# Bias Mitigation in ML Algorithms for Chronic Disease Detection

Comparative Analysis of IBM's AIF360

CS 5100 | Foundations of AI
Final Project | Fall 2024

Zadie Moon, Sibgha Ahmad, Yang Zhao

# Introduction

Bias in machine learning (ML) algorithms poses significant challenges, particularly in high-stakes domains such as healthcare and chronic disease detection. Inaccurate or biased predictions in these contexts can lead to disparities in patient outcomes, perpetuating health inequities and compromising the trustworthiness of AI-driven solutions. For instance, studies have demonstrated that models trained on unbalanced datasets often exhibit biases against certain demographic groups, leading to unequal access to accurate diagnoses and treatments (Mehrabi et al., 2021). Given the critical nature of healthcare decisions, ensuring fairness in ML algorithms is not only an ethical imperative but also a practical necessity for delivering equitable and effective care.

Chronic diseases, which account for the majority of deaths worldwide, disproportionately affect underserved populations due to social determinants of health such as access to care, socioeconomic status, and geographic location (World Health Organization, 2023). Incorporating ML into chronic disease prevention and detection has immense potential to enhance early intervention and improve outcomes. However, if ML models reinforce existing biases in healthcare data, they risk exacerbating health disparities rather than addressing them (Obermeyer et al., 2019). This highlights the urgent need for robust bias mitigation strategies that can promote fairness while maintaining or improving the accuracy of ML algorithms.

This paper explores the application of IBM's AI Fairness 360 (AIF360) toolkit in mitigating bias within a machine learning model for chronic disease detection, specifically focusing on heart disease prediction. AIF360 offers a comprehensive framework for detecting and addressing bias throughout the ML pipeline, providing tools for pre-processing, in-processing, and post-processing interventions. By comparing the results of a naive ML implementation with those adjusted using AIF360, this study aims to evaluate the impact of bias mitigation techniques on model fairness and accuracy and emphasize the importance of reducing bias in high-accuracy domains to ensure equitable implementation of life-saving technologies. By addressing these challenges, this work seeks to advance the development of ethical and effective AI tools that can truly transform healthcare for all populations.

# Motivation

Machine learning (ML) holds immense promise for revolutionizing healthcare by enabling early diagnosis, personalized treatment, and predictive modeling. However, the integration of ML into healthcare systems is not without challenges, particularly when it comes to addressing health disparities. Social determinants of health (SDOH)—the conditions in which people are born, grow, live, work, and age—play a critical role in shaping health outcomes (Marmot & Wilkinson, 2005). Factors such as socioeconomic status, education, access to healthcare, and geographic location influence not only individual health but also the quality and distribution of data used to train ML algorithms.

Healthcare datasets often reflect systemic inequities inherent in society. For instance, underrepresentation of certain demographic groups in medical research datasets can lead to biased ML models that disproportionately favor more privileged populations (Rajkomar et al., 2018). This bias in ML models perpetuates health disparities by limiting access to accurate diagnoses and effective treatments for marginalized groups. In high-stakes areas like chronic disease prevention—where early intervention can mean the difference between life and death—these biases can have profound and lasting consequences.

Chronic diseases such as heart disease, diabetes, and stroke are major public health challenges, disproportionately affecting low-income and minority populations (Braveman et al., 2011). These disparities are exacerbated by limited access to preventive care, lower health literacy, and environmental factors such as food deserts and exposure to pollutants. ML algorithms that fail to account for these SDOH risk providing inequitable recommendations, perpetuating cycles of poor health in underserved communities. The motivation for this study stems from the ethical imperative to address these inequities. By examining how bias manifests in ML models used for chronic disease detection and exploring methods to mitigate this bias, we aim to ensure that AI-driven healthcare solutions are equitable and inclusive.

# Literature Review

One notable approach to addressing algorithmic bias is IBM's AI Fairness 360 (AIF360) toolkit. This open-source library demonstrates how bias can be systematically detected and mitigated throughout the machine-learning pipeline.

AIF360 structures bias detection through standardized dataset handling and comprehensive metrics. The StandardDataset class enforces consistent data formatting and explicitly tracks protected attributes and privileged groups - an essential foundation for any bias detection system. This standardization allows for reliable computation of various fairness metrics, including disparate impact, statistical parity, and equal opportunity difference.

The toolkit's bias mitigation strategy operates at three critical stages: pre-processing, in-processing, and post-processing. Pre-processing techniques transform training data to remove inherent biases before model training begins. It uses techniques like Reweighing and Optimized Preprocessing. For example, Reweighing assigns weights to training examples to ensure equal representation across protected groups while maintaining the relationship between features and labels. In-processing methods incorporate fairness constraints directly into the model training process. The Prejudice Remover technique adds a regularization term to the learning objective. Post-processing approaches adjust model outputs to ensure fairness criteria are met. This multi-stage approach demonstrates the importance of considering bias throughout the entire machine learning lifecycle using techniques like Calibrated Equalized Odds. This approach learns threshold classifiers specific to each protected group to satisfy fairness constraints while minimizing accuracy loss.

The toolkit measures bias using metrics like disparate impact which is the ratio of favorable outcomes between unprivileged and privileged groups, equal opportunity difference which measures the difference in true positive rates between groups, and statistical parity which is the difference in selection rates between groups. These metrics are computed using confusion matrix elements specific to each protected group, allowing for quantitative bias assessment at each pipeline stage.

From an implementation perspective, AIF360 utilizes a modular design pattern. Each component - whether a dataset, metric, or mitigation algorithm - is implemented as a separate Python class with clear interfaces. This modularity allows for easy integration of new components and facilitates experimentation with different combinations of metrics and mitigation strategies. The toolkit's integration with popular machine learning frameworks like scikit-learn further showcases the importance of compatibility with existing ML ecosystems.

AIF360's approach to quantifying and visualizing bias metrics provides valuable insights for bias detection tool development. The toolkit implements a wide range of numerical measures and visualization capabilities, enabling detailed analysis of bias at both the dataset and model prediction levels. This comprehensive measurement approach helps identify specific areas where bias manifests and evaluate the effectiveness of mitigation strategies.

Building upon the insights from AIF360, another significant contribution to algorithmic fairness comes from Microsoft's Fairlearn toolkit. This open-source library approaches bias detection and mitigation through gradient-based optimization and constraint satisfaction.

Fairlearn structures bias detection through its MetricFrame class, which provides a foundation for computing and comparing fairness metrics across different subgroups. Similar to AIF360's StandardDataset, this framework enables disaggregated metric assessment, allowing developers to evaluate model performance across various demographic intersections. This approach facilitates the identification of disparities in model behavior that might not be apparent from aggregate statistics alone.

The toolkit's bias mitigation strategy focuses primarily on two approaches: constraints-based optimization and post-processing. While AIF360 operates at three stages, Fairlearn emphasizes the integration of fairness constraints directly into the optimization process. The toolkit employs the Exponentiated Gradient method, which iteratively adjusts model weights to satisfy fairness constraints while maximizing prediction accuracy. The Grid Search mechanism identifies optimal trade-offs between model performance and fairness objectives. Like AIF360's post-processing stage, Fairlearn implements ThresholdOptimizer to adjust decision boundaries for different groups to achieve fairness criteria.

Complementing these approaches, Themis-ML offers yet another perspective on bias detection and mitigation, with a particular emphasis on preprocessing techniques and statistical parity measures. The toolkit implements several preprocessing transformations that address discrimination at the data level, similar to AIF360's preprocessing stage but with different methodological approaches.

Themis-ML's bias mitigation strategy primarily operates through data transformation techniques such as Massaging, Reweighing, and Sampling methods. The Massaging technique, for instance, identifies instances near the decision boundary and modifies their labels to achieve statistical parity while preserving maximum possible accuracy. This approach differs from AIF360's Reweighing by focusing on label modification rather than instance weighting.

Both toolkits extend the measurement frameworks established by AIF360. Fairlearn implements metrics including demographic parity difference, equalized odds difference, and bounded group loss, providing additional perspectives on model fairness. Themis-ML focuses on statistical parity and disparate impact measures, computing discrimination scores based on outcome probability differences between privileged and unprivileged groups.

From an implementation perspective, both Fairlearn and Themis-ML adopt scikit-learn compatible APIs, following AIF360's emphasis on ecosystem integration. This design choice ensures seamless integration with existing machine learning workflows while maintaining familiar interfaces for practitioners. Like AIF360, both toolkits implement modular architectures that allow for independent use of their components.

The combined insights from these three toolkits - AIF360, Fairlearn, and Themis-ML - provide a comprehensive view of bias detection and mitigation strategies. While AIF360 offers a full-pipeline approach with extensive preprocessing options, Fairlearn contributes sophisticated optimization-based techniques, and Themis-ML adds specialized preprocessing transformations. Together, they demonstrate the importance of addressing algorithmic bias through multiple complementary approaches, each offering unique strengths and methodological perspectives.

These implementations collectively inform key design considerations for bias detection tools, including the value of integrating fairness constraints into model optimization, the importance of preprocessing techniques, and the benefits of maintaining compatibility with popular machine

learning frameworks. Their diverse approaches contribute valuable insights to the field of fair machine learning, highlighting the need for multiple strategies in addressing algorithmic bias effectively.

# Description of the dataset

The dataset utilized in this study is sourced from the Kaggle's **Personal Key Indicators of Heart Disease** dataset, which is designed to analyze key risk factors associated with heart disease and stroke. The dataset comprises **319,795 observations** and includes **18 features capturing demographic, behavioral, and clinical data**. This dataset provides a comprehensive view of cardiovascular health, making it ideal for building machine learning models aimed at heart disease detection and prevention.

## Key Features:

**Demographic Attributes:**

- **AgeCategory:** Categorical variable representing the participant's age group (e.g., 55-59, 80 or older).
- **Sex:** Binary classification (Male/Female).
- **Race:** Categorical variable representing the participant's race/ethnicity (e.g., White, Hispanic).

**Behavioral Factors:**

- **Smoking:** Whether the participant is a current smoker (Yes/No).
- **AlcoholDrinking**: Whether the participant drinks alcohol (Yes/No).
- **PhysicalActivity:** Whether the participant engages in physical activity outside of regular work (Yes/No).
- **DiffWalking:** Indicates whether the participant has difficulty walking (Yes/No).

**Clinical and Health Indicators:**

- **BMI:** Body Mass Index, a continuous variable indicating obesity and weight-related health risks.
- **PhysicalHealth:** Number of days in the past month the participant felt physically unwell.
- **MentalHealth:** Number of days in the past month the participant felt mentally unwell.
- **Stroke:** Whether the participant has had a stroke (Yes/No).
- **Diabetic:** Whether the participant has been diagnosed with diabetes (Yes/No).
- **GenHealth:** Self-reported general health (e.g., Very good, Good).
- **SleepTime:** Number of hours of sleep the participant gets on average per day.
- **Asthma, KidneyDisease, SkinCancer:** Binary indicators for whether the participant has asthma, kidney disease, or skin cancer.

### Analysis Focus

This dataset provides insights into various factors contributing to heart disease, including behavioral habits (e.g., smoking, alcohol consumption), clinical indicators (e.g., BMI, stroke history), and social determinants of health (e.g., race, age, physical activity). The large sample size and diverse features allow for a comprehensive analysis of cardiovascular disease risk factors across different population groups.

### Relevance to Study Goals

This dataset provides a comprehensive view of individual and population-level risk factors for heart disease and stroke. The inclusion of diverse demographic and socioeconomic indicators enables the exploration of biases in machine learning models. These biases often arise from the underrepresentation of certain groups or the unequal distribution of SDOH variables, making the dataset ideal for applying IBM's AIF360 fairness toolkit to detect and mitigate disparities.

### Preprocessing Considerations

Before analysis, the dataset will undergo preprocessing steps to address missing values, ensure consistency in feature representation, and define protected attributes (e.g., race, sex) and privileged/unprivileged groups for bias detection. This preparation ensures reliable evaluation and comparison of fairness metrics across models.

## Methods

This study aims to evaluate the effectiveness of bias detection and mitigation techniques in machine learning (ML) models for chronic disease prediction. Our methodology involves implementing a support vector machine (SVM), Neural Network (NN) and other ML models commonly used in healthcare prediction tasks, followed by applying IBM's AI Fairness 360 (AIF360) toolkit to assess and mitigate biases in model predictions.

### 1. Data Preprocessing

The dataset used, **Heart Disease and Stroke Prevention**, comprising demographic, behavioral, and clinical data. Preprocessing steps included:

- Handling missing values using mean/mode imputation techniques (Little & Rubin, 2019).
- Encoding categorical variables with one-hot encoding to ensure compatibility with ML algorithms (Pedregosa et al., 2011).
- Normalizing continuous variables, such as BMI and age, to enhance model convergence.
- Downsample imbalanced data.
- Defining protected attributes (*e.g., race, sex*) and privileged/unprivileged groups for bias detection.

### 2. Model Implementation

We implemented multiple ML models, including:

- **Logistic Regression**: A baseline model for comparison due to its simplicity and interpretability.
- **Decision Tree:** A tree-based model that splits data into branches based on feature a, offering interpretable rules for classification.
- **Random Forest**: Leveraged for its ensemble learning capabilities and robustness to overfitting.
- **Gradient Boosting:** An iterative ensemble method that builds models sequentially, optimizing errors of previous models to improve overall performance.
- **Support Vector Machine (SVM)**: Selected for its ability to handle high-dimensional data and binary classification tasks effectively (Cortes & Vapnik, 1995).
- **Neural Network:** Included to explore its capability to model complex, non-linear relationships in the data.

Each model was trained using an 80/20 train-test split. Performance metrics, including accuracy, precision, recall, and F1-score, confusion matrix were calculated using 5-fold cross validation to evaluate model effectiveness.

## 3. Bias Detection with AIF360

Using AIF360, we assessed bias in the models by calculating fairness metrics, including:

- **Disparate Impact**: Measures the ratio of favorable outcomes between privileged and unprivileged groups.
- **Equal Opportunity Difference**: Assesses the difference in true positive rates across groups.
- **Statistical Parity Difference**: Quantifies the difference in selection rates between groups (Bellamy et al., 2018).

These metrics were computed to identify the extent of bias in model predictions, focusing on disparities across race, gender, and income levels.

## 4. Bias Mitigation

Bias mitigation techniques were applied at various stages of the ML pipeline:

- **Pre-processing**: Techniques such as reweighing were used to adjust training data to achieve fairer representations across protected attributes (Feldman et al., 2015).
- **In-processing**: The Prejudice Remover algorithm was applied to incorporate fairness constraints during model training.
- **Post-processing**: Calibrated Equalized Odds was used to adjust model outputs to align with fairness criteria (Hardt et al., 2016).

# 5. Comparative Analysis

The results of the naive ML models and the bias-mitigated models were compared to evaluate the impact of mitigation techniques. Performance metrics, fairness metrics, and visualization tools provided by AIF360 were analyzed to assess trade-offs between accuracy and fairness.

To evaluate the impact of bias mitigation techniques, a comprehensive comparison was performed between naive machine learning models and their bias-mitigated counterparts. The analysis was structured around three core dimensions: model performance, fairness metrics, and the trade-offs between accuracy and fairness.

## 5.1 Model Performance Comparison

Three different machine learning algorithms—Neural Networks (NN), Singular Value Decomposition (SVD), and Logistic Regression—were employed to evaluate the dataset under varying configurations.

## 5.2 Bias mitigation techniques

Bias mitigation was applied using techniques such as reweighing, preprocessing transformations, and adversarial debiasing. These were evaluated for their ability to balance the trade-off between performance and fairness.

## 5.3 Performance Metrics and Fairness Metrics

Each model's performance was assessed based on:

- Performance Metrics: Accuracy, precision, recall, and F1-score.
- Fairness Metrics: Disparate impact, equal opportunity difference, and demographic parity.

## 5.4 Anticipated Trade-Offs

Based on the methodology, we anticipate the following key trade-offs:

Accuracy vs. Fairness: Models with applied bias mitigation may show slight decreases in accuracy while improving fairness metrics. The magnitude of these changes will depend on the specific technique used and the algorithm employed.

Model-Specific Trends: Neural Networks might demonstrate high accuracy but require significant tuning to achieve fairness. Logistic Regression, being simpler, might provide a more balanced trade-off between accuracy and fairness with minimal parameter adjustments. SVD is expected to reduce dimensionality effectively but will need evaluation to see how well it integrates fairness constraints.

## 5.5 Tools and Visualization Plans

**AIF360 Toolkit:** The fairness metrics and bias mitigation outcomes will be analyzed using tools from AIF360.

**Visualizations:** Scatter plots will be used to explore accuracy vs. fairness trade-offs, and heatmaps will visualize the impact of hyperparameter tuning on both metrics.

**Comparison Tables:** Summaries of model performances with and without bias mitigation will help in identifying the most effective combinations.

## 6. Tools and Frameworks

- **Python Libraries**: scikit-learn for model implementation, Pandas and NumPy for data preprocessing, and Matplotlib/Seaborn for data visualization (Pedregosa et al., 2011).
- **AIF360 Toolkit**: Used for bias detection and mitigation. Its modular architecture facilitated integration with our pipeline (Bellamy et al., 2018).

# Results / Discussion

**Model Performance Comparison**

A comprehensive comparison of various models was conducted to evaluate their performance in predicting heart disease. The results are summarized in the table below:

| Model | Recall | Precision | F1 Score | Accuracy | ROC-AUC |
|---|---|---|---|---|---|
| Multilayer Perceptron | 0.822319 | 0.738207 | 0.777996 | 0.762648 | 0.843023 |
| Support Vector Machine | 0.811304 | 0.749708 | 0.779291 | 0.767580 | 0.833215 |
| Gradient Boosting | 0.801011 | 0.751482 | 0.775457 | 0.765388 | 0.844505 |
| Random Forest | 0.781510 | 0.734182 | 0.757107 | 0.746393 | 0.814778 |
| Logistic Regression | 0.774287 | 0.762311 | 0.768252 | 0.763744 | 0.842523 |
| Decision Tree | 0.669556 | 0.678872 | 0.674182 | 0.672694 | 0.673098 |

- **Multilayer Perceptron (MLP)** emerged as the best-performing model, achieving the highest recall of 0.822 and a competitive F1 score of 0.778. This highlights its ability to effectively minimize false negatives, which is crucial for medical predictions.
- **Support Vector Machine (SVM)** and **Gradient Boosting** also performed well, with similar F1 scores and ROC-AUC values, indicating their potential for reliable predictions.
- **Logistic Regression** demonstrated strong precision and ROC-AUC, showing its capability to discriminate between classes despite its simpler structure.
- The **Decision Tree** underperformed compared to other models, likely due to its tendency to overfit on smaller datasets.

After evaluating various models, the Multilayer Perceptron (MLP), a type of neural network, was chosen for further exploration due to its best performance and ability to capture complex non-linear relationships in the data.

**Neural Network Tuning and Evaluation**

Grid search was employed to optimize hyperparameters for the MLP model. The best configuration included:

- **Activation**: tanh
- **Hidden Layer Sizes**: (128)
- **Solver**: adam
- **Learning Rate**: constant
- **Early Stopping**: Enabled

Using this configuration, the tuned MLP achieved the following evaluation metrics:

- **Recall**: 0.814
- **Precision**: 0.747
- **F1 Score**: 0.779
- **Accuracy**: 0.767

The high recall and balanced F1 score confirm the MLP's effectiveness in handling the imbalanced dataset, making it well-suited for detecting heart disease cases where sensitivity is critical.

**Model Performance Evaluation**

| MLP Model | Reweighted Model | Prejudice Remover Model |
|---|---|---|
| Recall: 0.814<br>Precision: 0.747<br>F1 Score: 0.779<br>Accuracy: 0.767 | Recall: 0.798<br>Precision: 0.751<br>F1-Score: 0.774<br>Accuracy: 0.766 | Recall: 0.767<br>Precision: 0.756<br>F1-Score: 0.761<br>Accuracy: 0.764 |

The post-mitigation performance reveals a shift in the balance between fairness and model performance. The ANN model experienced a slight reduction in recall (from 0.814 pre-mitigation to 0.798 post-mitigation) and a slight decrease in precision (from 0.747 to 0.751). Despite these small changes, the model still maintained good overall performance, with the F1-score dropping slightly from 0.779 to 0.774 and accuracy remaining relatively stable at 0.765. The decrease in recall, while present, can be attributed to a trade-off for reducing bias and improving fairness, especially across different subgroups.

The Reweighted Model demonstrated similar results, with a recall of 0.798 (down from 0.798 in the pre-mitigation ANN) and an F1-score of 0.774, maintaining a strong balance between performance and fairness. The Prejudice Remover model showed a reasonable compromise with a

recall of 0.767 and an F1-score of 0.761, which slightly reduced the recall but continued to preserve reasonable accuracy and fairness.

**Model Performance Evaluation (Male vs. Female)**

| Pre-mitigation | Post-mitigation |
|---|---|
| **MLP on Female:**<br><br>Recall: 0.765<br>Precision: 0.697<br>F1-Score: 0.729<br>Accuracy: 0.752 | **Prejudice Remover on Female:**<br><br>Recall: 0.767<br>Precision: 0.756<br>F1-Score: 0.761<br>Accuracy: 0.760 |
| **MLP on Male:**<br><br>Recall: 0.890<br>Precision: 0.760<br>F1-Score: 0.820<br>Accuracy: 0.778 | **Prejudice Remover on Male:**<br><br>Recall: 0.782<br>Precision: 0.756<br>F1-Score: 0.769<br>Accuracy: 0.765 |

In terms of gender-based performance, pre-mitigation results highlighted significant disparity between male and female performance, especially in recall and F1-score. Male outperformed females with a recall of 0.890 versus 0.765 for females and an F1-score of 0.820 compared to 0.729 for females. However, post-mitigation, the gender gap narrowed. The Prejudice Remover model improved female recall to 0.782, a notable increase, while slightly reducing male performance (with a post-mitigation F1-score of 0.761).

Overall, bias mitigation strategies such as reweighting and prejudice removal improved fairness across gender without a substantial sacrifice in model performance. This suggests that while there was a slight trade-off in recall and F1-score, the models became more balanced and fair, addressing disparities in performance across gender groups. The results reflect a successful application of bias mitigation, enhancing fairness without major losses in model accuracy and effectiveness.

# Conclusion

This study explored the application of various machine learning models for heart disease prediction, with a focus on achieving a balance between sensitivity and precision. Through extensive evaluation, the **Multilayer Perceptron (MLP)**, a type of neural network, emerged as the best-performing model, achieving a recall of 0.822 and an F1 score of 0.778. These metrics highlight its effectiveness in identifying heart disease cases while maintaining balanced predictions.

Post-mitigation performance, however, showed a slight decrease in recall (0.798) and precision (0.751), though the F1-score and accuracy remained strong, illustrating the MLP's robustness even after bias mitigation techniques were applied.

The decision to focus on recall is critical in medical applications, where false negatives could lead to missed diagnoses and severe consequences. The MLP's ability to minimize false negatives demonstrates its suitability for such high-stakes scenarios. Additionally, its performance underscores the importance of using advanced models capable of capturing non-linear relationships in complex datasets. With bias mitigation strategies such as Reweighing and Prejudice Remover, the model was able to improve fairness, particularly for gender groups, as demonstrated by the narrowing gender disparity post-mitigation. The Prejudice Remover model, for instance, improved female recall to 0.782 while reducing the performance gap, albeit with a slight sacrifice in male F1-score (0.761). This illustrates that bias mitigation techniques can improve fairness without a substantial loss in model accuracy.

Hyperparameter tuning using grid search played a pivotal role in optimizing the MLP's architecture, including selecting the appropriate activation function, solver, and hidden layer sizes. This process not only improved the model's predictive power but also showcased the potential of neural networks to adapt to imbalanced data when properly configured. Additionally, the Reweighted model and Prejudice Remover both showed that balancing fairness and performance is achievable, with F1-scores close to the original MLP and slight improvements in recall, especially for females.

While models like Support Vector Machine (SVM) and Gradient Boosting also delivered competitive results, the MLP's flexibility and scalability provided a distinct advantage. However, it is important to note that achieving optimal performance requires careful preprocessing, feature scaling, and parameter tuning, as demonstrated in this study. The application of bias mitigation techniques further emphasized that fairness can be integrated into machine learning models without sacrificing overall predictive power, ensuring that the model performs well across diverse groups.

Ensuring fairness in AI-driven healthcare models is essential for minimizing disparities and fostering trust in AI systems. When models are biased, they can perpetuate or exacerbate existing inequalities, leading to unfair treatment of certain groups, particularly those from historically marginalized communities. By addressing these biases, we help ensure that AI systems are not only technically effective but also ethically responsible, empowering all individuals to benefit from the advancements in healthcare technology. In doing so, we create more inclusive, trustworthy models that can support better health outcomes for everyone.

# References

1. https://aif360.res.ibm.com/
2. https://fairlearn.org/v0.11/user_guide/assessment/perform_fairness_assessment.html
3. https://themis-ml.readthedocs.io/
4. Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2021). A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6), 1-35.
5. 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.
6. World Health Organization. (2023). Noncommunicable diseases.
   https://www.who.int/news-room/fact-sheets/detail/noncommunicable-diseases
7. Marmot, M., & Wilkinson, R. G. (2005). Social determinants of health. *The Lancet*, 365(9464), 1099-1104.
8. Rajkomar, A., Hardt, M., Howell, M. D., Corrado, G., & Chin, M. H. (2018). Ensuring fairness in machine learning to advance health equity. *Annals of Internal Medicine*, 169(12), 866-872.
9. Braveman, P., Egerter, S., & Williams, D. R. (2011). The social determinants of health: Coming of age. *Annual Review of Public Health*, 32, 381-398.
10. https://www.kaggle.com/datasets/kamilpytlak/personal-key-indicators-of-heart-disease
11. Bellamy, R. K. E., Dey, K., Hind, M., et al. (2018). AI Fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic bias. *arXiv preprint arXiv:1810.01943*.
12. Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20, 273-297.
13. Feldman, M., Friedler, S. A., Moeller, J., et al. (2015). Certifying and removing disparate impact. *Proceedings of the 21st ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
14. Hardt, M., Price, E., & Srebro, N. (2016). Equality of opportunity in supervised learning. *Advances in Neural Information Processing Systems*, 29.
15. Little, R. J. A., & Rubin, D. B. (2019). *Statistical Analysis with Missing Data.*
16. Pedregosa, F., Varoquaux, G., Gramfort, A., et al. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830.
17. Flavio P. Calmon, Dennis Wei, Bhanukiran Vinzamuri, Karthikeyan Natesan Ramamurthy, and Kush R. Varshney, "Optimized Pre-Processing for Discrimination Prevention", Conference on Neural Information Processing Systems, 2017.
18. Elisa Celis, Lingxiao Huang, Vijay Keswani, Nisheeth Vishnoi, "Classification with Fairness Constraints: A Meta-Algorithm with Provable Guarantees", 2018
19. Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian, "Certifying and Removing Disparate Impact", ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2015.
20. Moritz Hardt, Eric Price, and Nathan Srebro, "Equality of Opportunity in Supervised Learning", Conference on Neural Information Processing Systems, 2016.
21. Faisal Kamiran and Toon Calders, "Data Preprocessing Techniques for Classification without Discrimination", Knowledge and Information Systems, 2012.
22. Faisal Kamiran, Asim Karim, and Xiangliang Zhang, "Decision Theory for Discrimination-Aware Classification", IEEE International Conference on Data Mining, 2012.

23. Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma, "Fairness-Aware Classifier with Prejudice Remover Regularizer", Joint European Conference on Machine Learning and Knowledge Discovery in Databases, 2012.
24. Geoff Pleiss, Manish Raghavan, Felix Wu, Jon Kleinberg, and Kilian Q. Weinberger, "On Fairness and Calibration", Conference on Neural Information Processing Systems, 2017.
25. Till Speicher, Hoda Heidari, Nina Grgic-Hlaca, Krishna P. Gummadi, Adish Singla, Adrian Weller, and Muhammad Bilal Zafar, "A Unified Approach to Quantifying Algorithmic Unfairness: Measuring Individual & Group Unfairness via Inequality Indices", ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2018.
26. Richard Zemel, Yu (Ledell) Wu, Kevin Swersky, Toniann Pitassi, and Cynthia Dwork, "Learning Fair Representations", International Conference on Machine Learning, 2013.
27. Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell, "Mitigating Unwanted Biases with Adversarial Learning", AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society, 2018.