

RESPONSIBLE AI

Week 7: Replication Project 3 – Bias Mitigation

Meira Gilbert and Nandita Rahman

TODAY'S OBJECTIVES

- Types of bias in model development
- Mitigating fairness issues in models through pre-, in-, and post-processing
- Re-learning and re-deploying models
- Explainability methods
- Replication Project Overview

BIAS DETECTION WITHIN AI SYSTEMS.

Types of Algorithmic Bias

Historical bias: misalignment between the world as-is and the values or objectives required from the ML model;

Representation bias: under-representation or failure for a population to generalize for groups in population;

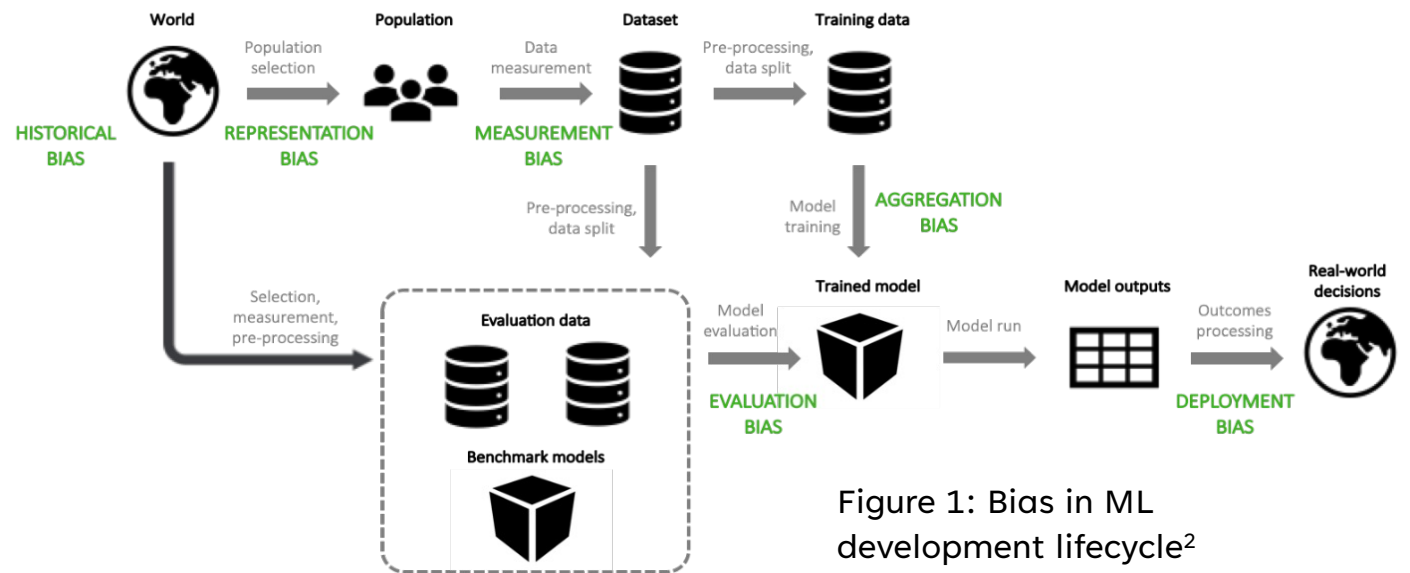
Measurement bias: choosing and utilizing features/labels that are noisy proxies for real-world quantities;

Aggregation bias: inappropriate combination of heterogeneous, distinct groups into a single model;

Evaluation bias: use of inappropriate performance metrics or the testing / external benchmark that does not represent the entire population; and

Deployment bias: inappropriate use or interpretation of model in a live environment.

Model bias is commonly thought to occur **primarily** in the data collection step.¹



Recent research suggests that bias can be introduced at **any stage** in the machine learning model life cycle (see Figure above).²

“ZOOMING OUT” - HEALTH EXPENDITURE USE CASE

Bias and Algorithmic Fairness

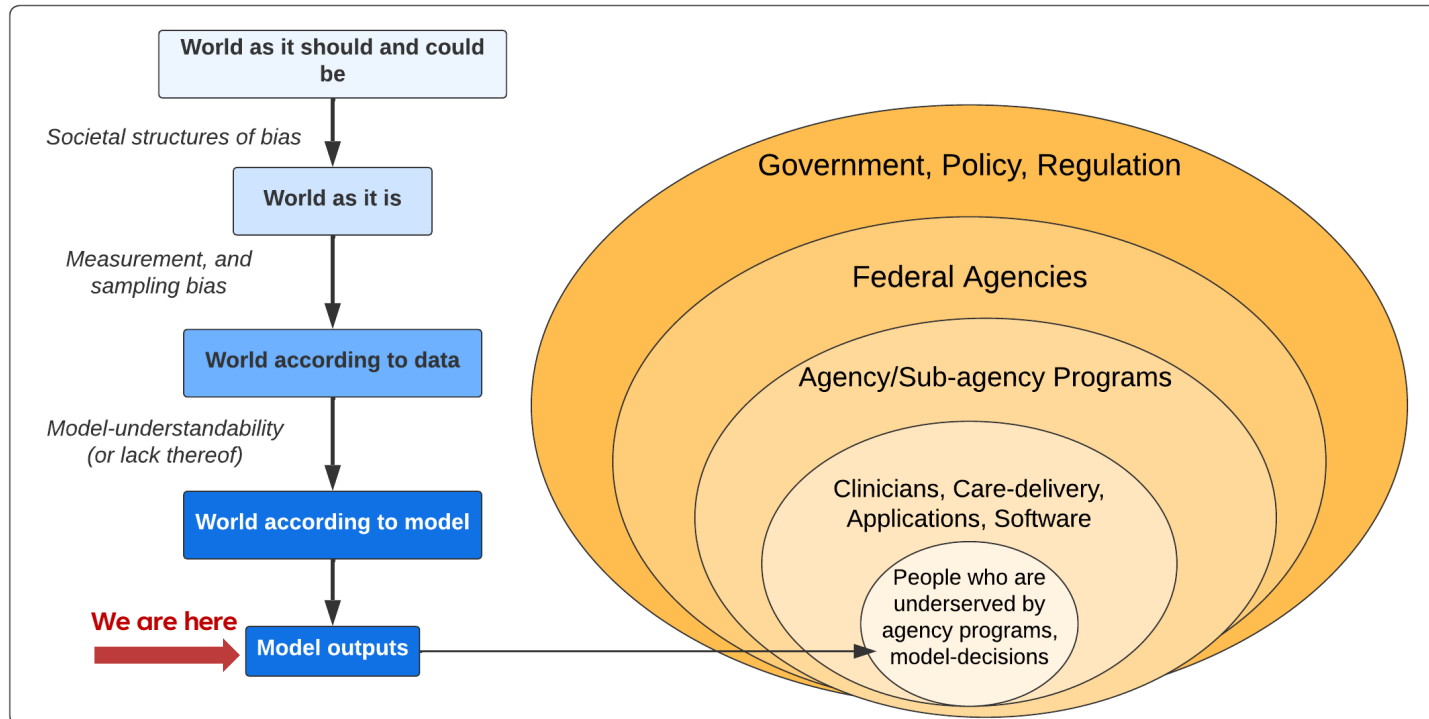


Figure above: Illustration of the types of bias that can enter a model's decision-space (stages of data and algorithm use), resulting in model-outputs that affect individuals downstream. (Adapted from: Mhasawadade 2020, Mitchell 2020)

If classifiers tied to model predictions are biased or not well-understood, this could lead to differentials in delivered action.

Composite scores for classifiers such as ‘**utilization**’ may indicate higher risk for certain groups over others – but is this a true representation of the real world?

Therefore, it is important to understand **different types of bias** and which de-biasing methods are best suited for your model.

Overrepresentation of a particular race in the training data will return more opportunity to learn fine-grained information about that group compared to others, despite **not explicitly including race as a feature**.

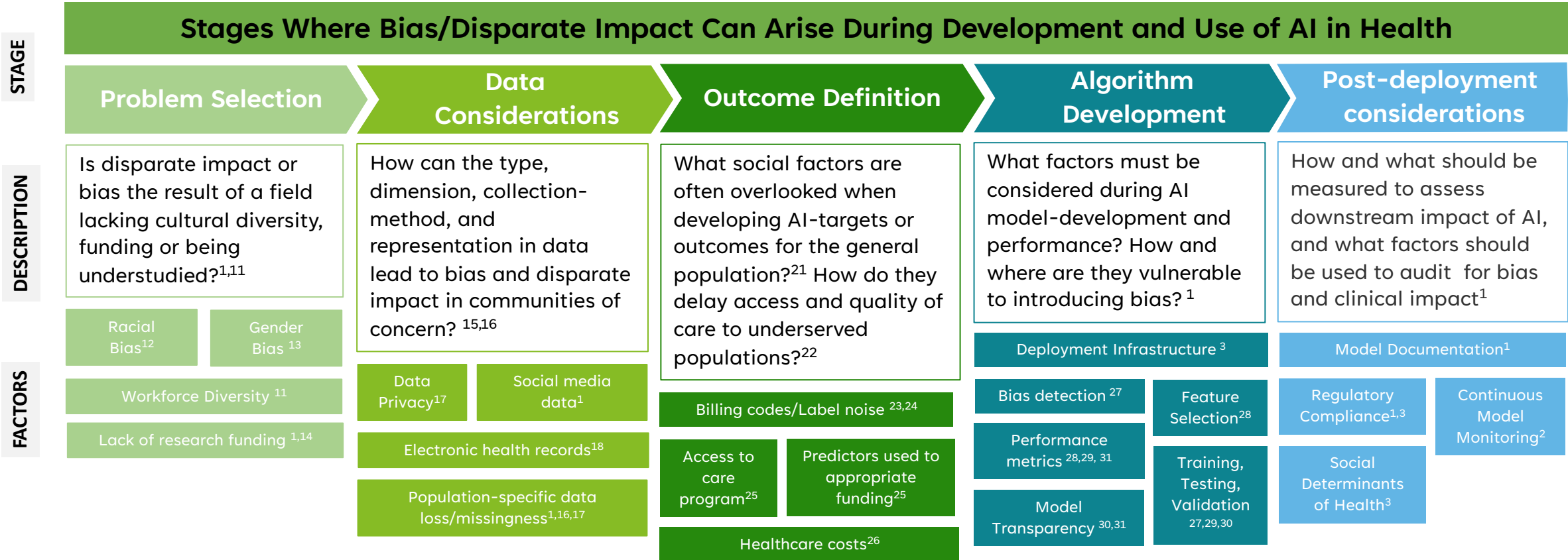
How would this affect Non-White beneficiaries who could have risk factor predictors that could inform a model to prioritize additional care? **How could it do the opposite?**

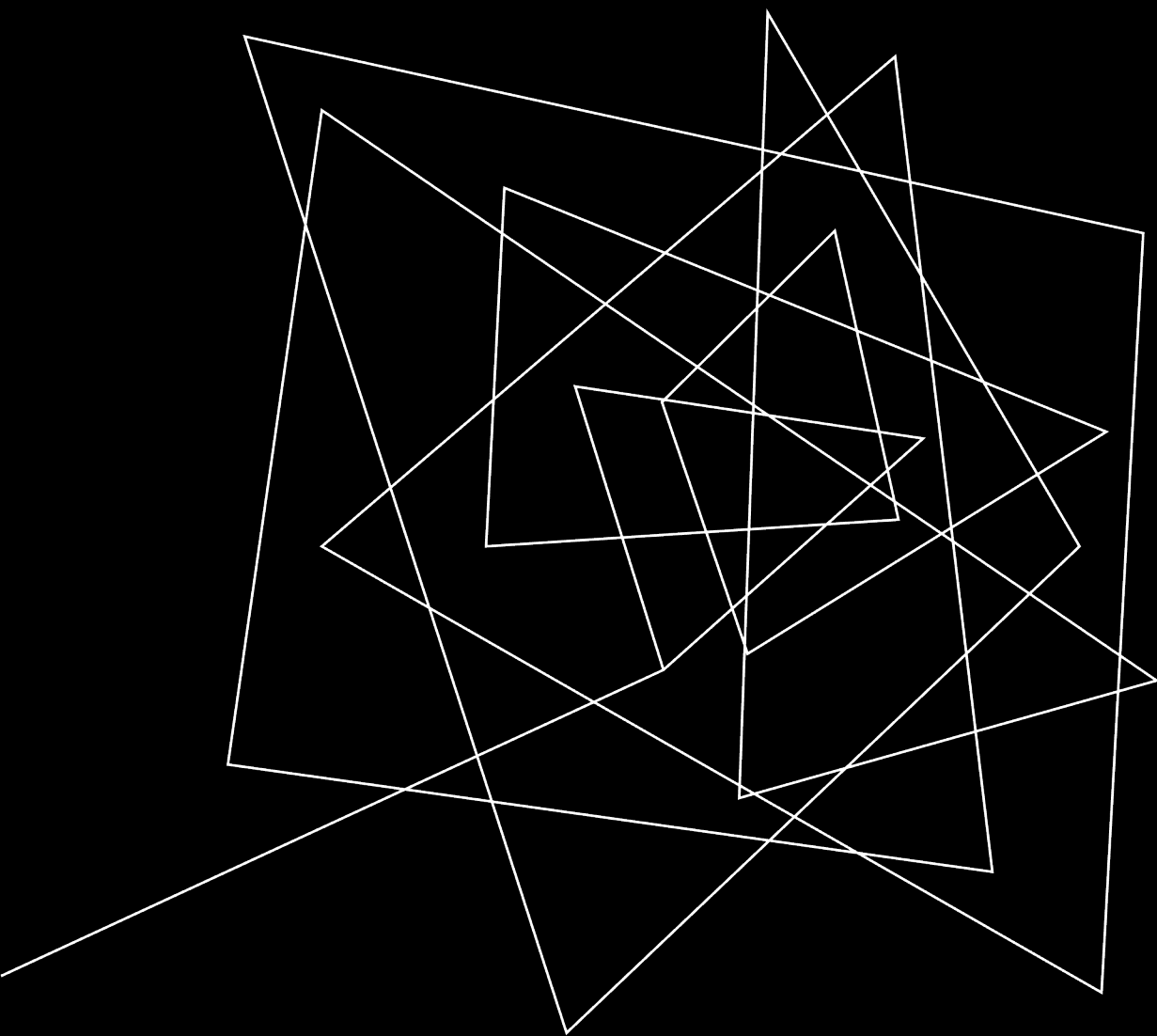
SOURCES OF BIAS^{1,2,3} IN AI AND HEALTH DATA – OPTIONAL SLIDE

Bias within AI in health can often hold co-existing sources both (technical and non-technical)^{3,4}.

(Taken from literature) are example “factors,” that have shown evidence of introducing bias during model development and use¹.

We prioritized factors that demonstrated contributing to:	
Negatively impacting group(s) disproportionately relative to general population ^{4,5}	Posed a barrier for advancing health equity ^{6,7} within historically marginalized group(s) or communities of concern ^{9,10}





READING PRESENTATION #2

Fairness Through Awareness
(Dwork et al.)

Which type of statistical fairness should you strive for?

RELEVANT DEFINITIONS

DISPARATE TREATMENT
Disparate treatment is a legal term defined as negative treatment of a person dependent on group of loan candidate size ability to that candidate's protected status (race, ethnicity, gender, etc).

DISPARATE IMPACT
Disparate impact is a legal term defined as unintentional but systemic negative treatment of a protected group of loan candidates. Because ML models lack a human decision maker to ask about their intent or reasoning, it's not always clear how disparate treatment and impact should apply to algorithms. Regulators should clarify this.

Credit to:

Valerie Cormier, "How to define fairness to detect and prevent discriminatory outcomes in machine learning," for more, good examples of when to use each type of fairness or when to use both.

Zhenyue Zhang, "A Tutorial on Fairness in Machine Learning" (for wrap-up of controversies around different types of fairness).

Moritz Hardt, Eric Price, Nathan Srebro, "Equality of Opportunity in Supervised Learning" (for comparison of equality of opportunity and odds).

Simon Barocas, Moritz Hardt, Arvind Narayanan, "On the (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

Allen Xiang, Indrakesh Debnath, "On the Legal (Im)possibility of Fairness Indicators" (for implications of disparate impact and disparate treatment).

You've got a machine learning system. You want it to be fair. But there are so many ways to be fair! Which should you choose?

By Samara Trilling and Madison Jacobs

START

Do you have, or is it possible to acquire, ground truth data on actual thing you want to predict?

E.g., You care about predicting crime, but you only have arrest data (a proxy). You care about predicting default on a loan, but you only have info on if someone was granted a loan (a proxy).

I have or can acquire ground truth data

I only have access to proxy data

Apply bias mitigation measures to your proxy training data. Then proceed with caution to making predictions about your actual target.

Bias Mitigation Measures:
See this [blog post](#) for a brief description of some bias mitigation measures and this [IBM toolkit](#) for a more comprehensive list, plus tools to help you implement them.

[Quickly revisiting last week's demo...](#)



DISCUSSION QUESTIONS

- How can we determine what bias mitigation techniques to use?
- What should we do if we encounter unexpected results from our bias mitigation?
- Given our use case and data, what bias mitigation techniques seem reasonable? Unreasonable?

BIAS MITIGATION TECHNIQUES



Preprocessing

Used to mitigate bias in training data, before building model



- Reweighting
- Optimized Preprocessing
- Learning Fair Representation
- Disparate Impact Remover



In-processing

Used to mitigate bias in classifiers, while developing model



- Adversarial Debiasing
- Prejudice Remover



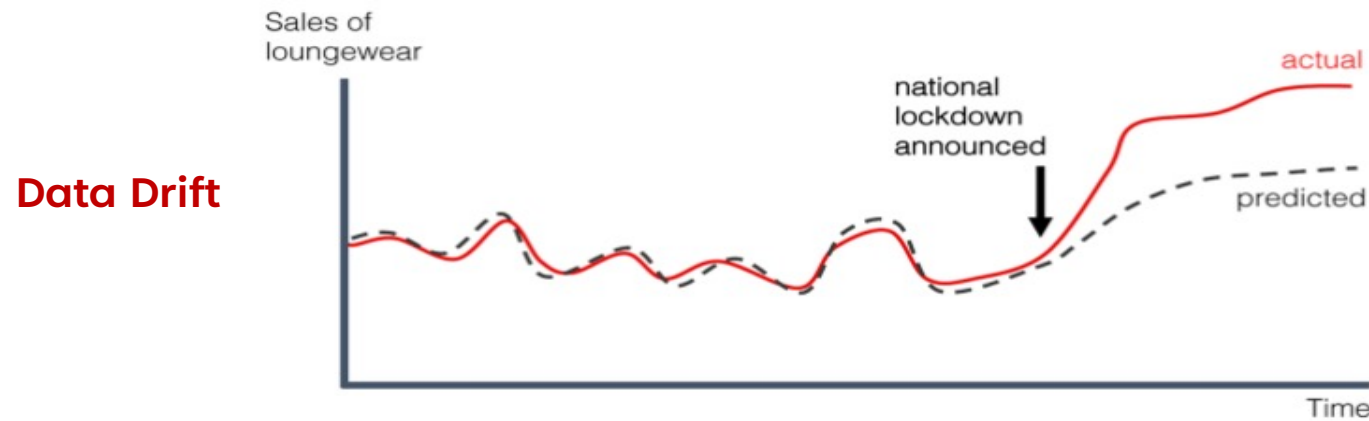
Post-processing

Used to mitigate bias in predictions/outcomes, after training the model



- Equalized Odds Postprocessing
- Calibrated Equalized Odds
- Reject Option Classification

RE-LEARNING AND RE-DEPLOYING MODELS



1. Deploy Model

Test model trained on 2014 (Panel 19) data on 2015 (Panel 20) data.

Does it exhibit fairness and maintain accuracy?

2. Re-Deploy Model

Test model trained on 2014 (Panel 19) data after **reweighing** on 2016 (Panel 21) data

Is there any drift?

3. Re-Learn Model

Re-learn model from 2015 (Panel 20) data.
Train and evaluate on **transformed** 2016 (Panel 21) data.

Is it relatively fair and accurate?

4. Re-Deploy Model

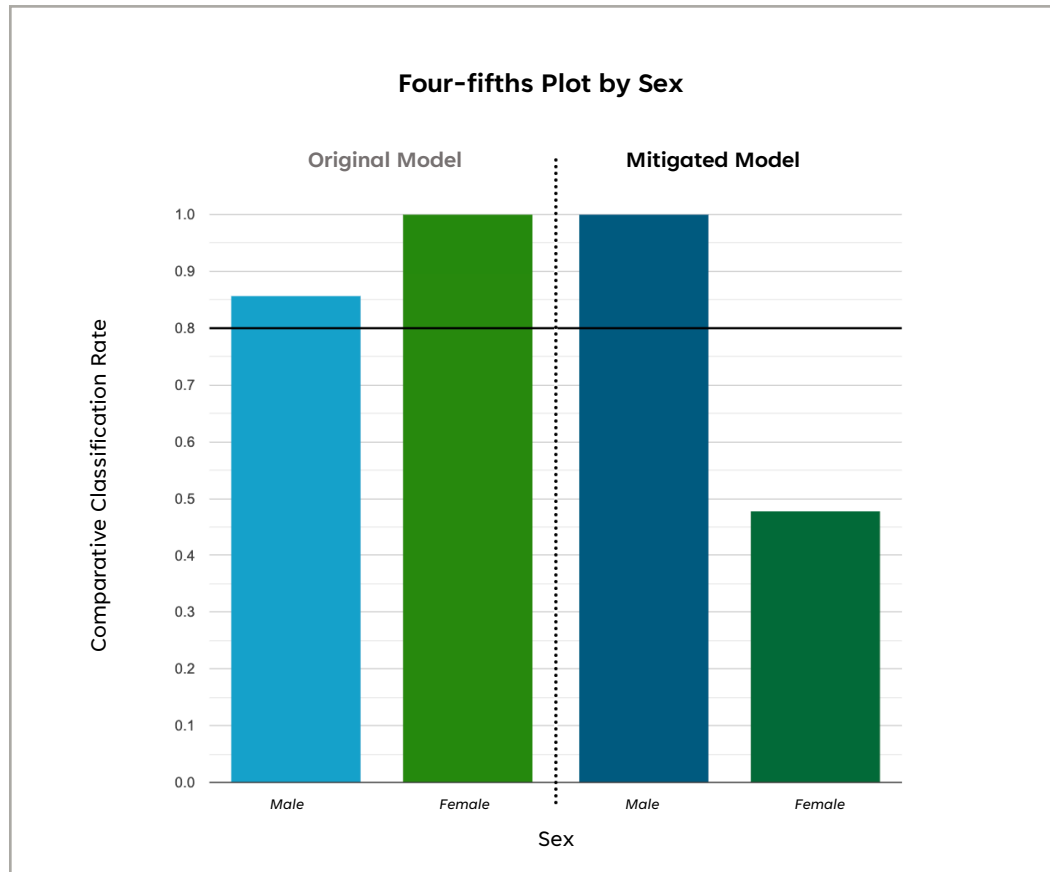
Test model trained on **transformed** 2015 (Panel 20) on 2016 (Panel 21) deployed data

Does it original fairness meet& accuracy specs?

FAIRNESS PLOT EXAMPLES WITH MITIGATED MODELS

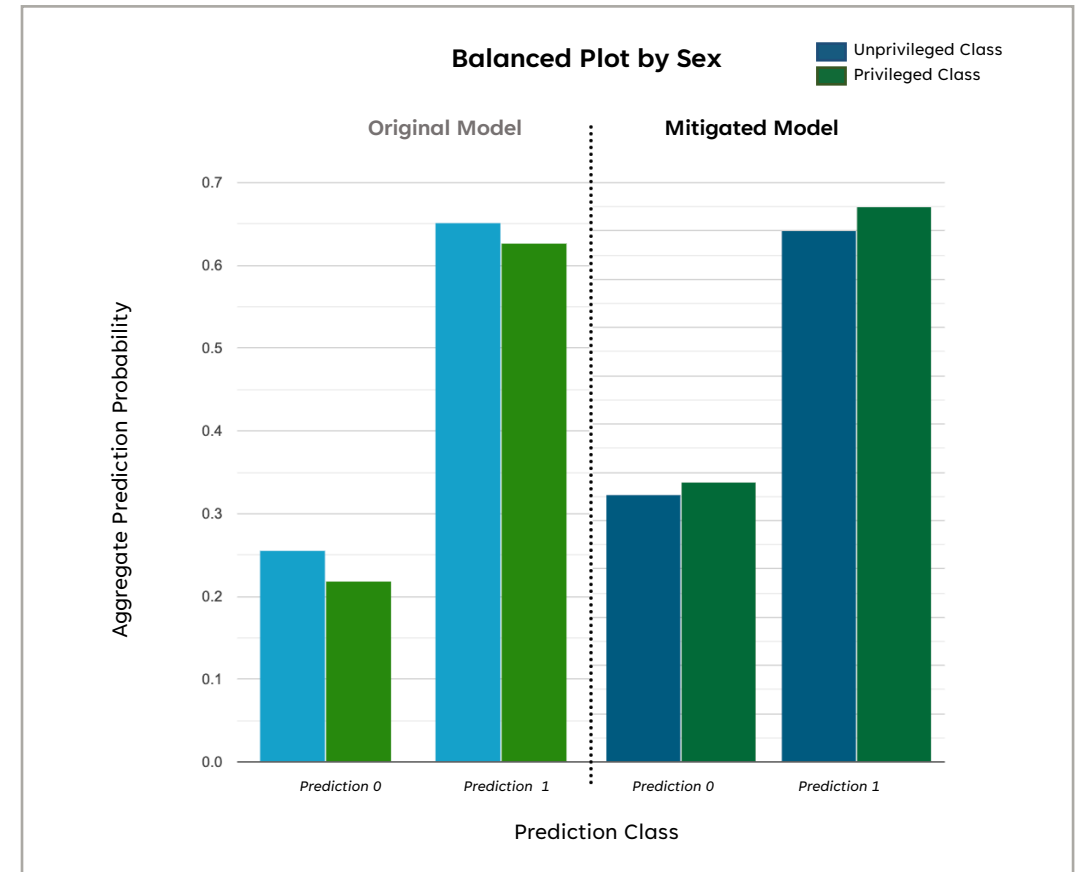
Four Fifths Plot

Identifies if there is adverse impact for unprivileged groups in comparison to the group with the highest selection rate.



Balance Plot

Examines whether average score received by individuals in positive and negative instances are similar regardless of sensitive attributes.

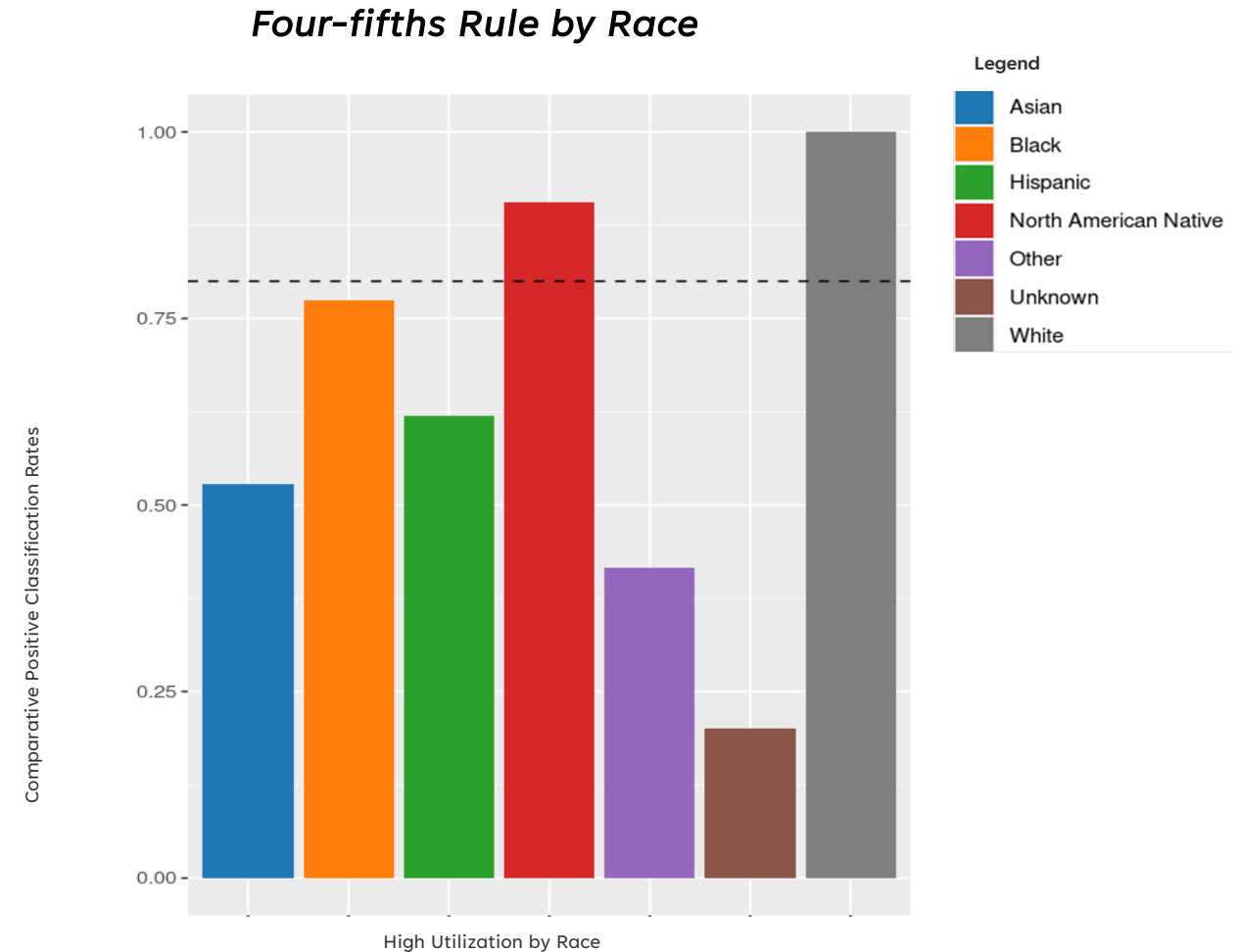


FAIRNESS PLOT EXAMPLES WITH RACE AND HEALTH EXPENDITURE

The four-fifths rule is a legal standard positing that if the selection rate for a certain group is less than 80 percent of that of the group with the highest selection rate, there is adverse impact on that group.

Whether White beneficiaries had the highest rates of 'high' utilization; could be largely reflective of different base rates in the data.

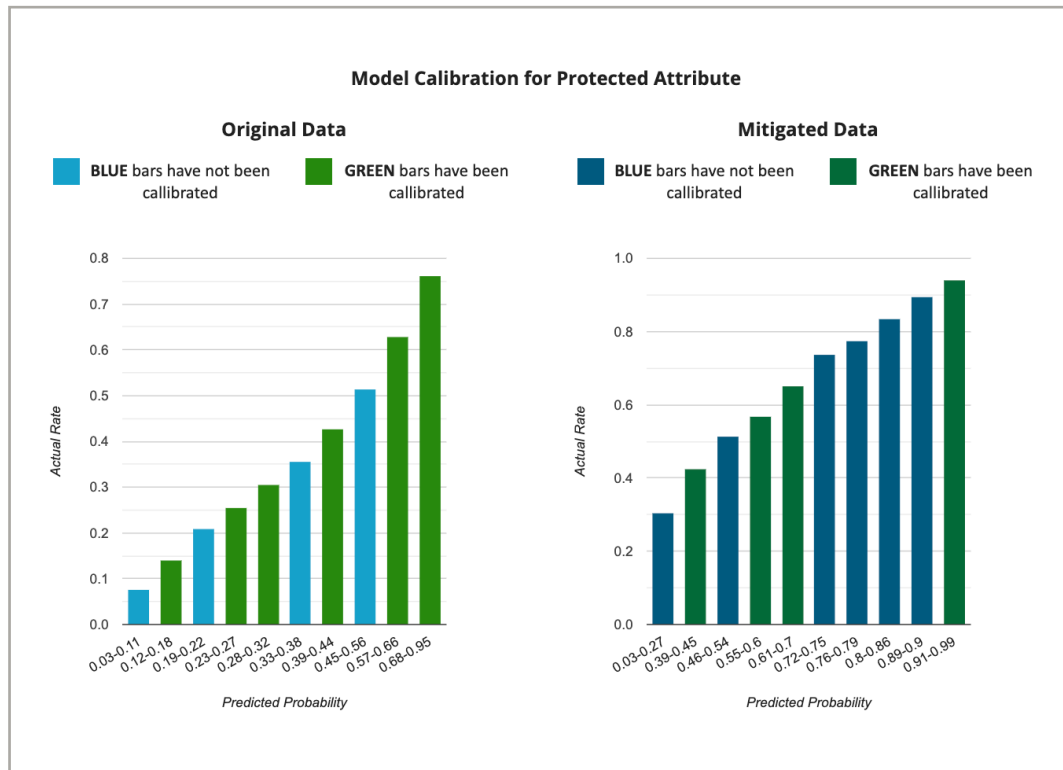
In the example on the right: the greatest discrepancy is for Other/Unknown groups – may be **indicative of missing data bias**.



OTHER FAIRNESS VISUALIZATIONS (*IMPUTED*) WITH MITIGATED DATA

Calibration Plots

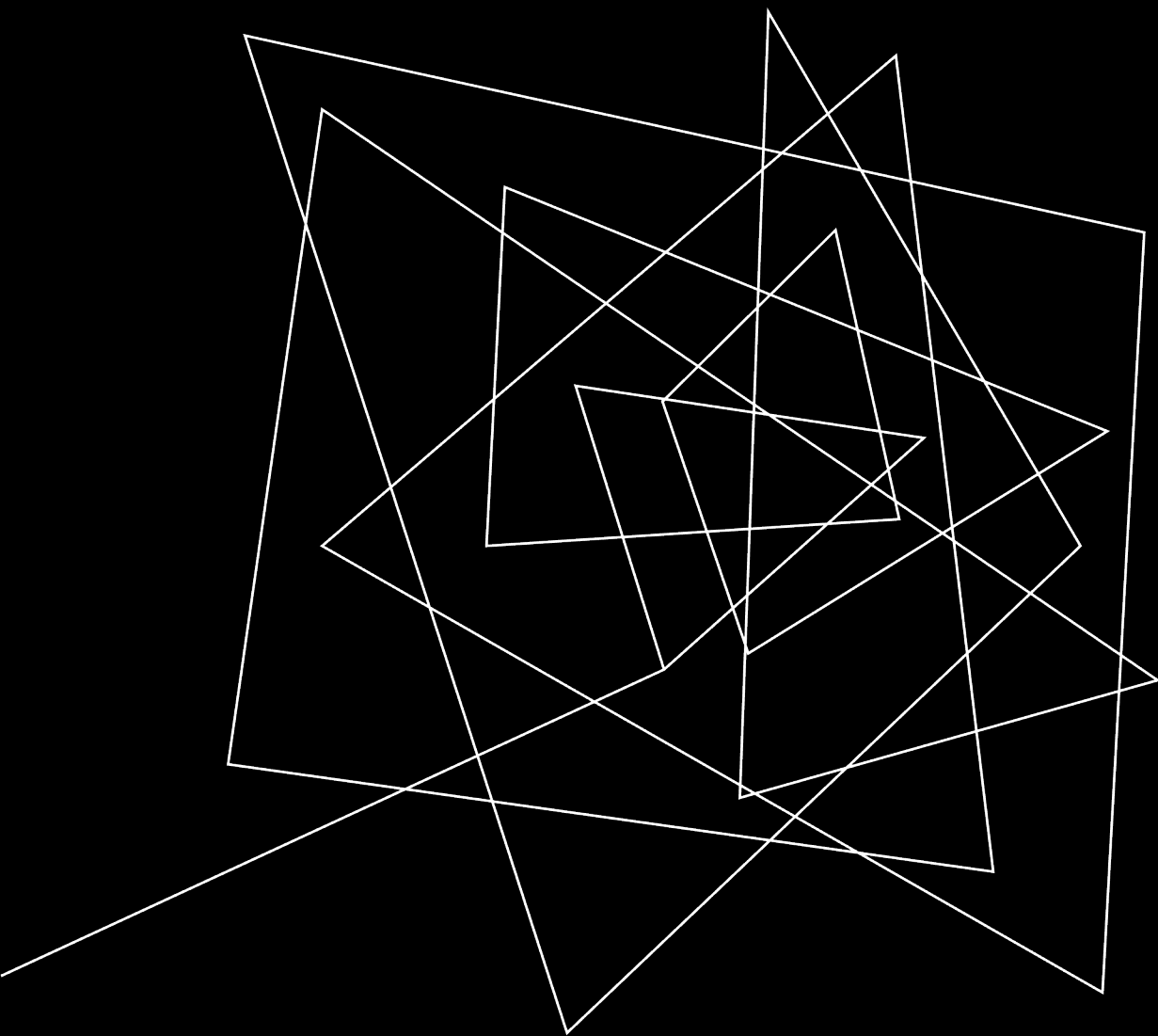
Checks if model makes accurate predictions in aggregate for members of each class.



Fairness Ratios

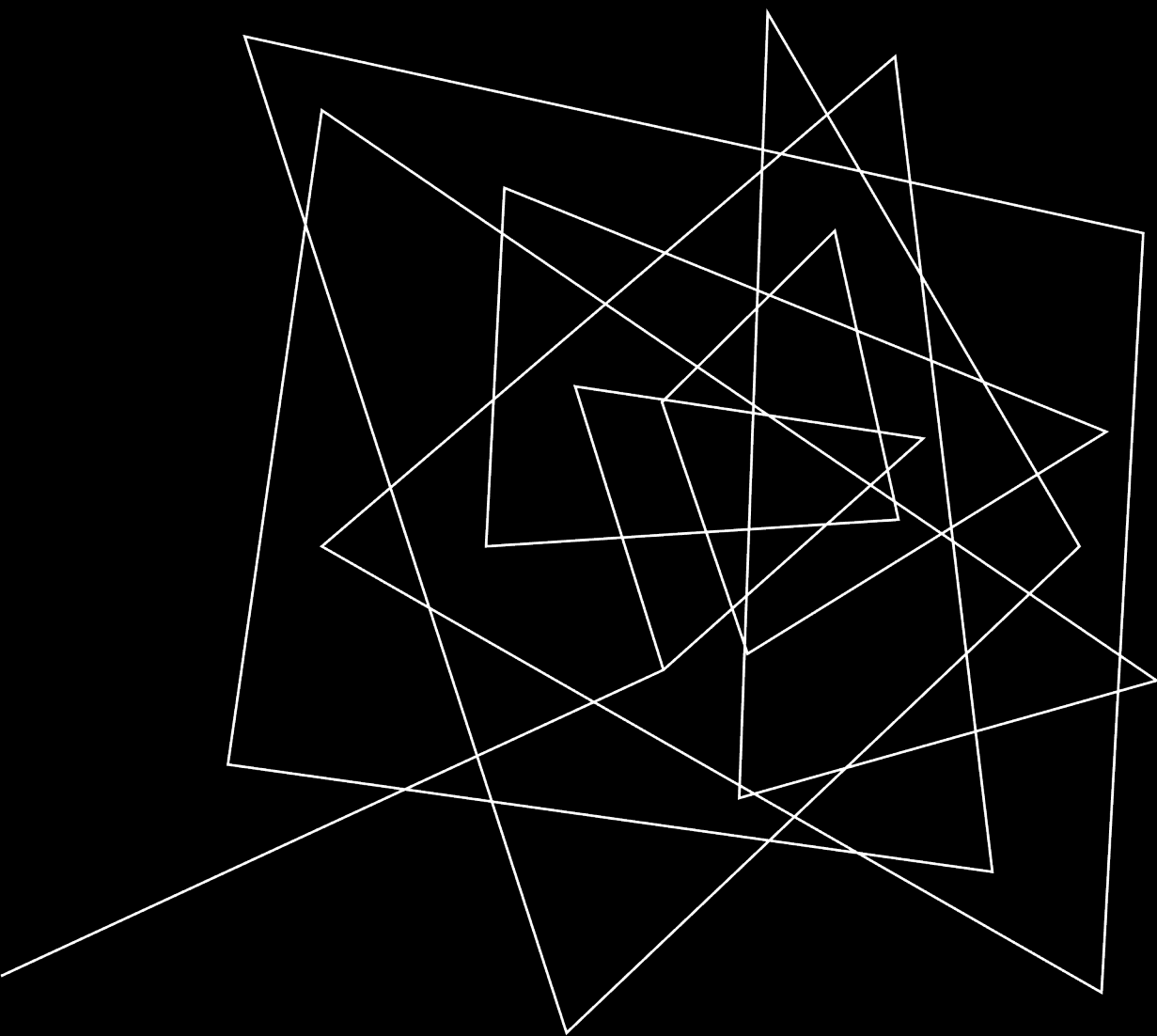
A selection of metrics is designed to help data scientists detect and evaluate bias within AI models.





READING PRESENTATION #1

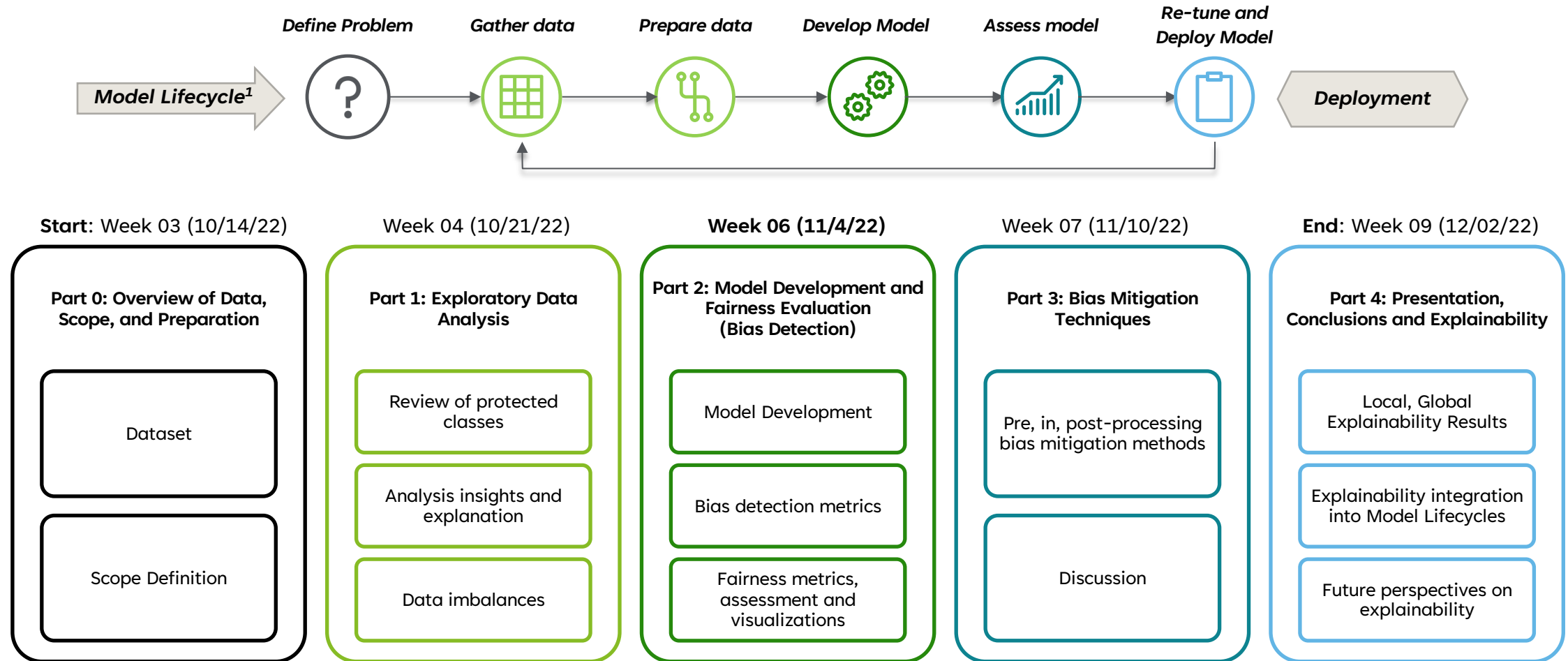
Model Cards for Model
Reporting (Mitchell et al.)



REPLICATION PROJECT PART 3

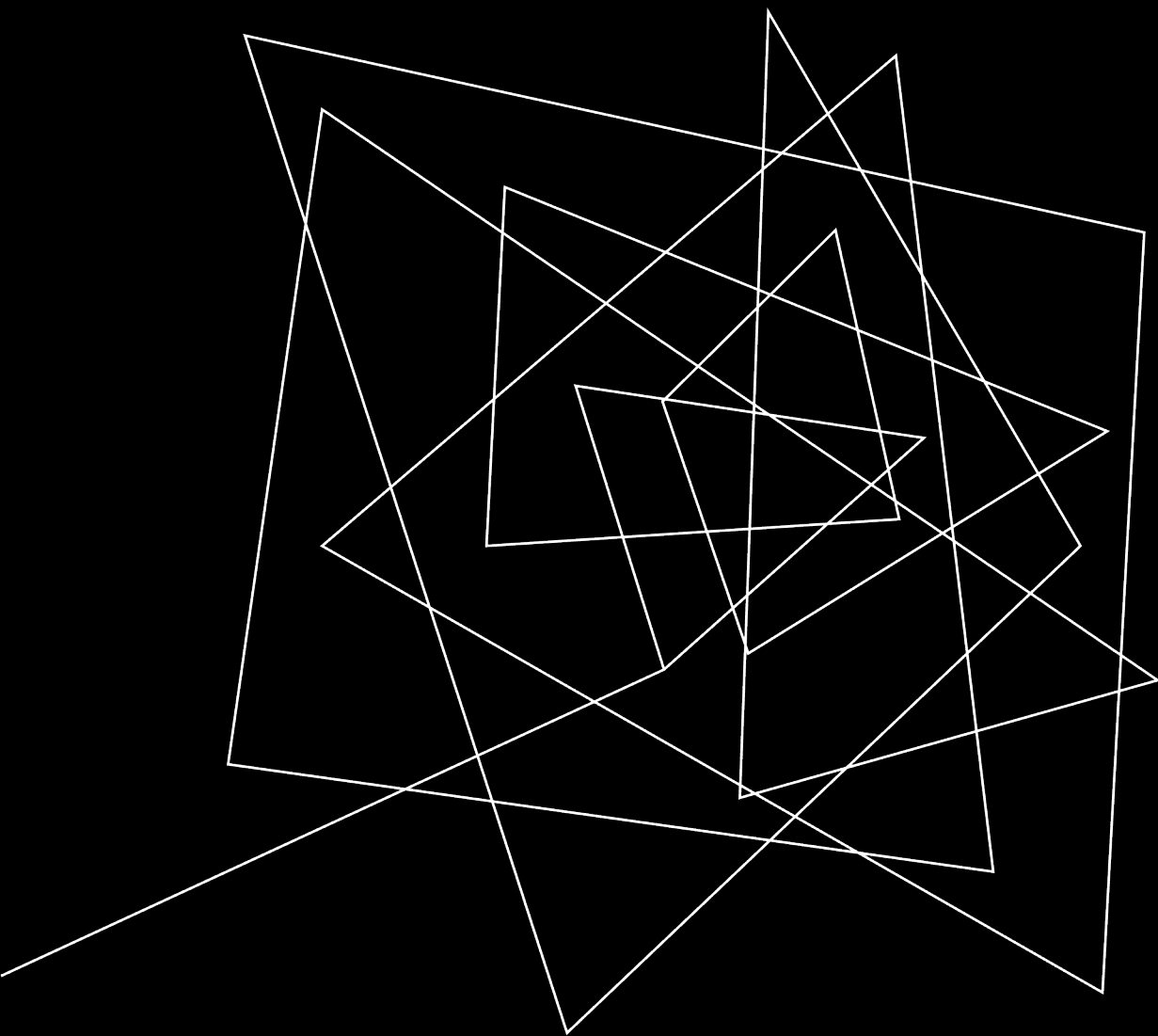
REPLICATION PROJECT: RESPONSIBLE AI IN ACTION

OPERATIONALIZING RESPONSIBLE AI ISN'T JUST USING CODE.
IT ALSO REQUIRES HUMANS IN THE LOOP FOR EACH STAGE OF THE MODEL LIFECYCLE



FOR NEXT WEEK

- Complete next week's **reading**
 - If you signed up to present **Algorithmic Fairness and Vertical Equity: Income Fairness with IRS Tax Audit Models (Black et al.)** come prepared to present next week and submit your presentation to Gradescope by 10 AM PT, Thursday, November 17th
- Submit your answers to next week's participation questions on Gradescope by 10 PM PT, Thursday, November 17th
- **Replication Part #3:** Notebook and writeup on bias mitigation techniques
 - Primary contact for replication project: Nandita Rahman (nanrahman@deloitte.com)
 - Office hours: Mondays 1-2pm PST
 - Notebook link: **Will be uploaded by Friday 11/11/22 + Instructions**
 - Nandita hosting virtually - [Zoom link](#)



MODEL EXPLAINABILITY METHODS

MODEL EXPLAINABILITY: LIME

Local Interpretable Model-agnostic Explanations

Using the LIME technique, a more complex model is interpreted by training less complex 'surrogate' models which provide explanations for individual data points.

LIME and SHAP are **attribution methods**, meaning that the prediction of a single instance is described as the **sum of feature effects**.

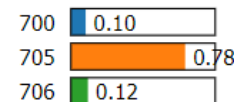
Patent Example:

This technique shows you which words were the most important. The numbers (700, 705, and 706) are different review groups. The NLP model works by connecting certain words to each of these groups.

1. Select the patent for which you'd like to interpret the explanation (e.g., why was patent X sorted into the review queue 705) .
2. LIME then creates fake patent data to train a new model, the **surrogate model**. This new model is **local**, meaning it can only be used to explain the one patent you selected.
3. The surrogate model uses a simpler algorithm (such as linear regression) and is more **explainable**, allowing for the creation of bar charts like the one show to the right.
4. By analyzing the bar charts, you can explain which words or **model features** contributed the most to its predictions or determination of *class* (review queue number).

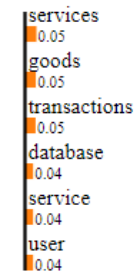
```
[64]: display(HTML(exp.as_html()))
```

Prediction probabilities



NOT 705

705



patent disclosure, as it appears in the United States Patent and Trademark Office files or records, but otherwise reserves all copyright rights whatsoever.

TECHNICAL FIELD

The present invention relates generally to a web-based procurement method and monitoring thereof. More particularly, the present invention relates to a method of procuring **goods** and/or **services** via a global communications network such as the Internet in which procurement **transactions** are monitored, costs accumulated and compared to predefined procurement (e.g., budgetary spending) thresholds, goals and/or constraints, and notifications based on the comparisons are communicated to persons of interest to facilitate the control of the procurement of the **goods** and/or **services**. In one embodiment, the procurement **transactions** (e.g., orders for **goods** and/or **services**) are automatically reviewed and at least one of approved, denied or placed on hold at a point of ordering based on the comparisons of accumulated costs to budgetary spending thresholds, goals and/or constraints.

BACKGROUND OF THE INVENTION

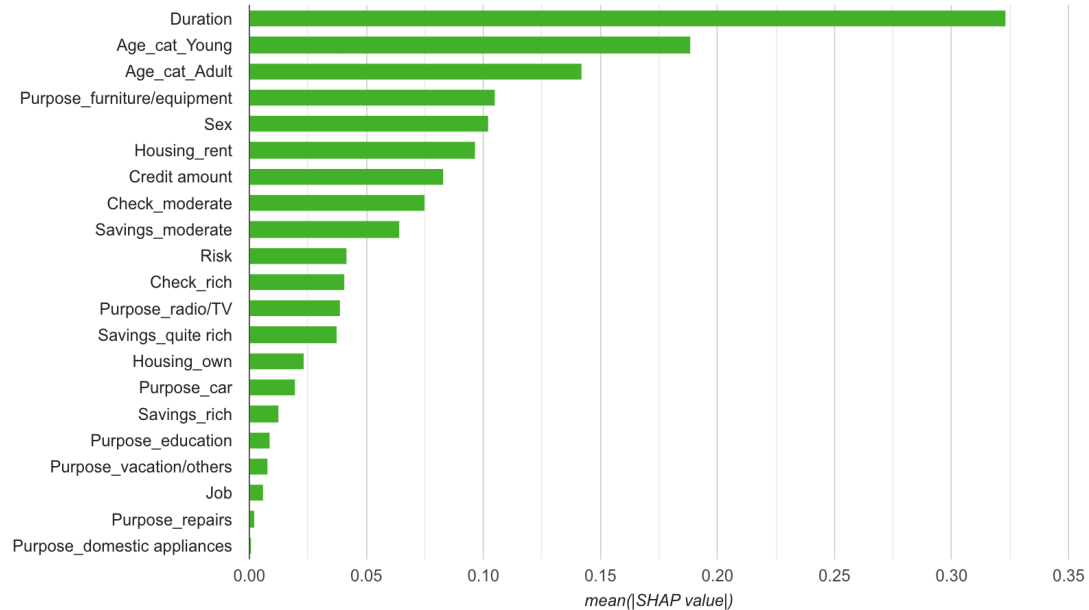
EXPLAINABILITY: SHAP

SHAP (**SH**apley **Ad**ditive **exP**lanations)

Global Explainability

Global explainability ensures the entire process of decision making is completely transparent

Global Explainability Plot



In the bar plot, features are sorted by decreasing overall importance to a model's ability to predict an outcome.

Features list at the top of Y-axis have the most predictive power

However, these features should not be interpreted as having casual impact on the model's prediction

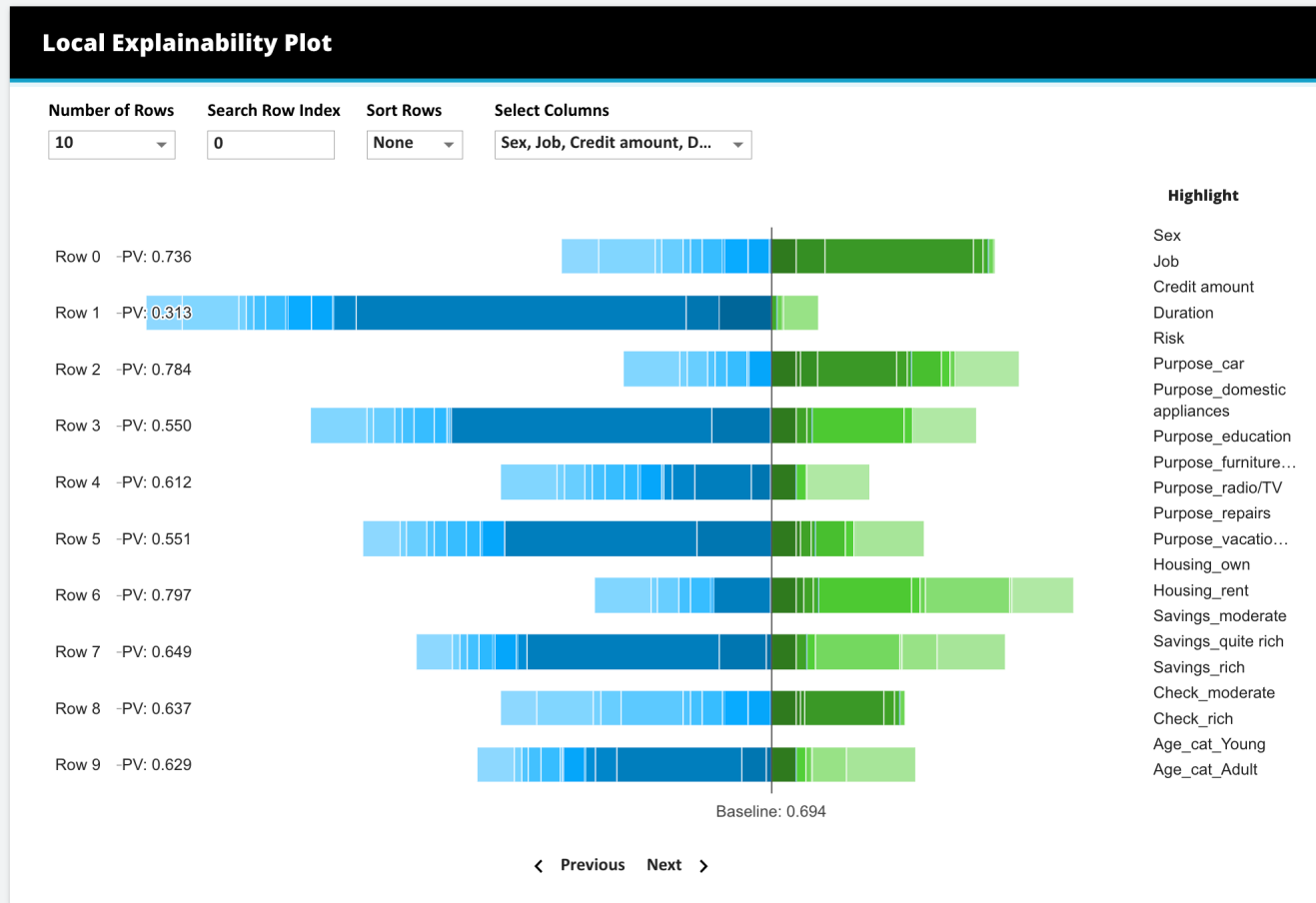
SHAP feature importance is measured as the "average Shapley value PER feature across ALL of the data."

Therefore, if a feature in this example plot has a mean Shapley value of 0.015, then that feature changes the model's predicted output on average by 1.5 percentage points

EXPLAINABILITY: SHAP

Local Explainability

- Local explainability provides explanations for each decision
- It takes a granular look at explaining how individual data points are used to make decisions



The **BASELINE** value represents the **average prediction across a dataset**

The **PREDICTION (PV)** value represents the **model's prediction** for that observation

The **BLUE** bars represent features that push the prediction value **higher**

The **GREEN** bars represent features that push the prediction value **lower**

What about other protected classes?

Potential things to consider when addressing bias for **Age**, **Race**, and **Sex**

Age



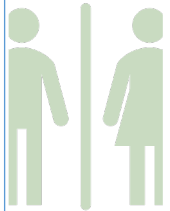
- Future work could consider identifying additional data sources to understand where proxies might pop-up, such as using age as a predictor – this would both increase accuracy and decrease the age-linked differential in error rates
- Youngest beneficiaries may be qualitatively different from older due to rates of disability.
- Consider separate models or hierarchical models which would explicitly learn the different dynamics of these types of subpopulations

Race



- Consider explicitly including race in models to correct for bias using pre-, in- and post- processing for bias mitigation
- Considerations of Impact (benefiting historically underserved groups by directing resources to group members at risk of mortality/hospitalization)
- Considerations for community stakeholders as “human values in the loop”¹, to suggest better data collection methods.
- “Unknown/other” missing data problem remediation requires investigating the source of upstream data collection

Sex



- Monitor for bias on an ongoing basis and examine intersectional bias for combinations of age and sex.
- Use visualizations such as the 4/5th plot to understand and communicate data remediation needs

