# The AI/Data Science Professional

AI Bias and Fairness

# Acknowledgement of Country

The AI/Data Science Professional

Course Coordinator: Flora Salim

What's next...

RMIT UNIVERSITY

RMIT University acknowledges the people of the Woi wurrung and Boon wurrung language groups of the eastern Kulin Nation on whose unceded lands we conduct the business of the University.

RMIT University respectfully acknowledges their Ancestors and Elders, past and present. RMIT also acknowledges the Traditional Custodians and their Ancestors of the lands and waters across Australia where we conduct our business.



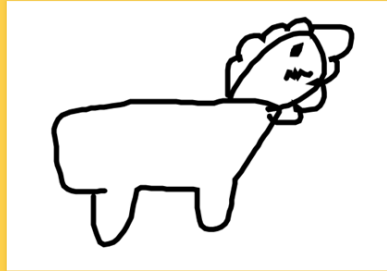Ngarara Place

# The AI/Data Science Professional – Week 6

AI Bias and Fairness

**What's next...**

RMIT
UNIVERSITY

# The Increasing Autonomy

# Real-world Examples

- COMPAS



VERNON PRATER
**Prior Offenses**
2 armed robberies, 1 attempted armed robbery
**Subsequent Offenses**
1 grand theft
LOW RISK **3**

BRISHA BORDEN
**Prior Offenses**
4 juvenile misdemeanors
**Subsequent Offenses**
None
HIGH RISK **8**

DYLAN FUGETT
LOW RISK **3**

BERNARD PARKER
HIGH RISK **10**

JAMES RIVELLI
LOW RISK **3**

ROBERT CANNON
MEDIUM RISK **6**

JAMES RIVELLI
**Prior Offenses**
1 domestic violence aggravated assault, 1 grand theft, 1 petty theft, 1 drug trafficking
**Subsequent Offenses**
1 grand theft
LOW RISK **3**

ROBERT CANNON
**Prior Offense**
1 petty theft
**Subsequent Offenses**
None
MEDIUM RISK **6**

# Real-world Examples

● COMPAS

# Real-world Examples

- Job Ads

# [Discussion]

Sometime ago, Verizon, one of the largest communication technology companies, placed an ad on Facebook to recruit applicants for a unit focused on financial planning and analysis. The ad showed a smiling, millennial-aged woman seated at a computer and promised that new hires could look forward to a rewarding career in which they would be "more than just a number." Some relevant numbers were not immediately evident. The promotion was set to run on the Facebook feeds of users 25 to 36 years old who lived in the nation's capital, or had recently visited there, and had demonstrated an interest in finance. For a vast majority of the hundreds of millions of people who check Facebook every day, the ad did not exist.

# Real-world Examples

- Facebook Career Ads

# Bias

# Bias in Data

# Types of Bias (Mehrabi et al. 2019)

23 types of bias !!
- Historical Bias
- Representation Bias
- Measurement Bias
- Evaluation Bias
- Aggregation Bias
- Population Bias
- Simpson's Paradox
- Longitudinal Data Fallacy
- Sampling Bias
- Behavioral Bias
- Content Production Bias

- Linking Bias
- Temporal Bias
- Popularity Bias
- Algorithmic Bias
- User Interaction Bias
- Social Bias
- Emergent Bias
- Self-Selection Bias
- Omitted Variable Bias
- Cause-Effect Bias
- Observer Bias
- Funding Bias

https://arxiv.org/pdf/1908.09635.pdf

# Common Types of Bias

- Historical bias

  ✤ Historical bias is the already existing bias and socio-technical issues in the world and can seep into from the data generation process even given a perfect sampling and feature selection.

  ✤ existing bias and socio-technical issues in the world, e.g. see below:



Source from https://arxiv.org/pdf/1908.09635.pdf

# Historical bias

- **Historical bias** occurs when the state of the world in which the data was generated is flawed.

- As of 2020, only <u>7.4%</u> of Fortune 500 CEOs are women. Research has shown that companies with female CEOs or CFOs are generally <u>more profitable</u> than companies with men in the same position, suggesting that women are held to higher hiring standards than men. In order to fix this, we might consider removing human input and using AI to make the hiring process more equitable. But this can prove unproductive if data from past hiring decisions is used to train a model, because the model will likely learn to demonstrate the same biases that are present in the data.

from <u>https://arxiv.org/pdf/1901.10002.pdf</u>

# Common Types of Bias

- Representation bias
  - Representation bias happens from the way we define and sample from a population feature selection.
  - Comes from the way we define and sample from a population, e.g.



ImageNet

| | |
|---|---|
| US | 45.4% |
| GB | 7.6% |
| IT | 6.2% |
| CA | 3.0% |
| AU | 2.8% |
| ES | 2.5% |
| AR | 1.0% |
| IE | 0.5% |
| CC | 0.0% |

Source from https://arxiv.org/pdf/1711.08536.pdf

# Representation bias

- **Representation bias** occurs when building datasets for training a model, if those datasets poorly represent the people that the model will serve.

- Data collected through smartphone apps will under-represent groups that are less likely to own smartphones. For instance, if collecting data in the USA, individuals over the age of 65 will be under-represented. If the data is used to inform design of a city transportation system, this will be disastrous, since older people have important needs to ensure that the system is accessible.

# Common Types of Bias

- Measurement bias
  - Measurement bias happens from the way we choose, utilize, and measure a particular feature.
  - Choosing and measuring the particular features of interest, e.g.,

# Measurement bias

**Measurement bias** occurs when the accuracy of the data varies across groups. This can happen when working with proxy variables (variables that take the place of a variable that cannot be directly measured), if the quality of the proxy varies in different groups.

# Measurement bias

- Your local hospital uses a model to identify high-risk patients before they develop serious conditions, based on information like past diagnoses, medications, and demographic data. The model uses this information to predict health care costs, the idea being that patients with higher costs likely correspond to high-risk patients. Despite the fact that the model specifically excludes race, it seems to demonstrate racial discrimination: the algorithm is less likely to select eligible Black patients. How can this be the case? It is because cost was used as a proxy for risk, and the relationship between these variables varies with race: Black patients experience increased barriers to care, have less trust in the health care system, and therefore have lower medical costs, on average, when compared to non-Black patients with the same health conditions.

# Common Types of Bias

- Aggregation bias
  - �֎ false conclusions are drawn for a subgroup based on observing other different subgroups or generally when false assumptions about a population affect the model's outcome and definition e.g.,

| Race/Ethnicity | Type 1 Diabetes Prevalence (per 1,000) | | Type 2 Diabetes Prevalence (per 1,000) | | Reference |
|---|---|---|---|---|---|
| | 0-9 years | 10-19 years | 0-9 years* | 10-19 years | |
| Non-Hispanic White | 1.03 | 2.89 | 0.0046 | 0.18 | [10] |
| Non-Hispanic Black | 0.57 | 2.04 | 0.0005 | 1.06 | [11] |
| Hispanic American | 0.44 | 1.59 | 0.0003 | 0.46 | [12] |
| Asian and Pacific Islanders | 0.26 | 0.77 | 0.014 | 0.52 | [13] |
| Native American | 0.08 | 0.28 | 0.021 | 1.45 | [14] |

Source from https://arxiv.org/pdf/1908.09635.pdf and https://www.ncbi.nlm.nih.gov/pubmed/22238408

# Aggregation bias

- **Aggregation bias** occurs when groups are inappropriately combined, resulting in a model that does not perform well for any group or only performs well for the majority group. (This is often not an issue, but most commonly arises in medical applications.)

- Hispanics have <u>higher rates</u> of diabetes and diabetes-related complications than non-Hispanic whites. If building AI to diagnose or monitor diabetes, it is important to make the system sensitive to these ethnic differences, by either including ethnicity as a feature in the data, or building separate models for different ethnic groups.

# Common Types of Bias

- Evaluation bias
  - Evaluation bias happens during model evaluation. This includes the use of inappropriate and disproportionate benchmarks for the evaluation of applications.
  - occurs during model iteration and evaluation, e.g.,





Source from
http://proceedings.mlr.press/v81/buolamwini18a/buolamwini18a.pdf

# Evaluation bias

- **Evaluation bias** occurs when evaluating a model, if the benchmark data (used to compare the model to other models that perform similar tasks) does not represent the population that the model will serve.

- The Gender Shades paper discovered that two widely used facial analysis benchmark datasets (IJB-A and Adience) were primarily composed of lighter-skinned subjects (79.6% and 86.2%, respectively). Commercial gender classification AI showed state-of-the-art performance on these benchmarks, but experienced disproportionately high error rates with people of color.

# Deployment bias

- **Deployment bias** occurs when the problem the model is intended to solve is different from the way it is actually used. If the end users don't use the model in the way it is intended, there is no guarantee that the model will perform well.

- The criminal justice system uses <u>tools</u> to predict the likelihood that a convicted criminal will relapse into criminal behavior. The predictions are <u>not designed for judges</u> when deciding appropriate punishments at the time of sentencing.

# Other Common Types of Bias

- **Population Bias.** Population bias arises when statistics, demographics, representatives, and user characteristics are different in the user population represented in the dataset or platform from the original target population.
- **Sampling Bias.** Sampling bias arises due to the non-random sampling of subgroups. As a consequence of sampling bias, the trends estimated for one population may not generalize to data collected from a new population.
- **Temporal Bias.** Temporal bias arises from differences in populations and behaviors over time.
- **Social Bias**. Social bias happens when other people's actions or content coming from them affects our judgment.

# Where do they exist?



(a) Data Generation

# Where do they exist?



(b) Model Building and Implementation

# Other classification and types of bias

- Population bias
- Simpson's Paradox
- Longitudinal data fallacy
- Sampling bias
- …...

**[Optional Reading]**

**<u>Social Data: Biases, Methodological Pitfalls, and Ethical Boundaries</u>**

# Algorithmic Fairness

# Algorithmic Fairness

- Types of Discrimination
  - ✂ Direct vs. Indirect discrimination
  - ✂ Systemic discrimination
  - ✂ Explainable vs. unexplainable discrimination
  - ✂ …...

# [Discussion]

Amazon's same-day delivery

Amazon recently rolled out same-day delivery across a select group of American cities. However, this service was only extended to neighborhoods with a high number of current Amazon users. As a result, predominantly non-white neighborhoods were largely excluded from the service.

# Algorithmic Fairness

- Definitions of Fairness
  - Equalised odds
  - Equal opportunity
  - Demographic parity
  - Fairness through awareness
  - Fairness through unawareness
  - Treatment equality
  - Test fairness
  - …...

# Algorithmic Fairness

- Categories of Fairness
  - Individual Fairness
  - Group Fairness
  - Subgroup Fairness

| Name | Group | Individual |
|---|---|---|
| Demographic parity | ✓ | |
| Conditional statistical parity | ✓ | |
| Equalized odds | ✓ | |
| Equal opportunity | ✓ | |
| Fairness through unawareness | | ✓ |
| Fairness through awareness | | ✓ |

# Fairness criteria: Demographic parity /statistical parity

- **Demographic parity** says the model is fair if the composition of people who are selected by the model matches the group membership percentages of the applicants.

- A nonprofit is organizing an international conference, and 20,000 people have signed up to attend. The organizers write a ML model to select 100 attendees who could potentially give interesting talks at the conference. Since 50% of the attendees will be women (10,000 out of 20,000), they design the model so that 50% of the selected speaker candidates are women.

# Fairness criteria: Equal opportunity

- **Equal opportunity** fairness ensures that the proportion of people who should be selected by the model ("positives") that are correctly selected by the model is the same for each group. We refer to this proportion as the **true positive rate** (TPR) or **sensitivity** of the model.

- A doctor uses a tool to identify patients in need of extra care, who could be at risk for developing serious medical conditions. (This tool is used only to supplement the doctor's practice, as a second opinion.) It is designed to have a high TPR that is equal for each demographic group.

# **Fairness criteria:** Equal accuracy

- Alternatively, we could check that the model has **equal accuracy** for each group. That is, the percentage of correct classifications (people who should be denied and are denied, and people who should be approved who are approved) should be the same for each group. If the model is 98% accurate for individuals in one group, it should be 98% accurate for other groups.

- A bank uses a model to approve people for a loan. The model is designed to be equally accurate for each demographic group: this way, the bank avoids approving people who should be rejected (which would be financially damaging for both the applicant and the bank) and avoid rejecting people who should be approved (which would be a failed opportunity for the applicant and reduce the bank's revenue).

# Fairness criteria: Group unaware / "Fairness through unawareness"

- **Group unaware** fairness removes all group membership information from the dataset. For instance, we can remove gender data to try to make the model fair to different gender groups. Similarly, we can remove information about race or age.

- One difficulty of applying this approach in practice is that one has to be careful to identify and remove proxies for the group membership data. For instance, in cities that are racially segregated, zip code is a strong proxy for race. That is, when the race data is removed, the zip code data should also be removed, or else the ML application may still be able to infer an individual's race from the data. Additionally, group unaware fairness is unlikely to be a good solution for historical bias.

# Example

- We'll work with a small example to illustrate the differences between the four different types of fairness. We'll use a **confusion matrix**, which is a common tool used to understand the performance of a ML model. This tool is depicted in the example below, which depicts a model with 80% accuracy (since 8/10 people were correctly classified) and has an 83% true positive rate (since 5/6 "positives" were correctly classified).

# Example (cont.)

- To understand how a model's performance varies across groups, we can construct a different confusion matrix for each group. In this small example, we'll assume that we have data from only 20 people, equally split between two groups (10 from Group A, and 10 from Group B).

# Example (cont.)

- The next image shows what the confusion matrices could look like, if the model satisfies **demographic parity** fairness. 10 people from each group (50% from Group A, and 50% from Group B) were considered by the model. 14 people, also equally split across groups (50% from Group A, and 50% from Group B) were approved by the model.

**Demographic parity**

20 applicants (50% from **Group A**)
14 approvals (50% from **Group A**)

**GROUP A**

**PREDICTED**

| TRUE | | Deny | Approve |
|---|---|---|---|
| | Deny | 1 | 2 |
| | Approve | 2 | 5 |

**GROUP B**

**PREDICTED**

| TRUE | | Deny | Approve |
|---|---|---|---|
| | Deny | 2 | 4 |
| | Approve | 1 | 3 |

# Example (cont.)

- For **equal opportunity** fairness, the TPR for each group should be the same; in the example below, it is 66% in each case.

**Equal opportunity**

**Group A:** 66% true positive rate: 4/(4+2)
**Group B:** 66% true positive rate: 2/(1+2)

**GROUP A**

**PREDICTED**

|  |  | Deny | Approve |
|---|---|---|---|
| **TRUE** | Deny | 3 | 1 |
|  | Approve | 2 | 4 |

**GROUP B**

**PREDICTED**

|  |  | Deny | Approve |
|---|---|---|---|
| **TRUE** | Deny | 6 | 1 |
|  | Approve | 1 | 2 |

# Example (cont.)

- Next, we can see how the confusion matrices might look for **equal accuracy** fairness. For each group, the model was 80% accurate.

**Equal accuracy**

**Group A**: 80% accurate: (6+2)/10
**Group B**: 80% accurate: (4+4)/10

# Example (cont.)

- Note that **group unaware** fairness cannot be detected from the confusion matrix, and is more concerned with removing group membership information from the dataset.

- Also note that none of the examples satisfy more than one type of fairness. For instance, the demographic parity example does not satisfy equal accuracy or equal opportunity.

To read more:
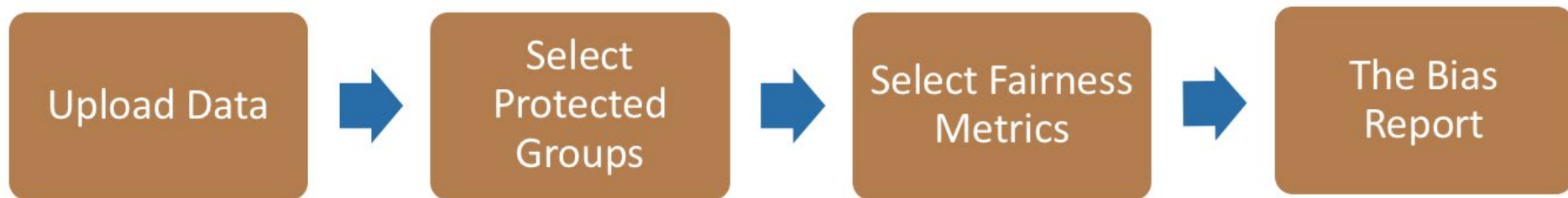https://arxiv.org/pdf/1710.03184.pdf

# What can we do about it?

# Accessing Fairness

Aequitas: Bias and Fairness Audit Toolkit



There is an example report on COMPASS
http://aequitas.dssg.io/example.html#audit-results-details-by-fairness-measures

# Accessing Fairness

## The AI Fairness 360 (AIF360 - IBM)



AI Fairness 360 - Demo

Data — Check — Mitigate — Compare

**2. Check bias metrics**

Dataset: Compas (ProPublica recidivism)
Mitigation: none

**Protected Attribute: Sex**

Privileged Group: *Female*, Unprivileged Group: *Male*
Accuracy with no mitigation applied is 66%
With default thresholds, bias against unprivileged group detected in 4 out of 5 metrics

| Statistical Parity Difference | Equal Opportunity Difference | Average Odds Difference | Disparate Impact | Theil Index |
|---|---|---|---|---|
| -0.36 | -0.3 | -0.35 | 0.59 | 0.21 |
| original | original | original | original | original |

# Mitigation



Original:

Balancing gender:

Balancing skin color:

Balancing age:

# Mitigation

**Optimized Pre-processing**

Use to mitigate bias in training data. Modifies training data features and labels.

→

**Reweighing**

Use to mitgate bias in training data. Modifies the weights of different training examples.

→

**Adversarial Debiasing**

Use to mitigate bias in classifiers. Uses adversarial techniques to maximize accuracy and reduce evidence of protected attributes in predictions.

→

**Reject Option Classification**

Use to mitigate bias in predictions. Changes predictions from a classifier to make them fairer.

→

**Disparate Impact Remover**

Use to mitigate bias in training data. Edits feature values to improve group fairness.

→

**Learning Fair Representations**

Use to mitigate bias in training data. Learns fair representations by obfuscating information about protected attributes.

→

**Prejudice Remover**

Use to mitigate bias in classifiers. Adds a discrimination-aware regularization term to the learning objective.

→

**Calibrated Equalized Odds Post-processing**

Use to mitigate bias in predictions. Optimizes over calibrated classifier score outputs that lead to fair output labels.

→

**Equalized Odds Post-processing**

Use to mitigate bias in predictions. Modifies the predicted labels using an optimization scheme to make predictions fairer.
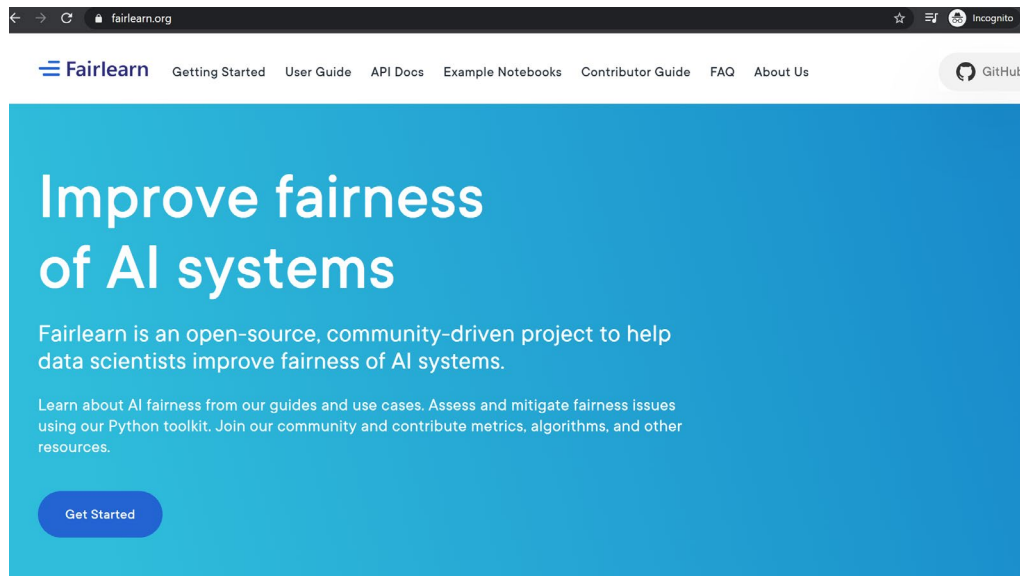
→

**Meta Fair Classifier**

Use to mitigate bias in classifier. Meta algorithm that takes the fairness metric as part of the input and returns a classifier optimized for that metric.
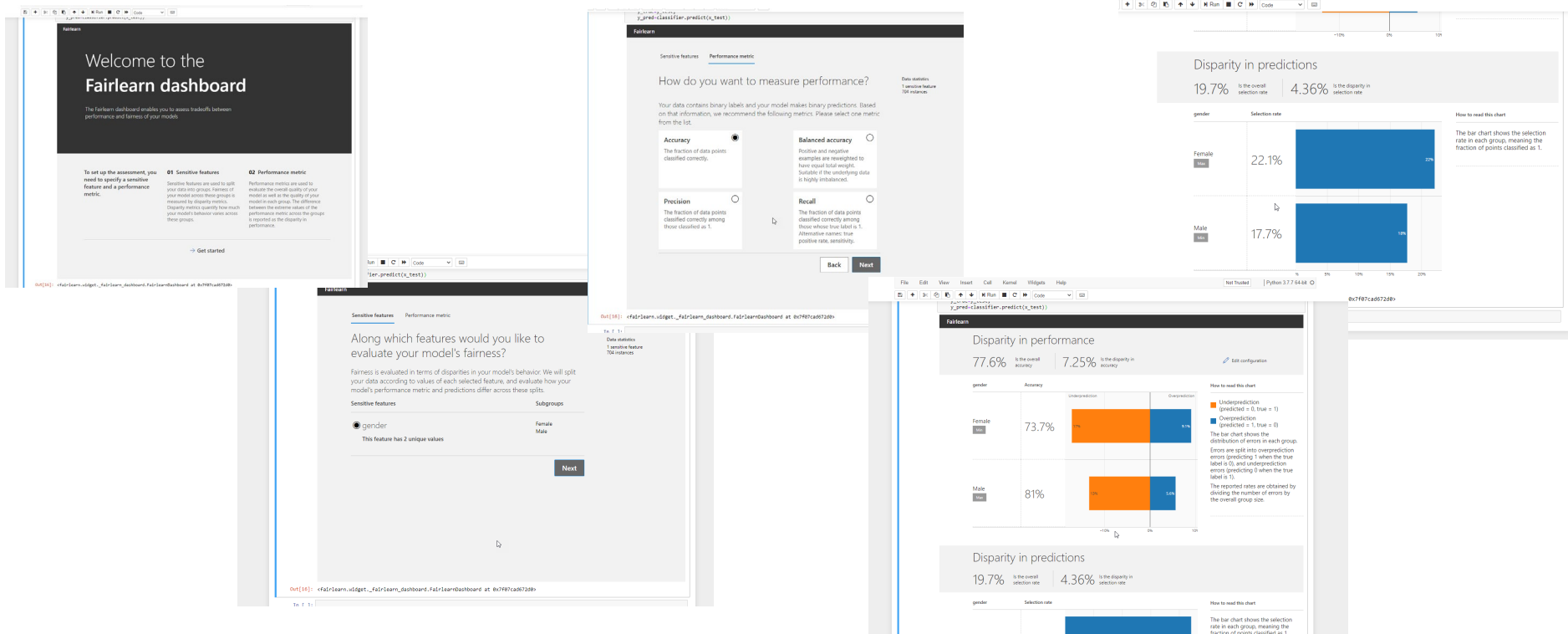
# Fairlearn (by Microsoft)



https://fairlearn.org/

https://www.microsoft.com/en-us/research/publication/fairlearn-a-toolkit-for-assessing-and-improving-fairness-in-ai/

# Fairlearn (by Microsoft)



https://fairlearn.org/

Credit: https://github.com/wmeints/fairlearn-demo
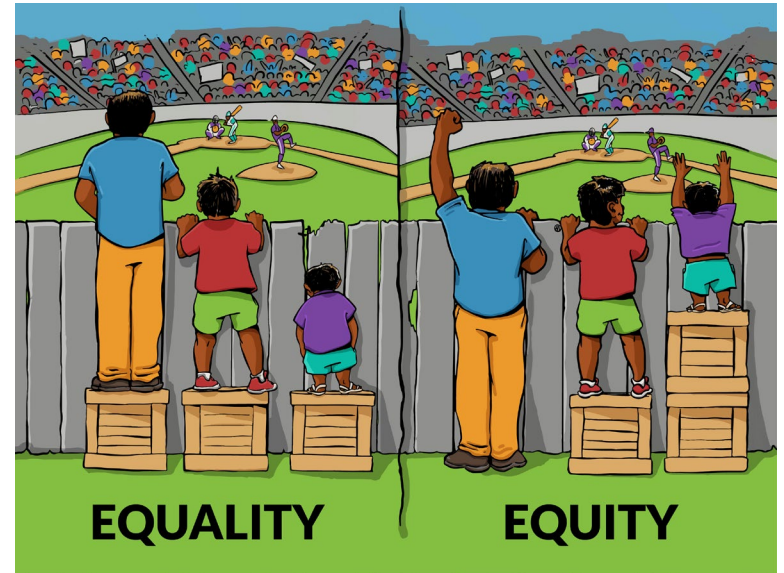
# Challenges?

# Challenges

- Synthesizing a definition of fairness.
- From Equality to Equity
- Searching for Unfairness



EQUALITY    EQUITY

https://interactioninstitute.org/illustrating-equality-vs-equity/

# What to do this week?

- Read materials

- Work on Task 1 (Case Study), due week 7 (next week)

- Must watch Guest Lecture by Kate Crawford

The Trouble with Bias - NIPS2017
https://youtu.be/fMym_BKWQzk

# Upcoming Weeks

- Week 7: Transparency and Explainability
- Week 8: Workshop – Presentation & Peer review of Task 1 (Case Study)
- Week 9: Guest Talk by A/Prof Richard Xu (**UTS**) on Data Science communication (stakeholder engagement) + Ethical ML as a Software Engineering enterprise (Charles Isbell's NeurIPS 2020 keynote)
- Week 10: Guest talk by Prof Milind Tambe (**Harvard**, Director "AI for Social Good", **Google Research India**) on "AI for social good"
- Week 11: Guest talk by Ron Tidhar (**Instagram**) on "Data Science and Product Development"
- Week 12: Workshop – Presentation & Peer review of Task 2 (AI/DS for Social Good project)

# References

[1] Mehrabi, Ninareh & Morstatter, Fred & Saxena, Nripsuta & Lerman, Kristina & Galstyan, Aram. (2019). A Survey on Bias and Fairness in Machine Learning.

[2] Harini Suresh and John V Guttag. 2019. A Framework for Understanding Unintended Consequences of Machine Learning. arXiv preprint arXiv:1901.10002 (2019).