


**CSE 312**

# **Foundations of Computing II**

**Lecture 27: Algorithmic Fairness**

# Agenda

- The context 
- Why fairness?
- What is “fair”?
- No fair lunch

## Algorithms deterministic and randomized

- Given:
  - a space of inputs  $\mathcal{X}$ 
    - possibly with an underlying distribution
  - A function  $f$  defined on  $\mathcal{X}$
- Find an algorithm  $A$  (possibly randomized) such that
  - $A(x) = f(x)$  for all  $x \in \mathcal{X}$
  - or  $|A(x) - f(x)|$  is small for all  $x \in \mathcal{X}$
  - or  $P_A(A(x) = f(x))$  is close to 1 for all  $x \in \mathcal{X}$
  - or  $P_x(A(x) = f(x))$  is close to 1 for  $x$  chosen randomly
  - or  $P_{A,x}(|A(x) - f(x)| \text{ is small})$  is close to 1 ...

## Joint Distributions and Prediction

- Underlying probability space  $(P, \Omega)$
- Joint distribution  $p_{Y, X_1, \dots, X_k}$  on  $\Omega$
- **Goal:** Predict  $Y$  given observations of  $X_1, \dots, X_k$
- **Examples:**
  - $Y$  = image classification, e.g. cat, facial match, stop sign, bicyclist
    - $X_1, \dots, X_k$  pixels of photo/video
  - $Y$  = success in college
    - $X_1, \dots, X_k$  might be high school GPA, SAT/ACT scores, personal statement, ZIP code, birthdate, gender, race/ethnicity, financial status/hardship.
  - $Y$  = likelihood of paying back a home loan
  - $Y$  = likelihood of re-offending if released

## Joint Distributions and Prediction

- Underlying probability space  $(P, \Omega)$
- Joint distribution  $p_{Y, X_1, \dots, X_k}$  on  $\Omega$
- **Goal:** Predict  $Y$  given observations of  $X_1, \dots, X_k$
- Algorithmic task:
  - Given  $X_1 = x_1, \dots, X_k = x_k$  produce a predicted value  $\hat{Y}$  for  $Y$
  - e.g.,  $Y$  is binary (+, -) “binary prediction”

## Joint Distributions and Prediction

- Underlying probability space  $(P, \Omega)$
- Joint distribution  $p_{Y, X_1, \dots, X_k}$  on  $\Omega$
- **Goal:** Predict  $Y$  given observations of  $X_1, \dots, X_k$
- What makes this different from typical algorithmic tasks?
  - The joint distribution  $p_{Y, X_1, \dots, X_k}$  is not explicitly available
    - Information is only accessible via *samples*  $(y, x_1, \dots, x_k)$ , a.k.a. *training data*
    - Samples may have missing data: e.g., attribute  $x_i$  might be noisy/not known. Even  $y$  might be missing (unlabeled examples).
  - The prediction function may not have access to all of  $x_1, \dots, x_k$
- What might constitute a good solution?

# Prediction Algorithms in Machine Learning (ML)

Prediction is one of the most common tasks in ML:

- Gather data from past about task [training data]
- Use an algorithm to find a model that (hopefully) makes accurate predictions [training / learning]
- Deploy the model to the “real world” and have it make predictions about new data [prediction / inference]

Traditionally, model that is “most accurate” is the one selected.

- Very difficult to find a perfect model in practice, mistakes are expected!


# Agenda

- The context
- Why fairness? ◀
- What is “fair”?
- No fair lunch



# Examples of problems with ML Systems



**INDEPENDENT**  
**GOOGLE'S ALGORITHM SHOWS  
PRESTIGIOUS JOB ADS TO MEN,  
BUT NOT TO WOMEN**



**REUTERS** Business Markets World Politics TV More

BUSINESS NEWS OCTOBER 9, 2018 / 8:12 PM / 5 MONTHS AGO

## Amazon scraps secret AI recruiting tool that showed bias against women

Jeffrey Dastin 8 MIN READ  

**The New York Times**

## Facebook Engages in Housing Discrimination With Its Ad Practices, U.S. Says

By Katie Benner, Glenn Thrush and Mike Isaac

March 28, 2019       168

## MIT Technology Review

### Intelligent Machines

# How to Fix Silicon Valley's Sexist Algorithms

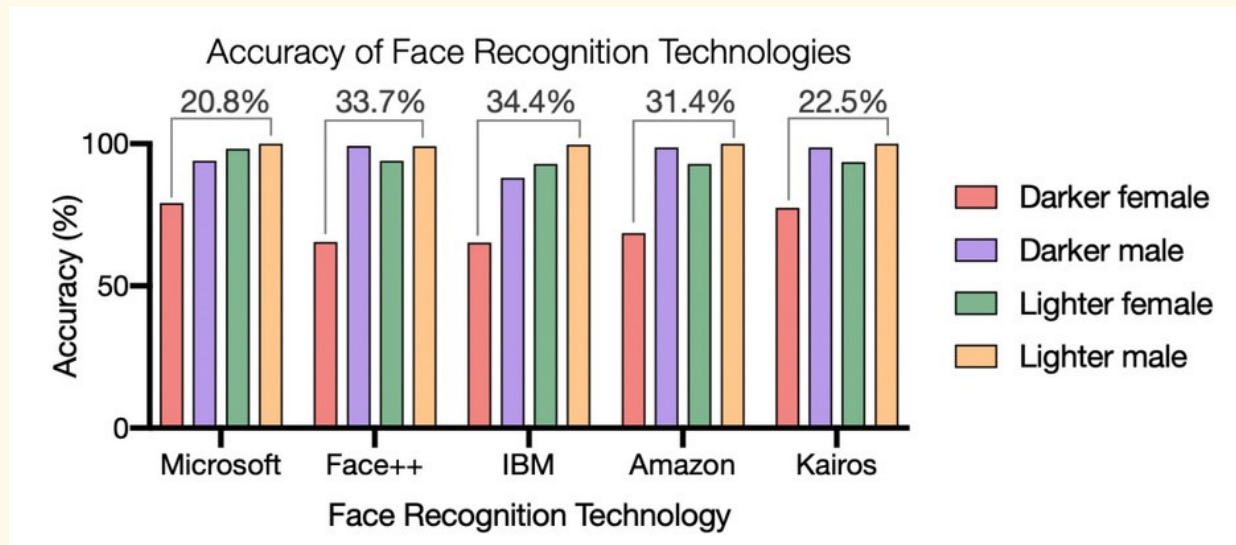
Computers are inheriting gender bias implanted in language data sets—and not everyone thinks we should correct it.

**PRO PUBLICA**

# Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

# Examples of problems with ML Systems



## A common theme...

Rates of predictions  $\{+, -\}$  vary based on factors that individuals don't control, such as race/ethnicity, gender, age, in ways that seem to produce less desirable outcomes for them more frequently if they have certain attributes...

... despite the fact that, in many settings (e.g., credit scores, housing, employment) there are clear legal protections against differential treatment based on some of these same 'protected' attributes.

## Why Different Treatment Happens?

More often than not, programmers aren't *intentionally* coding up racism or sexism into their models. Where does this difference in treatment come from?

Most often: the data

- Training a model on biased data will likely result in a biased model
  - Training data may not be representative of the actual distribution
    - Absence of data may reflect historical patterns and biases
    - Small sample sizes do not generalize as easily
  - Training data may include impacts of explicit biases

## Goal: A mathematical definition of fairness

Can we come up with a mathematical definition for what it means for the outcomes of a model to be fair/unfair?

- If we have a mathematical definition to spot when a model will be unfair, we can avoid deploying it before it can potentially harm people.

**Key Idea:** There is not going to be one definition of fairness. Every definition encodes its own set of values for what fairness means. Because there is not one universal definition, which one you choose should be carefully selected based on your values.

# Agenda

- The context
- Why fairness?
- What is “fair”? ◀
- No fair lunch

## Group Fairness

Today, we are going to be limiting our discussions to notions of **group fairness** which is also called **non-discrimination**.

- **Goal:** Avoid differential treatment based on membership in some protected group.

A **group** is some defined shared attribute between individuals

- Examples: race/ethnicity, gender, age, etc.
- Can also consider defining groups as combinations of the above
  - a.k.a. “intersectionality”

## A Toy Example: College Admissions

### Assumptions (unrealistic):

- There is a definition of “success” for college applicants, and the goal of an admissions decision is a prediction of “success”
- The only thing we will use as part of our decision is SAT Score
- To talk about group fairness, will assume everyone belongs to exactly one of two groups: Circles ( $2/3$ ) or Squares ( $1/3$ ).



## Notation

Example: College admissions only using SAT score

$X$  input about a person for prediction

– Example:  $X$  = SAT Score

$A$  variable indicating which group  $X$  belongs in

– Example:  $A = \square$  or  $A = \bigcirc$

$Y$  the “true label”

– Example:  $Y = +$  if truly successful in college,  $Y = -$  if not

$\hat{Y} = f(X)$  is our prediction for  $Y$  using a learned model  $f$

– Example:  $\hat{Y} = +$  if predicted successful,  $\hat{Y} = -$  otherwise

## First Attempt: “Shape-blind”

First attempt: To avoid unfair decisions, prevent the model from ever seeing the protected attribute (e.g., if an applicant is Circle/Square).

**Doesn't work:** In the real world, many things are correlated with group. Protected attribute can be unintentionally inferred from many other sources.

## Fairness Definition 1: **Statistical Parity**

**Idea:** “Admit decisions are equivalent across groups.”

E.g., Suppose  $P(Y = + | A = \square) = P(Y = + | A = \circ)$  then

$$P(\hat{Y} = + | A = \square) = P(\hat{Y} = + | A = \circ)$$

### **Pros:**

- Aligns with certain legal definitions of equity

### **Cons:**

- Rather weak in requirements (self-fulfilling prophecy)

## Types of Mistakes

A weakness of Statistical Parity comes from the fact it doesn't care about the true labels.

A stronger definition of fairness could require that the types of mistakes we make across groups are equal.

$$\begin{aligned} \text{True Positive Rate} &= \frac{\text{True Positive Prediction}}{\text{Positive}} \\ &= \frac{TP}{FN + TP} = 1 - \text{False Negative Rate} \end{aligned}$$

$$\begin{aligned} \text{True Negative Rate} &= \frac{\text{True Negative Prediction}}{\text{Negative}} \\ &= \frac{TN}{FP + TN} = 1 - \text{False Positive Rate} \end{aligned}$$

		prediction outcome		total
		<i>p</i>	<i>n</i>	
actual value	<i>p'</i>	True Positive	False Negative	<i>P'</i>
	<i>n'</i>	False Positive	True Negative	<i>N'</i>
total		<i>P</i>	<i>N</i>	

## Fairness Definition 2: Equal Opportunity

**Idea:** The true positive rate should be equivalent across groups.

$$P(\hat{Y} = + | A = \square, Y = +) = P(\hat{Y} = + | A = \circ, Y = +)$$

### Pros:

- Better controls for true outcome

### Cons:

- More complex to explain to non-experts
- Really only works in scenarios where there is a well defined  $Y = +$

## And many, many more

List of demographic fairness criteria			
Name	Closest relative	Note	Reference
Statistical parity	Independence	Equivalent	Dwork et al. (2011)
Group fairness	Independence	Equivalent	
Demographic parity	Independence	Equivalent	
Conditional statistical parity	Independence	Relaxation	Corbett-Davies et al. (2017)
Darlington criterion (4)	Independence	Equivalent	Darlington (1971)
Equal opportunity	Separation	Relaxation	Hardt, Price, Srebro (2016)
Equalized odds	Separation	Equivalent	Hardt, Price, Srebro (2016)
Conditional procedure accuracy	Separation	Equivalent	Berk et al. (2017)
Avoiding disparate mistreatment	Separation	Equivalent	Zafar et al. (2017)
Balance for the negative class	Separation	Relaxation	Kleinberg, Mullainathan, Raghavan (2016)
Balance for the positive class	Separation	Relaxation	Kleinberg, Mullainathan, Raghavan (2016)
Predictive equality	Separation	Relaxation	Chouldechova (2016)
Equalized correlations	Separation	Relaxation	Woodworth (2017)
Darlington criterion (3)	Separation	Relaxation	Darlington (1971)
Cleary model	Sufficiency	Equivalent	Cleary (1966)
Conditional use accuracy	Sufficiency	Equivalent	Berk et al. (2017)
Predictive parity	Sufficiency	Relaxation	Chouldechova (2016)
Calibration within groups	Sufficiency	Equivalent	Chouldechova (2016)
Darlington criterion (1), (2)	Sufficiency	Relaxation	Darlington (1971)

## Which definition to use?

We can't tell you! Each definition is a different take on what fairness means. Choosing a fairness measure is an explicit decision about of what values we hold when thinking about fairness.

**Takeaway:** Discrimination in ML models is not a problem that will only be solved algorithmically. We need people (e.g., policymakers, regulators, philosophers, developers) to be in the loop to determine the values we want to encode in the system.

## Brain Break





# Agenda

- The context
- Why fairness?
- What is “fair”?
- **No fair lunch** ◀

## (Im)possibility of Fairness

Four reasonable conditions we want in a real world ML Model:

1. Statistical Parity
2. Equality across false negative rates
3. Equality across false positive rates
4. Good accuracy of the model across subgroups

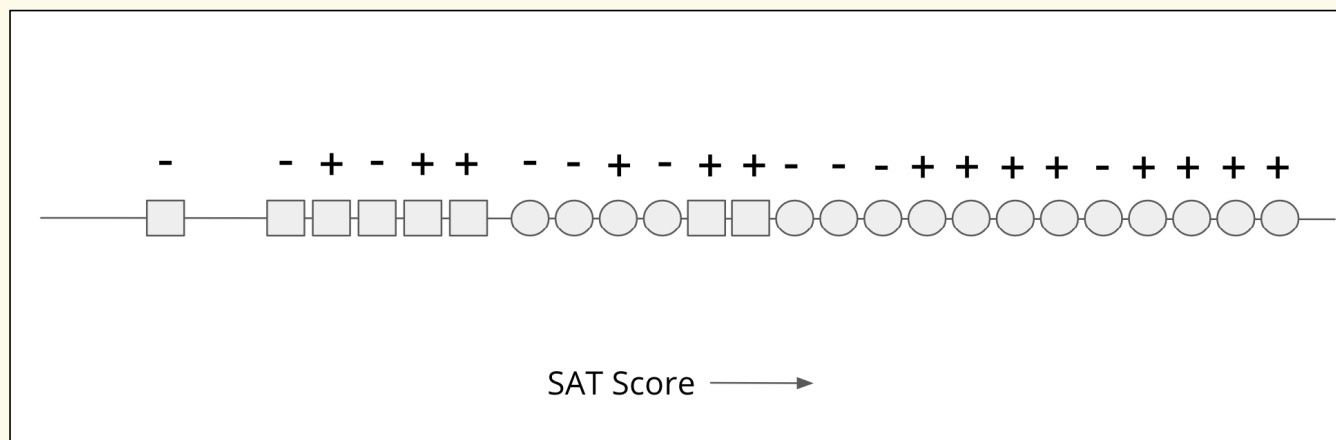
In general, can't satisfy all 4 simultaneously, *unless* groups have the exact same underlying distribution.

- This condition is rarely met in practice as we mentioned earlier when there are so many places for bias to enter our data collection.

## Example

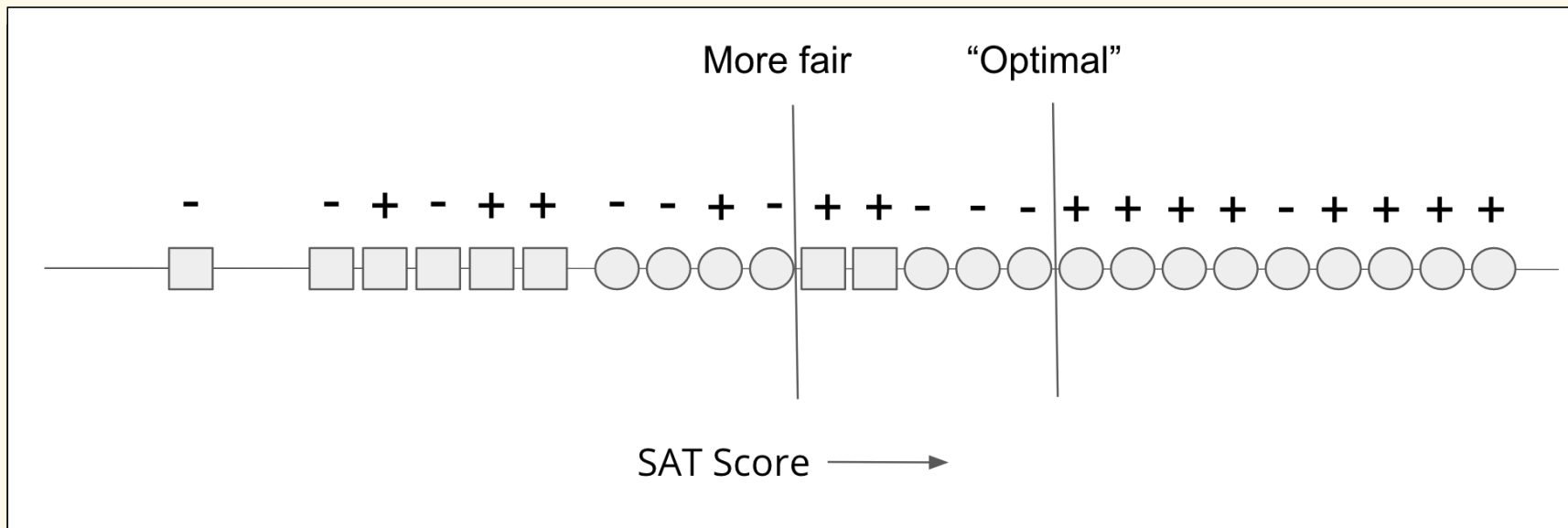
Consider a simplistic college admissions example, with a fake dataset.

- Majority (2/3) are Circle, the remaining 1/3 are Square
- Within each group, SAT is predictive, but score for Circles tends to be shifted higher when compared to Squares. E.g., Access to SAT Prep, Test retakes
- Even though we see statistical differences between groups in our data, the rate in which they are actually successful is the same.



## Fairness-Accuracy Tradeoff

In general, we find there is a tradeoff between accurate models and fair models. Making a model more fair tends to decrease accuracy by some amount.



## Notes on the Tradeoff

Overly simplistic but gives a sense of the complexity. Lots of examples of “accurate” models are indeed unfair.

This is not a statement that a tradeoff must exist, it just generally happens in real-world datasets.

- Originally just cared about finding the most accurate model, saw unfairness as a byproduct. Controlling for fairness will yield a different model than you found before.
- If data can encode biases and accuracy is determined in terms of that biased data, trying to achieve fairness will likely hurt accuracy.



## On the Pareto Frontier

This feels a bit cold-hearted, it's okay to like this is weird. Michael Kearns and Aaron Roth write in *The Ethical Algorithm*

While the idea of considering cold, quantitative trade-offs between accuracy and fairness might make you uncomfortable, the point is that there is simply no escaping the Pareto frontier. Machine learning engineers and policymakers alike can be ignorant of it or refuse to look at it. But once we pick a decision-making model (which might in fact be a human decision-maker), there are only two possibilities. Either that model is not on the Pareto frontier, in which case it's a “bad” model (since it could be improved in at least one measure without harm in the other), or it is on the frontier, in which case it implicitly commits to a numerical weighting of the relative importance of error and unfairness. Thinking about fairness in less quantitative ways does nothing to change these realities—it only obscures them.

Making the trade-off between accuracy and fairness quantitative does **not** remove the importance of human judgment, policy, and ethics—it simply focuses them where they are most crucial and useful, which is in deciding exactly which model on the Pareto frontier is best (in addition to choosing the notion of fairness in the first place, and which group or groups merit protection under it, [...]). Such decisions should be informed by many factors that cannot be made quantitative, including what the societal goal of protecting a particular group is and what is at stake. Most of us would agree that while both racial bias in the ads users are shown online and racial bias in lending decisions are undesirable, the potential harms to individuals in the latter far exceed those in the former. So in choosing a point on the Pareto frontier for a lending algorithm, we might prefer to err strongly on the side of fairness—for example, insisting that the false rejection rate across different racial groups be very nearly equal, even at the cost of reducing bank profits. We'll make more mistakes this way—both false rejections of creditworthy applicants and loans granted to parties who will default—but those mistakes will not be disproportionately concentrated in any one racial group.

## Recap

- ML (and human) systems can exhibit unfair behavior
- There are tools to define mathematically how to spot if a model is unfair. There are many definitions, and which you choose is a decision about values.
- In practice, there is generally a tradeoff between the fairness and accuracy of an ML model. Viewing the “Pareto Frontier” can help you visualize what this tradeoff looks like, but deciding which tradeoff is appropriate is yet another decision about values.

This is not something that will be solved “magically” by technology.