

RESPONSIBLE AI

Week 6: Replication Project 2 – Fairness Assessments

Credit: Meira Gilbert and Nandita Rahman

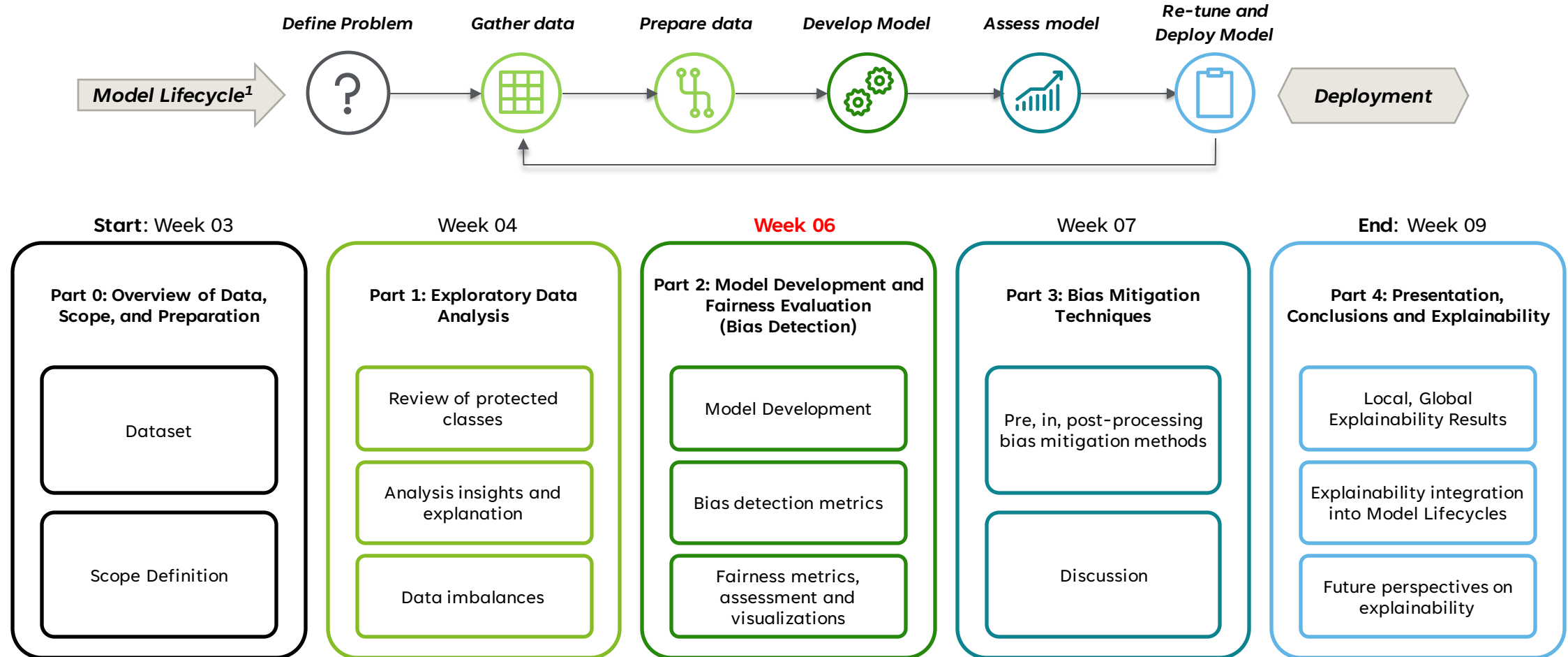


TODAY'S OBJECTIVES

- How to build a model?
- How can we assess algorithmic fairness? What metrics are commonly used?
- What can't be captured in data, and what are the limitations of fairness metrics?

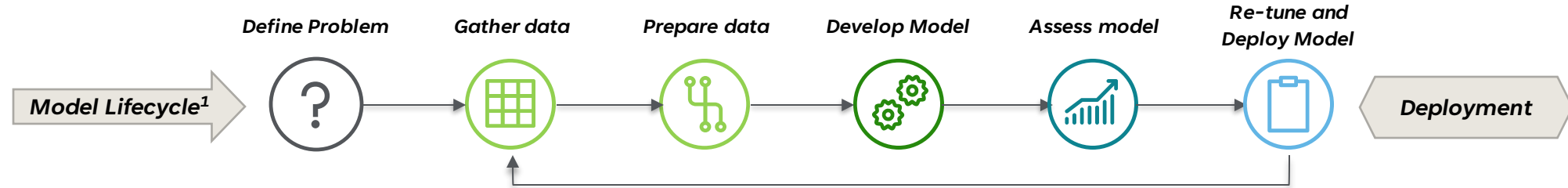
REPLICATION PROJECT: RESPONSIBLE AI IN ACTION

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



REPLICATION PROJECT PART 02

MODEL DEVELOPMENT & FAIRNESS EVALUATION



- **Model Development and Fairness Evaluation:**
After your Exploratory Data Analysis (Part 01), you'll now train your data and assess fairness metrics without de-biasing.
- There will be **TWO PARTS** to this portion of the replication project.
 - (1) Training models without de-biasing, using IBM's tutorial
 - (2) Training models without de-biasing, using any insights gained from your EDA step in Part 01. In addition, apply any model development techniques including (1) Feature Selection, (2) Encoding, and others that you may have learned in your methodology portion of this course.

Part 1

Use IBM AIF360's Tutorial to run training and testing, and capture fairness metrics

Evaluate whether the model performs sufficiently for production.

Does the model answer the question with sufficient confidence given the test data?

Part 2

Based on your EDA results should you collect additional data, do feature engineering, or experiment with other algorithms?

Use visualizations to understand your model

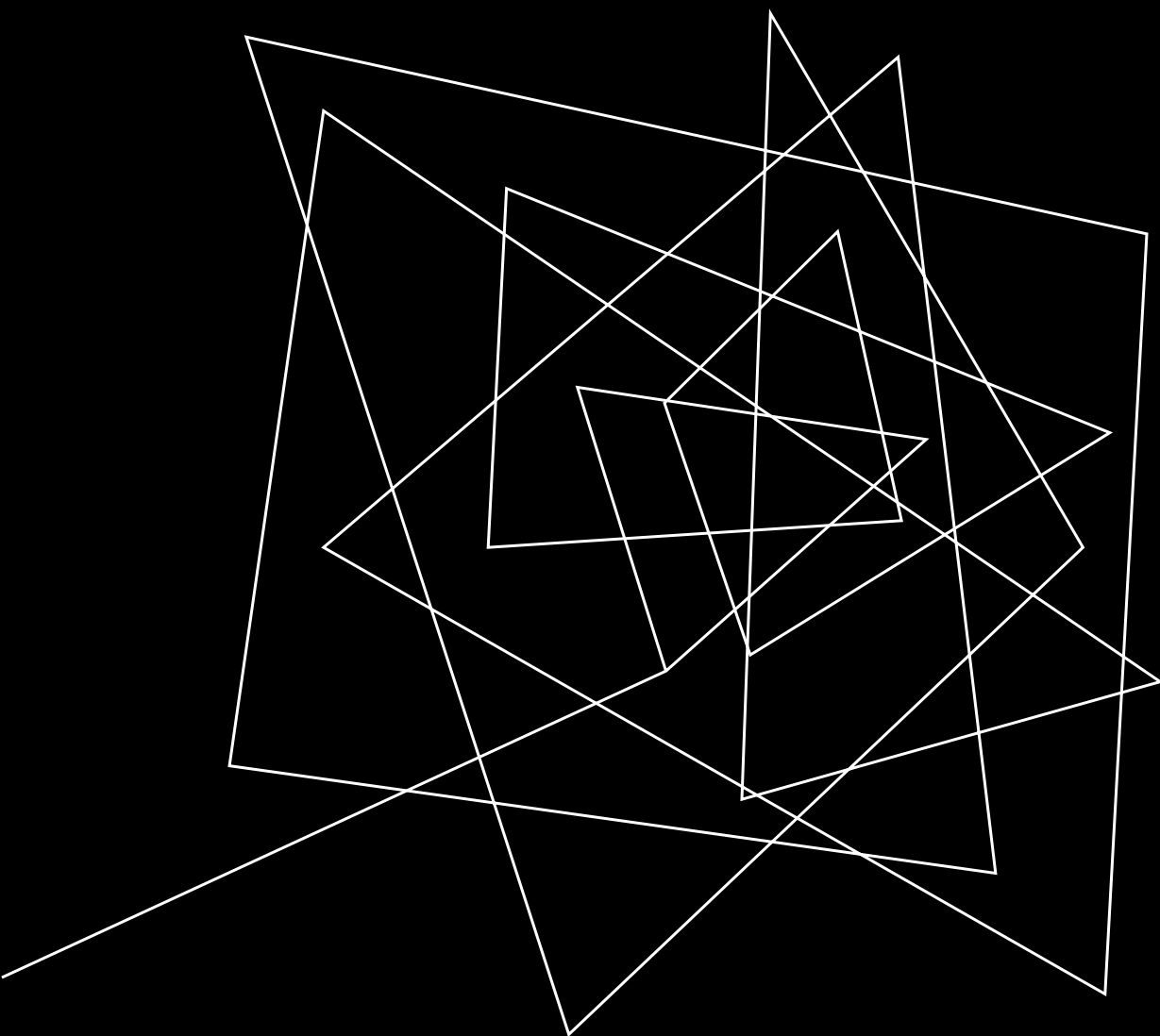
Rerun the same analysis you completed in Part 1 to capture fairness metrics.

Assessing Fairness

Explain the fairness of your model predictions using your fairness metric results.

Compare and contrast the differences between the fairness metrics results for Part 1 and Part 2.

Based on your results, provide your initial judgement for which type of **model** and **fairness metrics** seem appropriate to use.



READING PRESENTATION #1

Ethical Machine Learning in
Health Care (Chen et al.
2021)

FAIRNESS METRICS

FAIRNESS CRITERIA DESCRIBES THE CONNECTION BETWEEN SENSITIVE ATTRIBUTES AND TRUE/PREDICTED LABELS MATHEMATICALLY.

Statistical Parity / Demographic Parity	Equal proportion of outcomes between groups regardless of other factors
Statistical Parity Difference	Computed as the difference of the rate of favorable outcomes received by the unprivileged group to the privileged group
Equalized Odds	Equal false negative rates and false positive rates
Equal Opportunity	Equal false negative or false positive rates between groups
Equal Opportunity Difference	This metric is computed as the difference of true positive rates between the unprivileged and the privileged groups. The true positive rate is the ratio of true positives to the total number of actual positives for a given group.
Average Odds Difference	Computed as average difference of false positive rate (false positives / negatives) and true positive rate (true positives/positives) between unprivileged and privileged groups.
Disparate Impact	Computed as the ratio of rate of favorable outcome for the unprivileged group to that of the privileged group.
Theil Index	Computed as the generalized entropy of benefit for all individuals in the dataset, with $\alpha = 1$. It measures the inequality in benefit allocation for individuals.

How can we determine which metrics to use, given our data and use case?

When you have competing fairness metrics, how to pick which to prioritize?

What do you do when you encounter different definitions for similar metrics?

DISPARATE IMPACT

Disparate Impact:

$$= \frac{P(\text{favorable} = 1) \mid C = \text{unprivileged}}{P(\text{favorable} = 1) \mid C = \text{privileged}}$$

$$P(\text{favorable} = 1) \mid C = \text{privileged}$$

In other words:

$$\frac{\# \text{ favorable from underprivileged}}{\# \text{ underprivileged}}$$

$$\frac{\# \text{ favorable from privileged}}{\# \text{ privileged}}$$

General Consensus (industry standard):

If your disparate impact is LESS THAN 80% - your model is unfair and should be revised.

STATISTICAL PARITY

Something is "fair" if there isn't a difference on whether or not you have a given attribute (A):

$$\text{disparity} = Pr(Y = y \mid A = 0) - Pr(Y = y \mid A = 1)$$

```
Statistical parity difference =  $\frac{\text{num\_positives(privileged=False)}}{\text{num\_instances(privileged=False)}} - \frac{\text{num\_positives(privileged=True)}}{\text{num\_instances(privileged=True)}}$ 
```

- **Under 0:** Higher benefit for the monitored group.
- **At 0:** Both groups have equal benefit.
- **Over 0** Implies higher benefit for the reference group.

General Consensus (industry standard):

There is no disparity if disparity = 0. This (obviously) will not occur – so we introduce a threshold.

Which type of statistical fairness should you strive for?

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 things 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).

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 [film](#) (pdf) for a more comprehensive guide to help you implement them.

I have or can acquire ground truth data

I only have access to proxy data

Do you care more that two similar individuals are treated the same, or that overall the groups will be treated the same?
E.g., that 2 people with a similar credit score get a similar interest rate, or that average interest rates are about the same for men and women?

The stories about Apple Card giving a woman a lower credit limit than her partner and LendingTree charging a Howard grad more than an NYU grad are both examples of caring about individual fairness.

Groups

Individuals

Do any of the features that predict your outcome correlate with race, gender, age, or other protected classes?
Spoiler: many features do!

Maybe OK to use Fairness through Unawareness (Don't tell your model about race, gender, age, etc.)

Conditional Statistical Parity (Require demographic groups to get the same number of loans - but only if they have the same creditworthiness)

Yes, and I want to change it

No

Yes, but I'm ok with that

Is the thing you're predicting subjective?

Yes

No

You probably have historical bias in your training data

Conditional Statistical Parity is the current fair lending standard and a legally accepted way to address disparate impact. But it can still discriminate if the creditworthiness features correlate with race or gender. And it requires everyone to agree on what 'creditworthiness' means, which makes it harder to add new data (like next repayment history) that might help some groups.

Apply bias mitigation measures to your proxy training data. Then proceed with caution.

Do you want to use this system to correct for existing structural bias in the world?

Yes

No

E.g., offering college students from undersourced high schools extra tutoring, writing promotion processes that don't disadvantage women

Do you have a plan to support the underprivileged group and prevent reinforcement of historical biases?

Yes

No

Use Demographic Parity

is one kind of error (false positives or false negatives) more OK with you than another?

Yes

No

Your model will probably perpetuate or exacerbate historical biases

Use Equal Opportunity

Equal false negative rate or Equal false positive rate

Equal false negative rate and false positive rate

Use Equalized Odds

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

Equal false negative rate and false positive rate

RELEVANT DEFINITIONS

DISPARATE TREATMENT

Disparate treatment is a legal term defined as negative treatment of a loan candidate or group of loan candidates due solely 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 - but 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:

Valeria Cortez, "How to define fairness to detect and prevent discriminatory outcomes in Machine Learning" (for many good examples of when to use each type of fairness)

Zhuang Zhang, "A Tutorial on Fairness in Machine Learning" (for examples of controversies around different types of fairness)

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

Solon Barocas, Moritz Hardt, Arvind Narayanan, Fair ML Book (for predictive parity)

Alice Xiang, Indolene Deborah Raj, "On the Legal Compatibility of Fairness Definitions" (for applications of disparate impact and disparate treatment)

Can everyone agree how to quantify exactly how similar two people are?

E.g. age, race, gender, education, previous job experience, financial history, or other demographic factors

Yes

No

Sorry! You can measure individual fairness, but you can't guarantee it.

Are you OK with explicitly using race, gender, age, or other protected classes in the model?

Yes

No

Individual Fairness

Fairness through awareness

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

Disparate Impact Standard

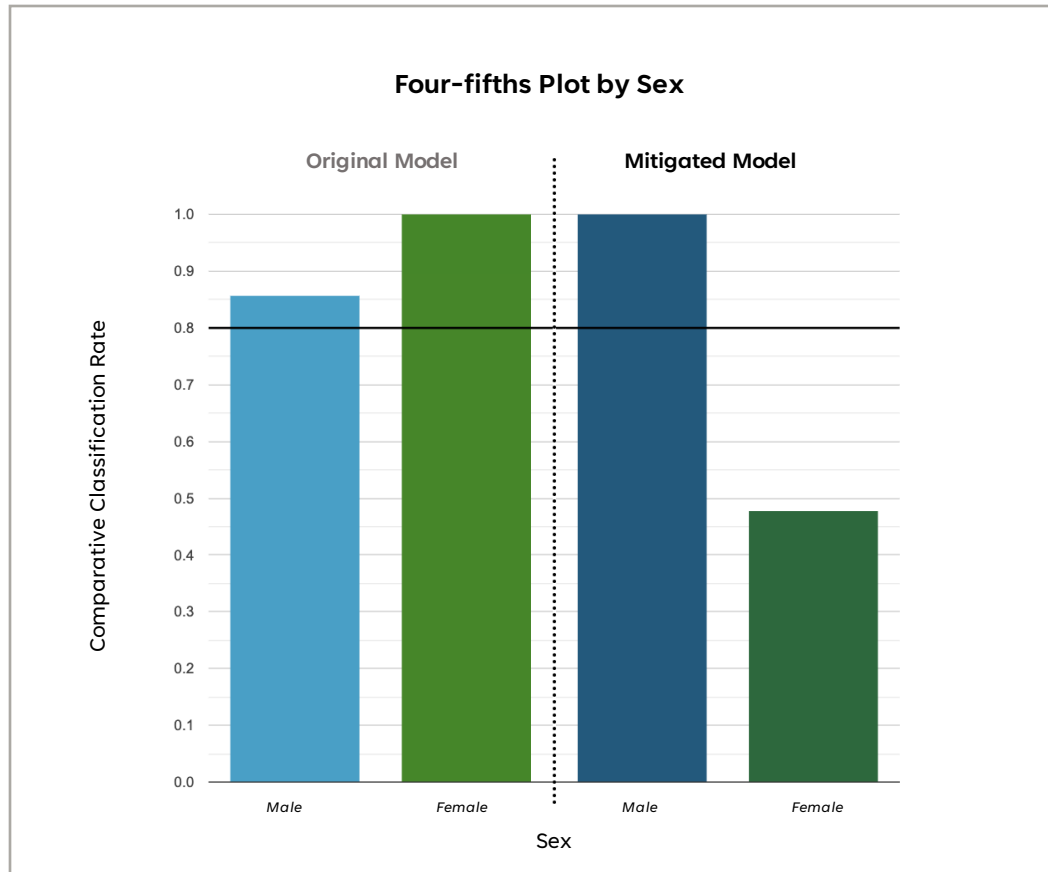
WHY DO I HAVE TO PICK?

Predictive parity, demographic parity, and equalized odds are mutually exclusive—you can't satisfy more than one. (Except in specific cases: E.g., if both groups are "actually" equally likely to default, then you can satisfy both demographic parity and equalized odds). Read more [here](#).

FAIRNESS PLOT EXAMPLES

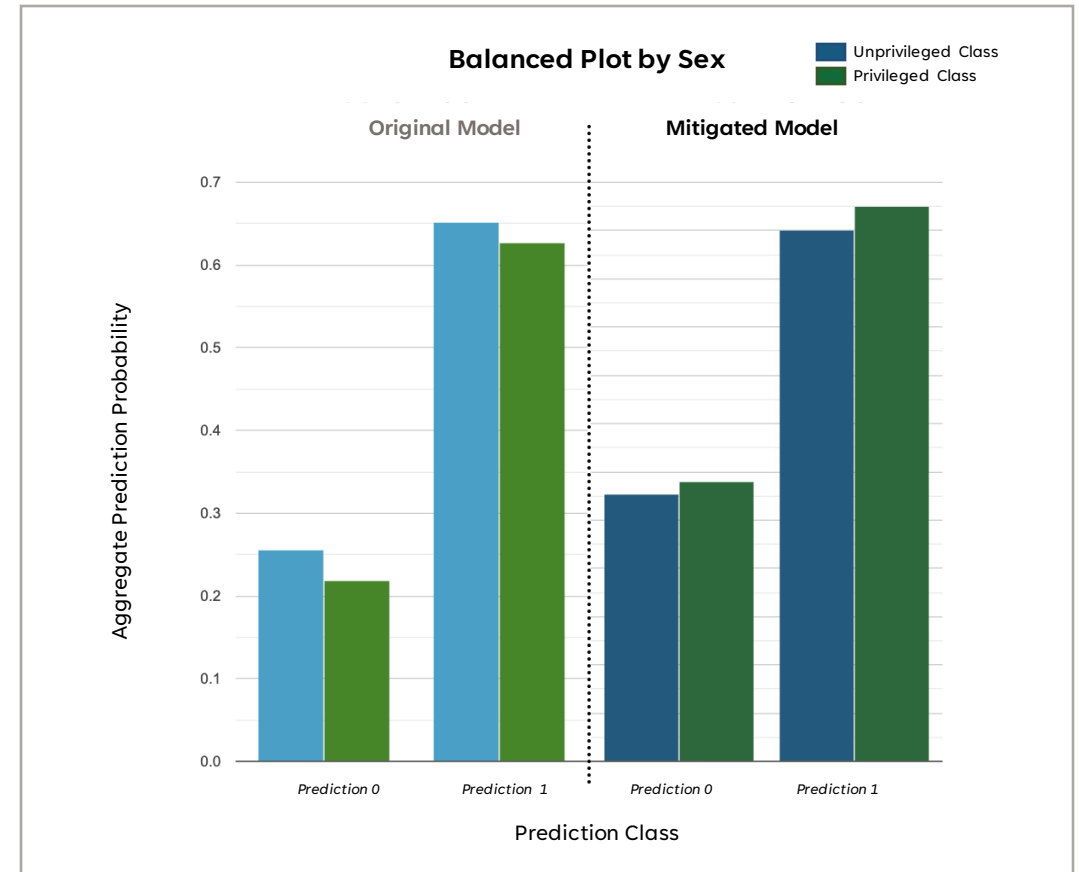
Four Fifths Plot

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



Balance Plot

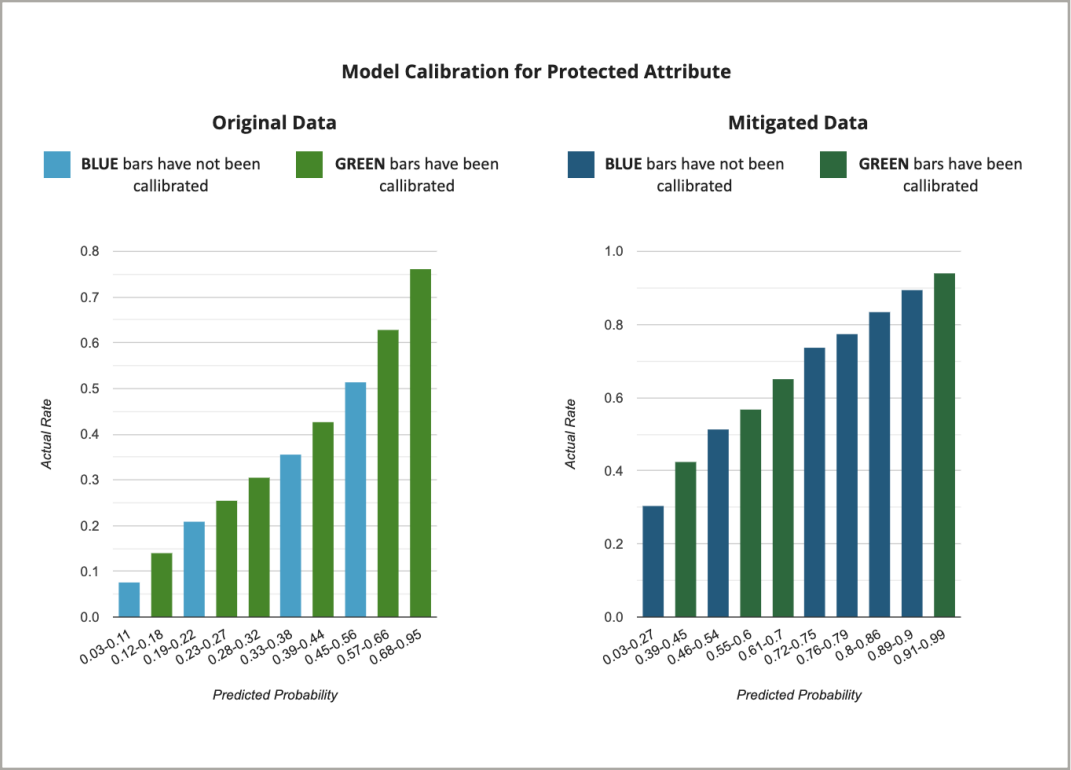
Examines whether average score received by individuals in positive and negative instances are similar regardless of sensitive attributes.



OTHER FAIRNESS VISUALIZATIONS (IMPUTED)

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.





DISCUSSION QUESTIONS

- Which fairness metrics do you think would be most valuable for our use case (predicting utilization)?
- Which fairness metrics are you still confused or concerned about? Why?
- What visualizations are you interested in creating, based off your EDA and understanding of the case?

“ZOOMING OUT” - Health Expenditure Use Case

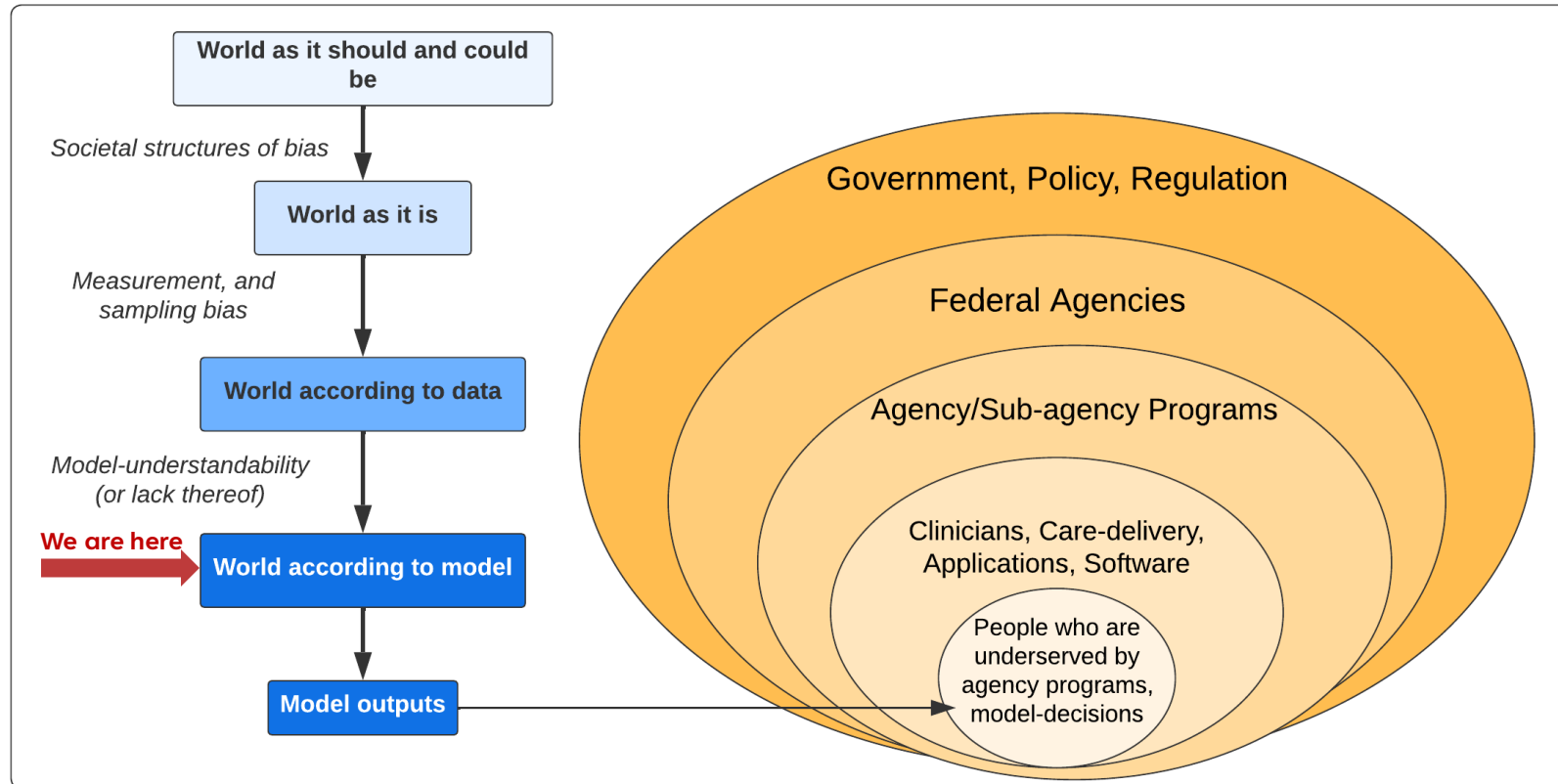


Figure above: Illustration of the types of bias that can enter a model's decision-space (stages of data and algorithm use), resulting in model-outputs that affect individuals downstream. (Adapted from: Mhasawadade 2020, Mitchell 2020)

Bias and Algorithmic Fairness

Use Cases

- Why are we implementing this? Who are affected parties?

Build

- What are the environmental costs of AI? Labor? Human Rights?

Deployment

- Once used, are the outcomes of the AI fair, equitable, and impartial? How can we measure this?

FOR NEXT WEEK

- Next week's class will be held as normal
- Complete next week's **readings**
 - If you signed up to present **Model Cards for Model Reporting (Mitchell et al.)** or **Fairness Through Awareness (Dwork et al.)** come prepared to present next week by 11 AM PT, Thursday, November 16th
- Submit your answers to next week's participation questions on Gradescope by 11 AM PT, Thursday, November 16th
- **Replication Part #2:** Notebook and writeup on model development and fairness metric assessments
 - Primary contact for replication project: Emily Ramond & Parker Addison
 - Office hours: Tuesdays 3:30-4:30pm
 - Start early! You'll need to learn/research fairness metrics.