to healthcare, concerns about their fairness in decision-making processes arise. This paper introduces a robust quantitative methodology to measure fairness in healthcare-oriented machine learning algorithms. By evaluating diverse datasets, we identified notable performance disparities across patient subgroups, such as gender and ethnicity. These findings highlight that even models optimized for accuracy can inadvertently perpetuate systemic biases. To counteract these imbalances, we propose specific mitigation strategies, demonstrating their efficacy in enhancing fairness without compromising overall performance. Our research underscores the importance of ensuring equitable AI applications in healthcare, emphasizing that accuracy and fairness must coexist for the optimal benefit of all patients. While there's broad recognition of the need to address fairness, the healthcare domain lacks a comprehensive quantitative metric to assess and counteract it. This paper introduces a novel quantitative measure designed to evaluate fairness in ML algorithms, emphasizing its applicability to healthcare scenarios. We formulate the metric by grounding it in the intricacies of healthcare data and its multifaceted challenges. Our empirical analysis, conducted on multiple healthcare datasets, showcases the utility of our measure in identifying and mitigating biases. The results underscore the metric's potential in aiding the development of more equitable ML models, ensuring that advancements in healthcare ML are both transformative and just for all patient demographics.*

Keywords: Algorithmic fairness, machine learning, healthcare, quantitative measure, bias, Adversarial Debiasing

1. Introduction

The transformative power of machine learning (ML) in healthcare is undeniable. From predicting patient readmissions to aiding in diagnostics, ML algorithms have demonstrated the potential to revolutionize healthcare outcomes and efficiencies. However, with the increased adoption of these algorithms, concerns about their fairness and potential biases have become paramount. Unlike applications in other sectors, the stakes in healthcare are exceptionally high-biased algorithms could inadvertently prioritize or disadvantage certain patient groups, leading to suboptimal or even harmful medical interventions.

Fairness in machine learning, particularly in healthcare, is not just a computational challenge but a moral imperative. While there's a consensus on the urgency to address fairness, the community lacks a universally accepted quantitative metric to assess and ensure it, especially in the intricate domain of healthcare with its diverse patient populations and multifaceted data sources.

This paper introduces a novel quantitative measure tailored to evaluate fairness in machine learning algorithms, specifically in the context of healthcare. We contend that a nuanced understanding of healthcare data, combined with a rigorous fairness metric, can lead to more equitable algorithms, thereby ensuring that the benefits of ML in healthcare are shared across all demographics. Through this paper, we aim to bridge the gap between the theoretical aspirations of fairness and its practical implementation in healthcare machine learning applications.

2. Background & Related Work

2.1 Background

Machine Learning in Healthcare

Machine Learning (ML) has rapidly permeated the healthcare domain over the last decade. Its applications range from predictive analytics, where algorithms forecast patient health events, to diagnostics, where models assist physicians in identifying diseases from medical images or genomic sequences. The promise of ML in healthcare is to augment human expertise, personalize patient care, and optimize hospital operations.

Fairness in Machine Learning

The concept of fairness in ML transcends mere algorithmic intricacies. It deals with ensuring that the models' predictions do not discriminate against particular groups, especially in contexts where biases can have severe societal ramifications. The absence of fairness can lead to skewed predictions, often disadvantaging already marginalized groups. In healthcare, this is especially concerning given the potential for real-world health disparities based on biased algorithmic recommendations.

2.2 Related Work

Existing Fairness Metrics

Various metrics have been proposed in the realm of ML to quantify fairness, such as demographic parity, equal opportunity, and disparate impact, among others. While these metrics offer valuable insights into specific aspects of fairness, they often fall short in capturing the multifaceted nature of healthcare data. For instance, Hardt et al. (2016) introduced the notion of equal opportunity, which ensures equal false positive rates across protected groups, but its applicability can be limited in complex healthcare scenarios.

Fairness in Healthcare ML

A growing body of research has emerged emphasizing the importance of fairness in healthcare ML. Obermeyer et al. (2019) demonstrated how a widely - used healthcare algorithm manifested racial bias, proving the vital need for fairness measures tailored to the domain. Several efforts have been made to adapt general fairness principles to healthcare. However, these adaptations often require a deeper understanding of domain - specific intricacies, emphasizing the need for a dedicated metric.

Mitigating Biases in ML Models

Techniques such as re - sampling, re - weighting, and adversarial training have been employed to mitigate biases in training data or model predictions. For healthcare, Zafar et al. (2017) presented a method to ensure fairness constraints in clinical prediction models, emphasizing the delicate balance between fairness and model utility.

3. Quantitative Measure for Fairness: Definition & Methodology

3.1 Definition: Healthcare Equitable Impact Score (HEIS)

The Healthcare Equitable Impact Score (HEIS) quantifies the extent to which a machine learning model's predictions are equitable across different protected groups within healthcare contexts. The HEIS is defined as a value between 0 and 1, where a score of 1 indicates perfect fairness, and a score closer to 0 suggests significant disparities in model outcomes across the evaluated groups.

Mathematically, given a set of protected groups $G = \{g_1, g_2, \dots, g_n\}$ and a set of outcomes $O = \{o_1, o_2, \dots, o_m\}$, the HEIS is defined as:

$$HEIS = 1 - \frac{\sum_{i=1}^n \sum_{j=1}^m |P(o_j | g_i) - P(o_j)|}{n \times m}$$

Where:

$P(o_j | g_i)$ is the probability of outcome o_j given protected group g_i .

$P(o_j)$ is the overall probability of outcome o_j across all groups.

3.2 Methodology

Data Partitioning

- Protected Groups Identification:** Identify and categorize protected groups within the dataset (e. g., based on race, gender, age groups).
- Outcome Analysis:** Analyze different outcomes or predictions made by the model (e. g., likelihood of disease, readmission risk).

HEIS Calculation

- Individual Group Probabilities:** For each protected group g_i , compute the probability of each outcome $P(o_j | g_i)$.
- Overall Probabilities:** Calculate the overall probability of each outcome $P(o_j)$ across the entire dataset.
- Disparity Computation:** For each combination of protected group and outcome, compute the absolute

difference between the group - specific outcome probability and the overall outcome probability.

- Aggregation:** Sum all the disparities and normalize by the number of groups and outcomes to get the final HEIS.

Interpretation

- Thresholding:** While a perfect score of 1 indicates absolute fairness, in real - world scenarios, a threshold (e. g., 0.95) may be set as the acceptable fairness limit.
- Contextual Analysis:** It's essential to analyze HEIS in the context of the specific healthcare application. For example, certain disparities may arise due to legitimate medical reasons and not due to model bias.

Adjustments and Iterations

Based on the HEIS, one can iteratively adjust the ML model, employ debiasing techniques, and recalibrate to enhance fairness while maintaining accuracy.

4. Dataset & Preprocessing

4.1 Dataset Description

Name: Health Equity Dataset (HED)

Source: This dataset was sourced from multiple hospitals and healthcare institutions across the country, ensuring a diverse representation of patient populations.

Size: The dataset comprises 500,000 patient records.

Features:

Demographic Data: Age, Gender, Race, Socioeconomic Status, and Zip Code.

Medical History: Past Diagnoses, Medications, Allergies, Surgical History.

Recent Health Metrics: Vital Signs, Lab Results, Imaging Results.

Outcome Variable: The likelihood of readmission within 30 days.

4.2 Preprocessing Steps

1) Data Cleaning:

Missing Values: Rows with missing outcome variables were dropped. For other missing data, we employed k - nearest neighbors imputation to predict and fill in the missing values.

Outliers: Extreme outliers, identified using the IQR method, were addressed to prevent skewness.

2) Feature Engineering:

Aggregation: Certain features, like vital signs taken over time, were aggregated to create meaningful metrics, such as average blood pressure in the last year.

One - Hot Encoding: Categorical variables like race and gender were one - hot encoded to convert them into binary columns, making them suitable for ML models.

3) Data Splitting:

The dataset was split into a training set (80%) and a test set (20%). The splitting ensured that the distribution of the outcome variable and key demographic features was consistent across both sets.

4) Handling Class Imbalance:

Given that the number of readmissions is typically lower than non-readmissions, Synthetic Minority Over-sampling Technique (SMOTE) was employed to balance the classes in the training set.

5) Normalization:

Continuous features were normalized using Z-score normalization to ensure they're on the same scale, aiding in convergence during model training.

4.3 Fairness Considerations in Data

- **Representation:** The dataset was examined to ensure that minority groups were adequately represented. This is essential for fairness considerations, as underrepresented groups can lead to models that are less accurate for those groups.
- **Bias Analysis:** Initial exploratory data analysis was conducted to check for any inherent biases in the dataset, especially concerning the outcome variable. Patterns indicating potential bias were noted for further exploration during model evaluation.

5. Experimental Setup

5.1 Objectives

The main objectives of our experiments are:

- To evaluate the fairness of different machine learning algorithms using the proposed Healthcare Equitable Impact Score (HEIS).
- To compare the performance of these algorithms in terms of traditional metrics like accuracy, recall, and precision.

5.2 Algorithms Evaluated

For our experiments, we considered a mix of both classical and deep learning models, known for their widespread use in healthcare applications:

- 1) **Logistic Regression (LR):** A baseline model, known for its simplicity and interpretability.
- 2) **Decision Trees (DT):** Chosen for its non-linear decision-making capability.
- 3) **Random Forest (RF):** An ensemble method for increased accuracy and robustness.
- 4) **Gradient Boosting Machines (XGBoost):** Known for high performance in structured data tasks.
- 5) **Convolutional Neural Networks (CNN):** Incorporated to process imaging data embedded within our dataset.
- 6) **Recurrent Neural Networks (RNN):** Used to process sequential data like patient health metrics over time.

5.3 Training Configuration

- **Epochs:** For deep learning models, we trained for 50 epochs with early stopping to prevent overfitting.
- **Batch Size:** 128 for deep learning models.
- **Optimizer:** Adam optimizer with a learning rate of 0.001 for deep learning models.
- **Regularization:** L2 regularization was applied to prevent overfitting across all models.

5.4 Fairness Enhancing Interventions

Given our focus on fairness, we also evaluated variations of the above algorithms incorporating fairness-enhancing interventions:

- **Adversarial Debiasing:** Trained models in an adversarial setting where a secondary network tries to predict the protected attribute from the model's predictions.
- **Pre-processing Techniques:** Re-sampling and re-weighting strategies were explored, especially for underrepresented groups.
- **Fairness Constraints:** Integrated fairness constraints during model optimization to ensure equitable predictions.

5.5 Evaluation Metrics

- a) **Primary Metric:** Healthcare Equitable Impact Score (HEIS) – our proposed metric for fairness.
- b) **Secondary Metrics:**
 - **Accuracy:** The proportion of correctly classified instances.
 - **Precision, Recall, and F1 - Score:** Especially crucial given the potential class imbalance in readmission datasets.
 - **AUC - ROC:** Useful for understanding the trade-off between sensitivity and specificity.

5.6 Experimental Environment

- **Hardware:** Experiments were conducted on a workstation with an Intel i9 processor, 64GB RAM, and NVIDIA RTX 3090 GPU.
- **Software:** Models were implemented using TensorFlow 2.x and Scikit-learn libraries in a Python 3.8 environment.

6. Results

In this section, we elucidate the performance and fairness of different algorithms on the Health Equity Dataset (HED), assessed through our experimental setup.

6.1 Performance Metrics

The following table summarizes the performance metrics for each evaluated algorithm:

Algorithm	Accuracy	Precision	Recall	F1 - Score	AUC - ROC
LR	0.89	0.85	0.80	0.82	0.92
DT	0.87	0.83	0.79	0.81	0.90
RF	0.92	0.89	0.88	0.88	0.95
XGBoost	0.93	0.91	0.89	0.90	0.96
CNN	0.90	0.87	0.85	0.86	0.93
RNN	0.91	0.88	0.86	0.87	0.94

Gradient Boosting Machines (XGBoost) achieved the highest accuracy and AUC - ROC. However, high performance does not necessarily equate to fairness, which led us to evaluate models through the lens of the Healthcare Equitable Impact Score (HEIS).

6.2 Healthcare Equitable Impact Score (HEIS) Analysis

The HEIS for each model is presented below:

Algorithm	HEIS
LR	0.82
DT	0.79
RF	0.85
XGBoost	0.83
CNN	0.81
RNN	0.84

Random Forest exhibited the highest HEIS, indicating that among the evaluated models, it was the most equitable across different demographic groups. Notably, despite XGBoost's performance lead, its HEIS was marginally lower than that of the Random Forest.

6.3 Fairness Enhancing Interventions

Employing adversarial debiasing on the Random Forest model boosted its HEIS to 0.90 while retaining an accuracy of 0.91. Pre - processing techniques, especially re - sampling, improved the HEIS for Decision Trees and Logistic Regression by approximately 3%. Fairness constraints during model optimization proved beneficial for XGBoost, raising its HEIS by 4% but at a minor performance trade - off, reducing accuracy by 2%.

6.4 Dissecting Disparities

A deeper analysis revealed that initial disparities were most pronounced for older age groups and certain racial categories. With fairness interventions, these disparities reduced significantly, as evidenced by the improved HEIS across models.

7. Discussion

7.1 Interpretation of Results

Our experiments have reaffirmed the dual challenge of ensuring both high performance and fairness in machine learning models designed for healthcare applications.

- **Performance vs. Fairness:** While XGBoost demonstrated superior performance in traditional metrics, it was not the most equitable model. This highlights the trade - off that sometimes exists between optimizing for accuracy and ensuring fairness. However, the success of Random Forest in achieving a balance between the two demonstrates that they're not mutually exclusive objectives.
- **Fairness Interventions:** The positive impact of fairness - enhancing interventions, especially adversarial debiasing and fairness constraints, underscores the value of these techniques in practice. However, their varying degrees of success across different models suggests that there's no one - size - fits - all approach to debiasing.

7.2 Implications for Healthcare

- **Ethical Decision - making:** Achieving high HEIS values in models is not just a technical triumph but an ethical necessity. As ML models increasingly influence clinical

decisions, their fairness or lack thereof can have real - world consequences, potentially reinforcing existing health disparities.

- **Trust in AI:** For healthcare professionals to trust and adopt AI tools, these tools need to be both accurate and fair. Our research can serve as a roadmap for developing models that fulfill both criteria, promoting wider acceptance and integration of ML in healthcare.

7.3 Limitations & Challenges

- **Dataset Representativeness:** While the Health Equity Dataset (HED) is diverse, it may not capture all nuances of global patient populations. Results might vary on datasets from other regions or demographics.
- **Model Generalization:** The fairness measures and interventions evaluated in this study were specific to the models and dataset at hand. Their applicability and effectiveness might differ in other contexts or with newer algorithms.

7.4 Future Directions

- **Expanding Fairness Metrics:** While the HEIS offers a holistic view of fairness, future work can explore more granular metrics targeting specific aspects of fairness, such as equality of opportunity or treatment.
- **Cross - Domain Validation:** It would be valuable to validate our findings across different healthcare domains, such as diagnostics, treatment recommendation, or patient management, to ensure the universality of our conclusions.
- **Active Learning & Feedback Loops:** Integrating real - world feedback loops where clinicians validate or correct model predictions can be a promising avenue to iteratively improve both model accuracy and fairness.

8. Limitations & Future Work

8.1 Limitations

- **Dataset Constraints:** Our study leveraged the Health Equity Dataset (HED), which, while diverse, may not represent global patient demographics. Variations in datasets from different regions or institutions can influence model fairness and performance.
- **Model Specificity:** The fairness interventions and measures evaluated were tailored to specific algorithms. Their effectiveness may vary with different or newer models, implying the need for continuous assessment as the ML landscape evolves.
- **Bias Blindspots:** While we endeavored to address algorithmic biases, inherent biases in the data—stemming from historical or systemic disparities—might still influence model predictions. Our approach mitigates but may not entirely eliminate such deeply rooted biases.

9. Conclusion

The integration of machine learning in healthcare holds unparalleled promise for improving patient outcomes, operational efficiencies, and medical discoveries. However,

as our study emphasizes, the journey from algorithmic potential to real - world impact must be navigated with fairness at its core. Through the introduction of the Healthcare Equitable Impact Score (HEIS), we provided a novel quantitative measure to assess and ensure fairness in healthcare ML applications.

Our findings underscore a dual imperative: achieving high algorithmic performance while ensuring that predictions are equitable across diverse patient groups. While certain algorithms showcased commendable performance, they sometimes fell short on the fairness scale, underscoring the delicate balance researchers and practitioners must strike.

As ML continues to play an ever - growing role in healthcare decisions, ensuring its fairness becomes not just a computational challenge, but an ethical and societal one. Our research serves as a foundational step in this direction, offering tools, insights, and methodologies to build ML models that are both accurate and just.

Beyond the metrics and models, our study underscores a broader message: In the realm of healthcare, where stakes are inherently high, we must prioritize fairness to ensure that the benefits of AI - driven innovations are accessible and equitable for all.

References

- [1] Hardt, M., Price, E., & Srebro, N. (2016). Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems* (pp.3315 - 3323).
- [2] Obermeyer, Z., Powers, B., Vogeli, C., & Mullainathan, S. (2019). Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366 (6464), 447 - 453.
- [3] Zafar, M. B., Valera, I., Rodriguez, M. G., & Gummadi, K. P. (2017). Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *Proceedings of the 26th International Conference on World Wide Web* (pp.1171 - 1180).
- [4] Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic minority over - sampling technique. *Journal of Artificial Intelligence Research*, 16, 321 - 357.
- [5] Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR*.
- [6] Buolamwini, J., & Gebru, T. (2018). Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on Fairness, Accountability and Transparency* (pp.77 - 91).
- [7] Doshi - Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. arXiv preprint arXiv: 1702.08608.