# RESPONSIBLE AI

Week 6: Replication Project 2 – Fairness Assessments
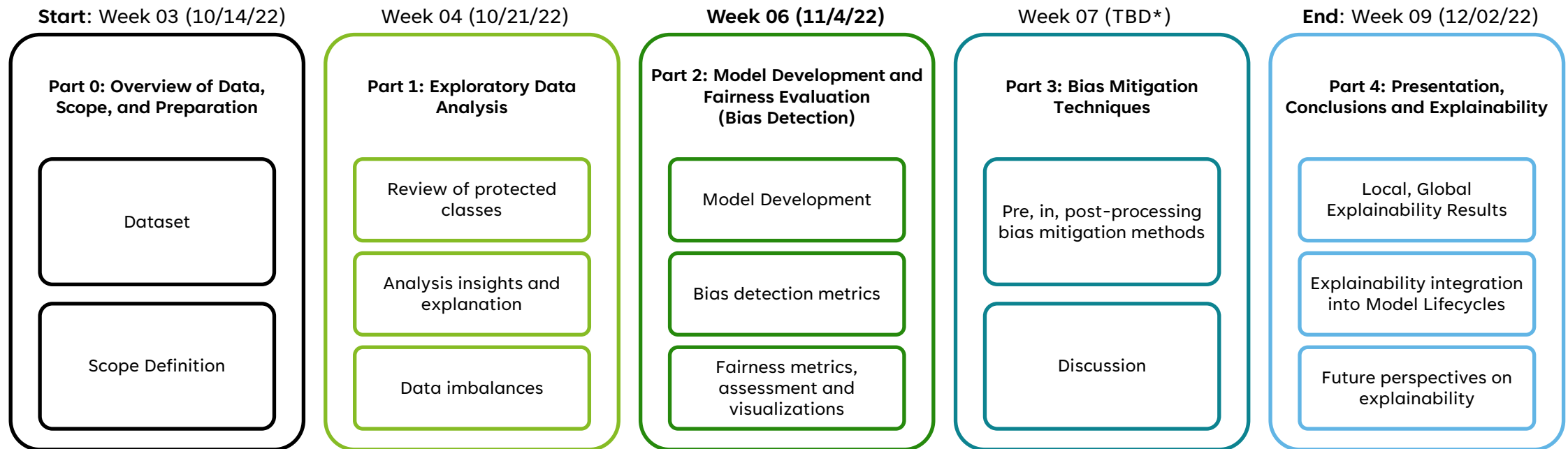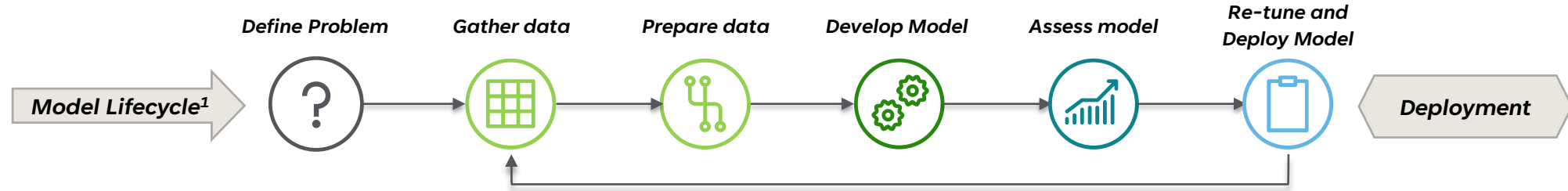
Meira Gilbert and Nandita Rahman

# 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?
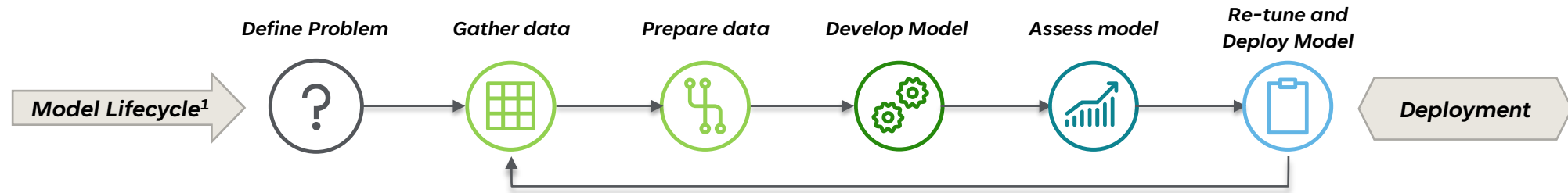
# REPLICATION PROJECT: RESPONSIBLE AI IN ACTION

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



Model Lifecycle[1]

Define Problem · Gather data · Prepare data · Develop Model · Assess model · Re-tune and Deploy Model · Deployment

| **Start**: Week 03 (10/14/22) | Week 04 (10/21/22) | **Week 06 (11/4/22)** | Week 07 (TBD*) | **End**: Week 09 (12/02/22) |
|---|---|---|---|---|
| **Part 0: Overview of Data, Scope, and Preparation** | **Part 1: Exploratory Data Analysis** | **Part 2: Model Development and Fairness Evaluation (Bias Detection)** | **Part 3: Bias Mitigation Techniques** | **Part 4: Presentation, Conclusions and Explainability** |
| Dataset | Review of protected classes | Model Development | Pre, in, post-processing bias mitigation methods | Local, Global Explainability Results |
| Scope Definition | Analysis insights and explanation | Bias detection metrics | Discussion | Explainability integration into Model Lifecycles |
| | Data imbalances | Fairness metrics, assessment and visualizations | | Future perspectives on explainability |

# REPLICATION PROJECT PART 02

## MODEL DEVELOPMENT & FAIRNESS EVALUATION

Define Problem | Gather data | Prepare data | Develop Model | Assess model | Re-tune and Deploy Model

Model Lifecycle[1]

Deployment

- **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?

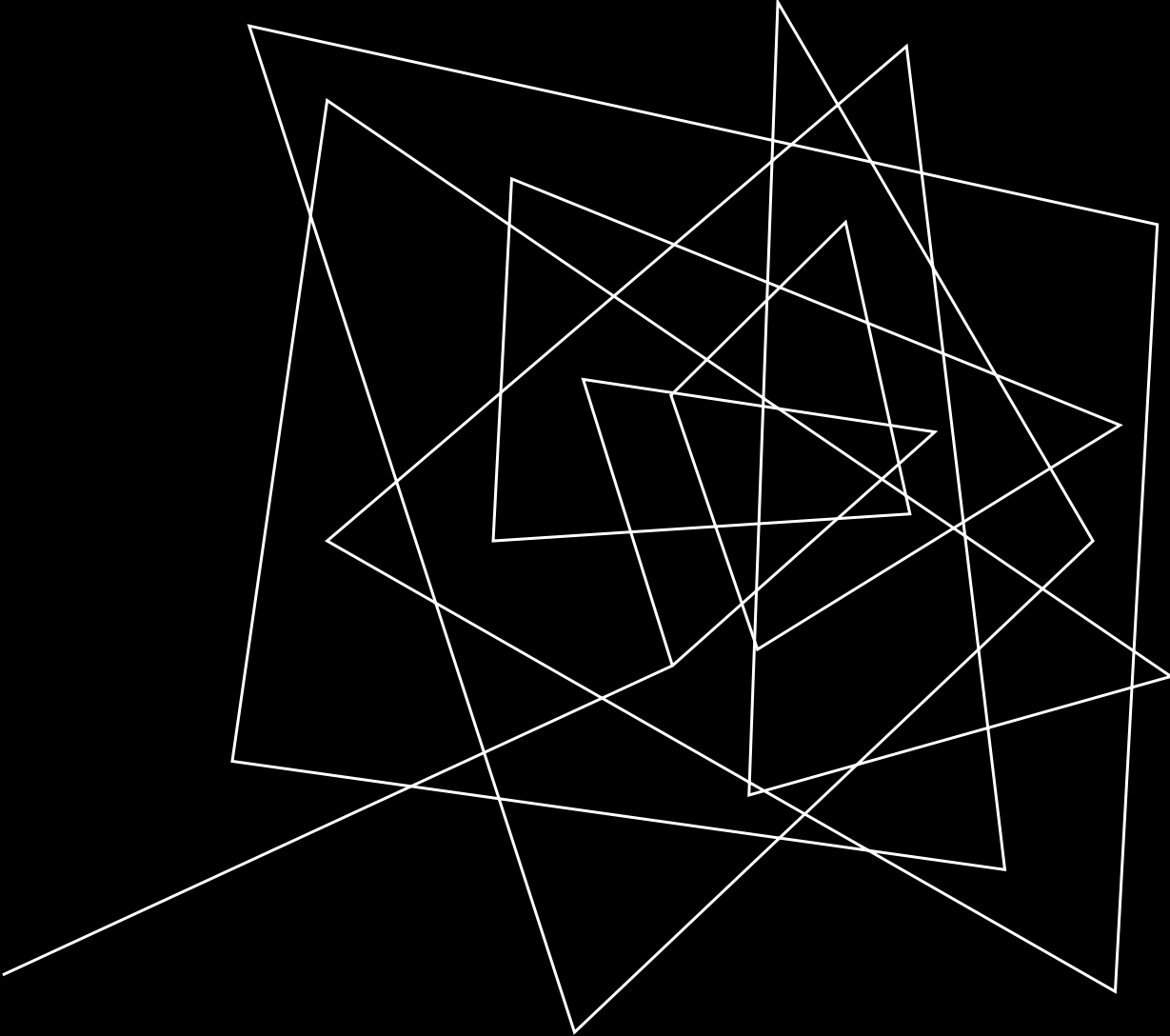Use visualizations to understand your model

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

### Assessing Fairness

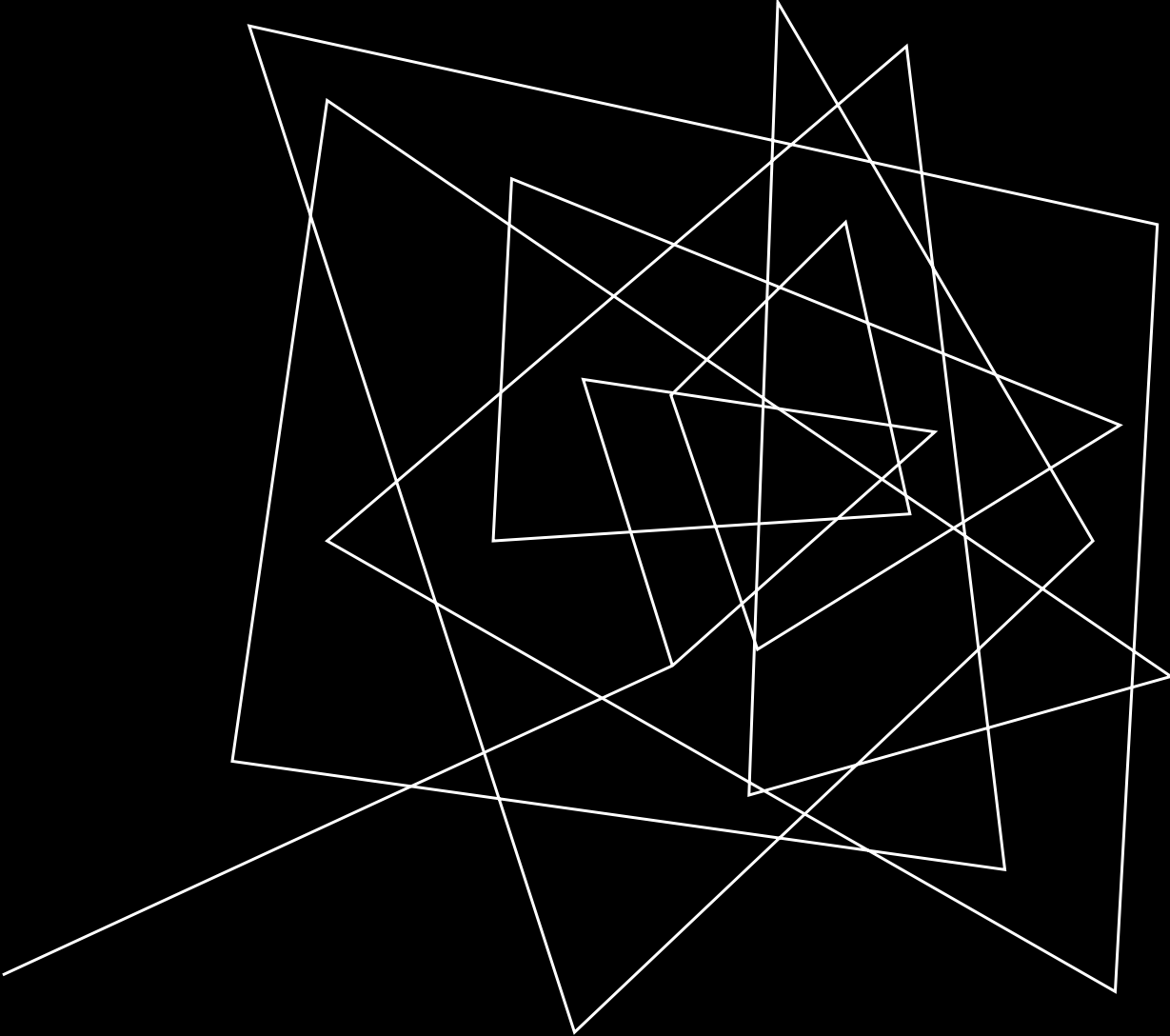Explain the fairness of your model predictions using your fairness metric results.

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)

# READING PRESENTATION #2

Millions of black people affected by racial bias in health-care algorithms (Nature)

# FAIRNESS METRICS

| | |
|---|---|
| **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?**

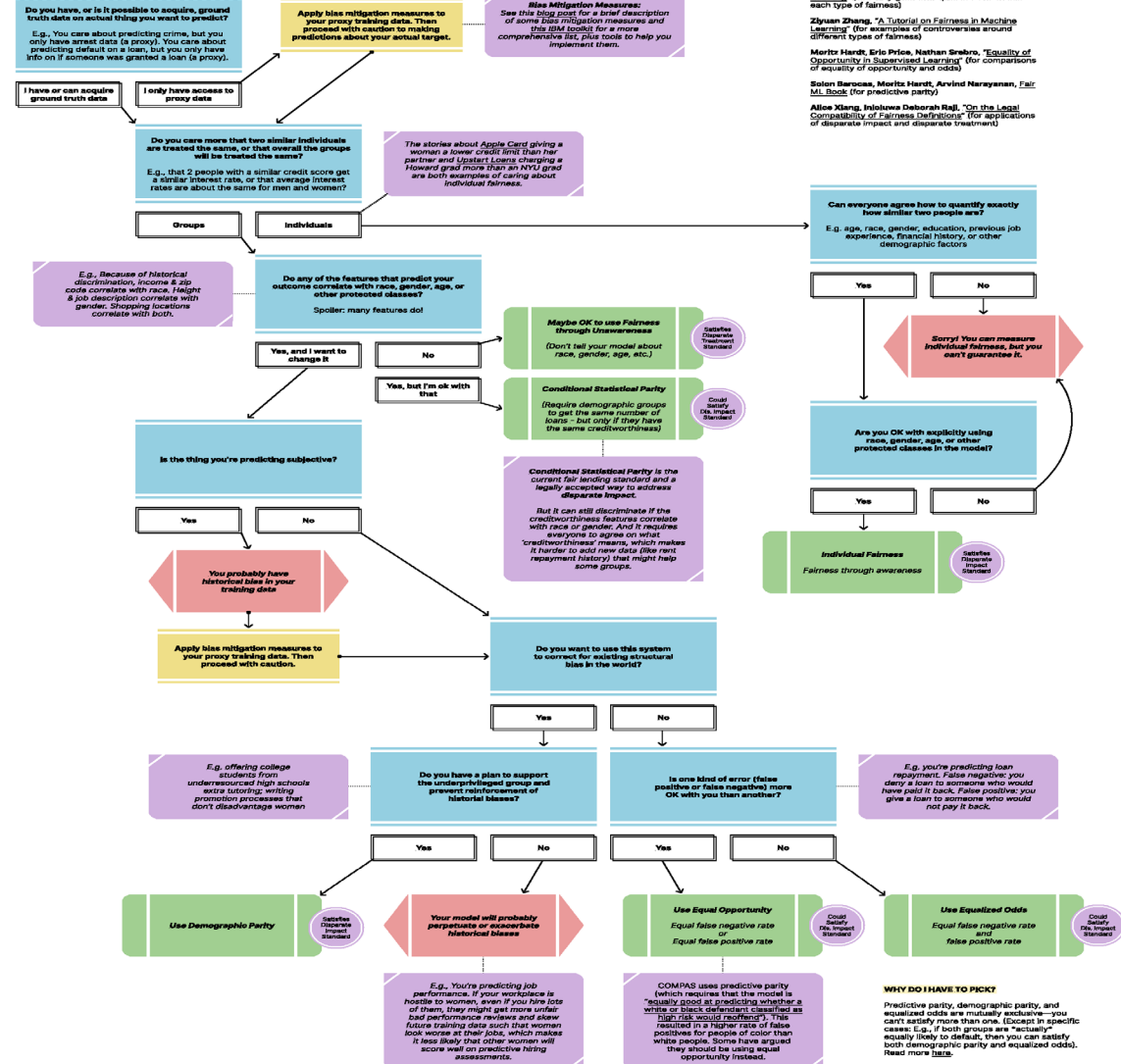**When you have competing fairness metrics, how to pick which to prioritize?**

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

# 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

**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)

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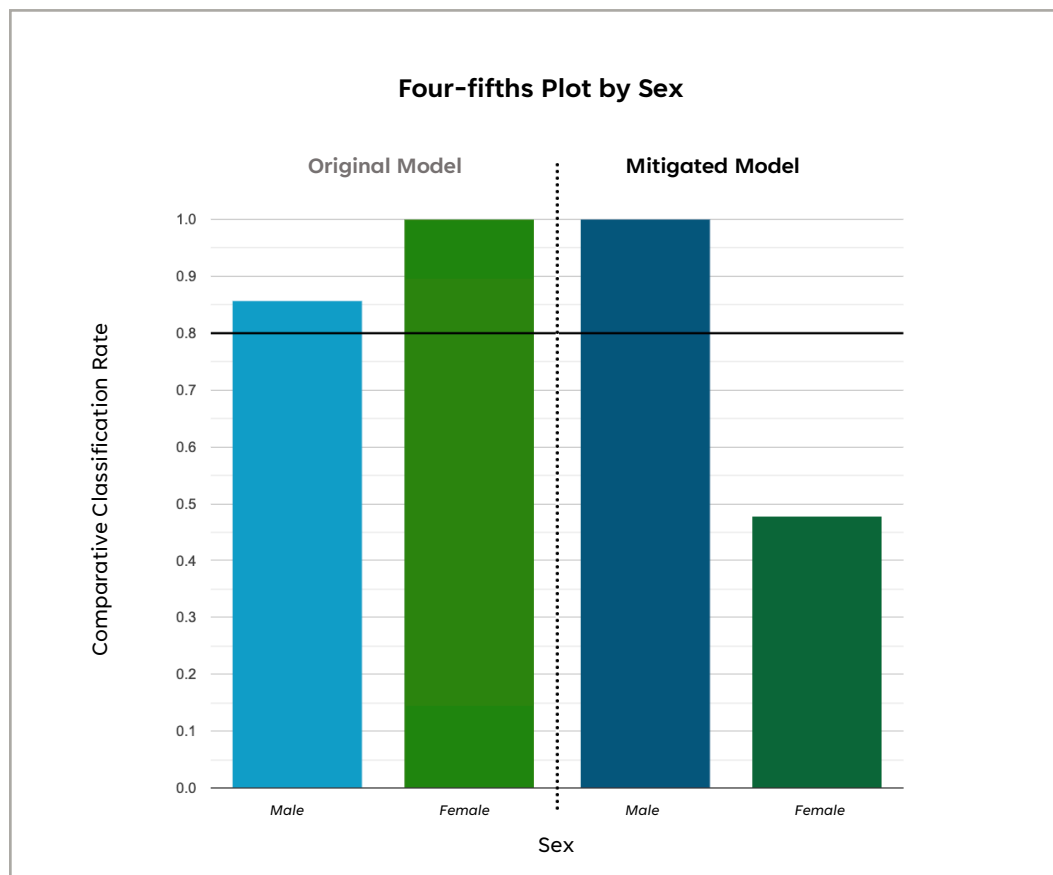
Solon Barocas, Moritz Hardt, Arvind Narayanan, Fair ML Book (for predictive parity)

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

https://www.aspentechpolicyhub.org/wp-content/uploads/2020/07/FAHL-Tree.pdf
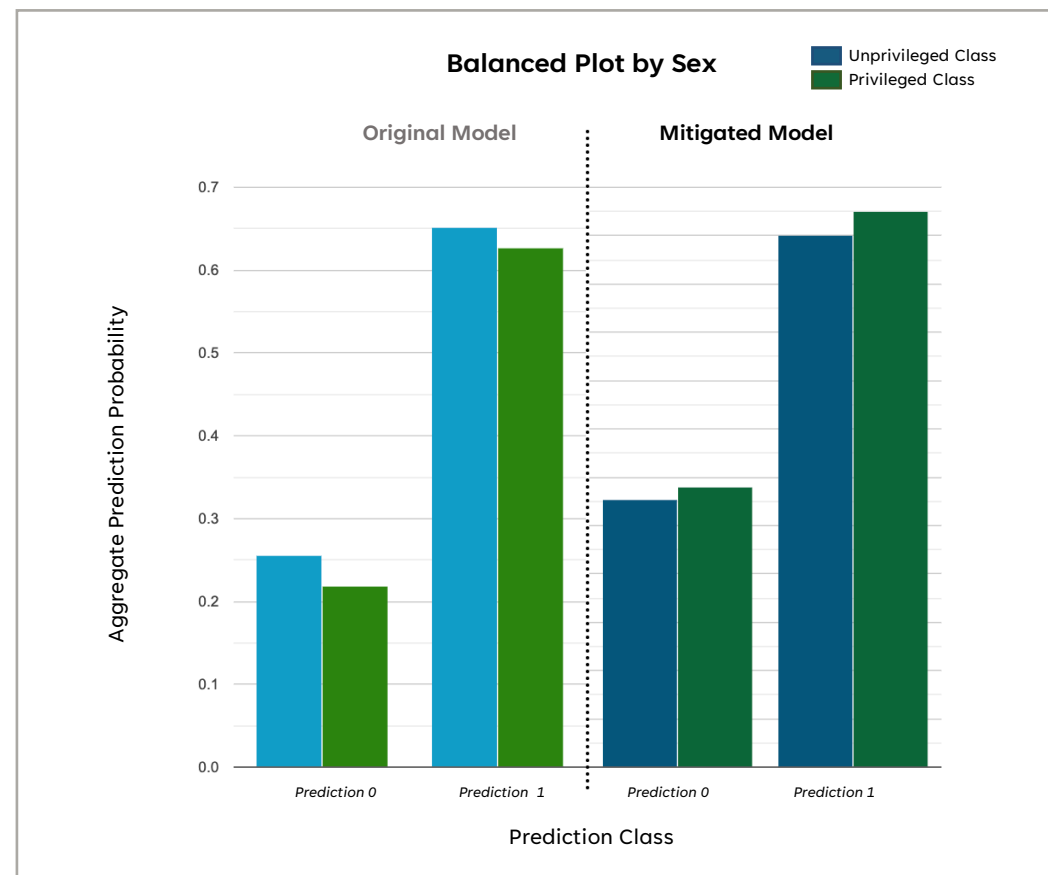
# FAIRNESS PLOT EXAMPLES

**Four Fifths Plot**

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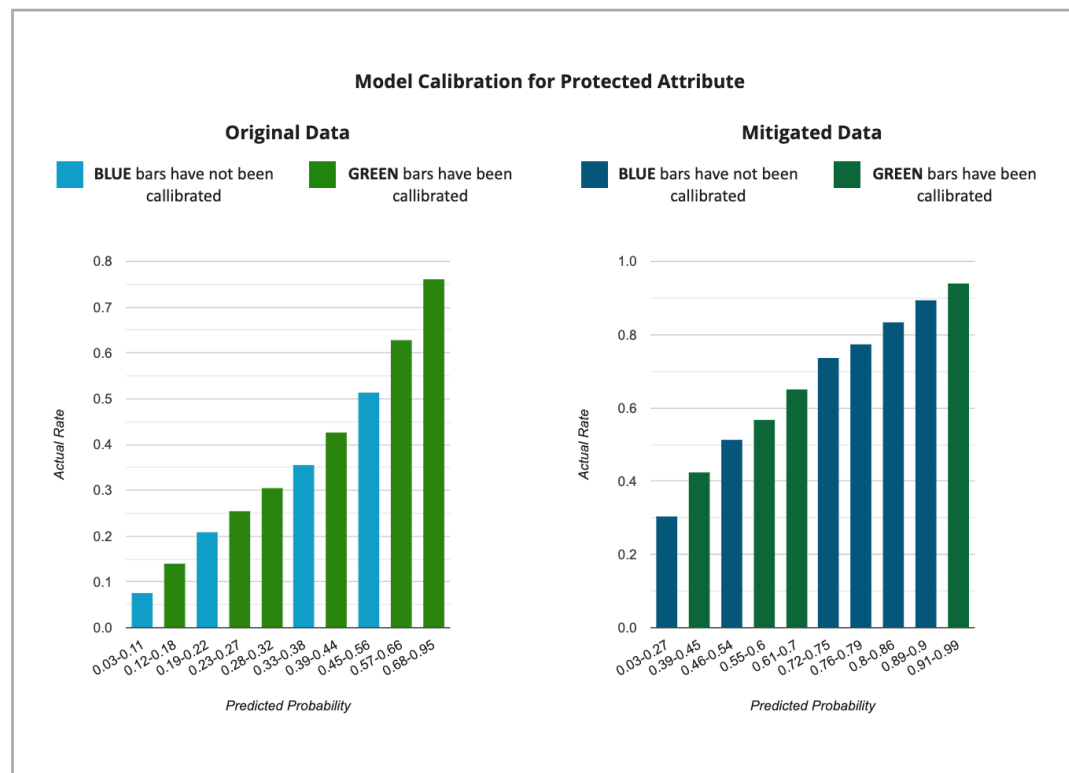
**Balance Plot**

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



Four-fifths Plot by Sex



Balanced Plot by Sex

# 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?

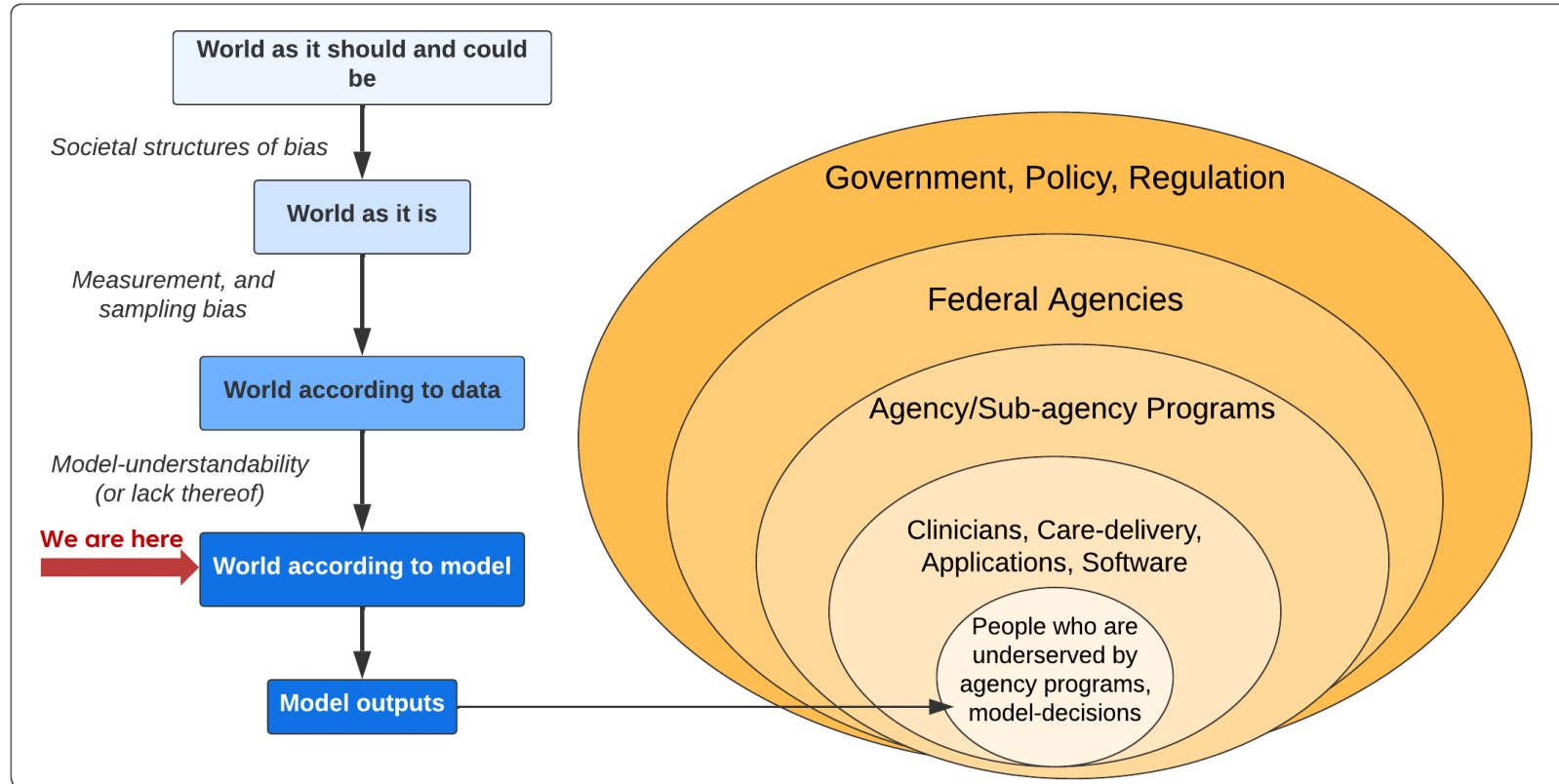# "ZOOMING OUT" – *Health Expenditure Use Case*



**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)*

# FOR NEXT WEEK

- Next week's class will be held **on THURSDAY, November 10th 12:30 PST!**
- 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 and submit your presentation to Gradescope by 10 AM PT, Thursday, November 10th

- Submit your answers to next week's participation questions on Gradescope by 10 PM PT, Thursday, November 10th

- **Replication Part #2:** Notebook and writeup on model development and fairness metric assessments
  - Primary contact for replication project: Nandita Rahman (nanrahman@deloitte.com)
  - Office hours: Mondays 1-2pm PST
  - Notebook link here

- Office Hours for Replication Project
  - Mondays from 1-2pm PT (4-5pm ET)
  - Nandita hosting virtually - Zoom link