

# Trustworthy Machine Learning

## Fairness in Learning 1

**Sangdon Park**

POSTECH

# Contents from

## FAIRNESS AND MACHINE LEARNING

### Limitations and Opportunities

Solon Barocas, Moritz Hardt, Arvind Narayanan

## Equality of Opportunity in Supervised Learning

Moritz Hardt      Eric Price      Nathan Srebro

October 11, 2016

### Abstract

We propose a criterion for discrimination against a specified sensitive attribute in supervised learning, where the goal is to predict some target based on available features. Assuming data about the predictor, target, and membership in the protected group are available, we show how to optimally *adjust* any learned predictor so as to remove discrimination according to our definition. Our framework also improves incentives by shifting the cost of poor classification from disadvantaged groups to the decision maker, who can respond by improving the classification accuracy.

In line with other studies, our notion is *oblivious*: it depends only on the joint statistics of the predictor, the target and the protected attribute, but not on interpretation of individual features. We study the inherent limits of defining and identifying biases based on such oblivious measures, outlining what can and cannot be inferred from different oblivious tests.

We illustrate our notion using a case study of FICO credit scores.

- and contents partially from slides by Roger Grosse at University of Toronto.

# Why Fairness in Learning?

The image shows two screenshots of the Google Translate interface. The top screenshot shows the source text "She is a doctor. He is a nurse." being translated into Turkish as "O bir doktor. O bir hemşire." The bottom screenshot shows the source text "O bir doktor. O bir hemşire" being translated back into English as "He is a doctor. She is a nurse". This illustrates how the translation process introduces gender bias.

English Turkish Spanish Detect language ▾ English Turkish Spanish ▾ Translate

She is a doctor.  
He is a nurse. 31/5000

O bir doktor.  
O bir hemşire.

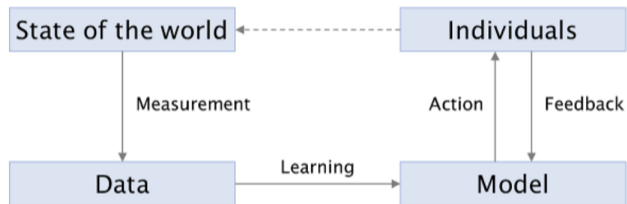
English Turkish Spanish Turkish - detected ▾ English Turkish Spanish ▾ Translate

O bir doktor.  
O bir hemşire 28/5000

He is a doctor.  
She is a nurse ✓

- Translation from English to Turkish, then back to English injects gender bias.

# Why Fairness in Learning?



- The machine learning loop
- Biased models enforce the bias of the world.

# Fairness in Learning: Overview

## Goal

Identify and mitigate “bias” in ML-based decision making.

## Source of bias:

- Data
  - ▶ imbalanced data (e.g., rare data, gender-biased data)
  - ▶ incorrect data (e.g., noisy data, data with historical bias)
- Model
  - ▶ modeling error
  - ▶ bias in loss

Credit: Richard Zemel

# Fairness in Learning: Definitions

- Known definitions
  - ▶ Demographic parity
  - ▶ Equalized odds
  - ▶ Equal opportunity
  - ▶ Equal (weak) calibration
  - ▶ Equal (strong) calibration
  - ▶ Fair subgroup accuracy
  - ▶ ...
- Definitions are controversial and should be used depending on applications.

# Setup

- Supervised learning for binary classification
- $f$ : a classifier
- $Y \in \{0, 1\}$ : an outcome
- $X$ : features
- $A \in \{0, 1\}$ : a protected attribute
- $\hat{Y} := f(X, A) \in \{0, 1\}$ : a prediction

# Demographic Parity

Definition (demographic parity)

$$\mathbb{P} \left\{ \hat{Y} = 1 \mid A = 0 \right\} = \mathbb{P} \left\{ \hat{Y} = 1 \mid A = 1 \right\}$$

- Its variants appears in many papers.



# Demographic Parity

Definition (demographic parity)

$$\mathbb{P} \left\{ \hat{Y} = 1 \mid A = 0 \right\} = \mathbb{P} \left\{ \hat{Y} = 1 \mid A = 1 \right\}$$

- Its variants appears in many papers.
- Is this definition okay?

# Demographic Parity

## Definition (demographic parity)

$$\mathbb{P} \left\{ \hat{Y} = 1 \mid A = 0 \right\} = \mathbb{P} \left\{ \hat{Y} = 1 \mid A = 1 \right\}$$

- Its variants appears in many papers.
- Is this definition okay?
  - ✗ Actually not quite fair (in some common sense)
    - ★ A classifier accepts qualified applicants in  $A = 0$  but unqualified applicants in  $A = 1$ .
    - ★ e.g., when we don't have enough training samples for  $A = 1$ , this constraint forces to have  $\hat{Y} = 1$  for  $A = 1$ .

# Demographic Parity

## Definition (demographic parity)

$$\mathbb{P} \left\{ \hat{Y} = 1 \mid A = 0 \right\} = \mathbb{P} \left\{ \hat{Y} = 1 \mid A = 1 \right\}$$

- Its variants appears in many papers.
- Is this definition okay?
  - ✗ Actually not quite fair (in some common sense)
    - ★ A classifier accepts qualified applicants in  $A = 0$  but unqualified applicants in  $A = 1$ .
    - ★ e.g., when we don't have enough training samples for  $A = 1$ , this constraint forces to have  $\hat{Y} = 1$  for  $A = 1$ .
  - ✗ This definition does not allow the perfect predictor  $\hat{Y} = Y$ .

## Better Fairness Definitions

### Definition (equalized odd)

We say that a predictor  $\hat{Y}$  satisfies equalized odds with respect to the protected attribute  $A$  and outcome  $Y$  if  $\hat{Y}$  and  $A$  are independent conditional on  $Y$ , e.g.,

$$\mathbb{P} \left\{ \hat{Y} = 1 \mid A = 0, Y = y \right\} = \mathbb{P} \left\{ \hat{Y} = 1 \mid A = 1, Y = y \right\} \quad \forall y \in \{0, 1\}.$$

- The definition is applicable to other setups, e.g., multi-class classification.

## Better Fairness Definitions

### Definition (equalized odd)

We say that a predictor  $\hat{Y}$  satisfies equalized odds with respect to the protected attribute  $A$  and outcome  $Y$  if  $\hat{Y}$  and  $A$  are independent conditional on  $Y$ , e.g.,

$$\mathbb{P} \left\{ \hat{Y} = 1 \mid A = 0, Y = y \right\} = \mathbb{P} \left\{ \hat{Y} = 1 \mid A = 1, Y = y \right\} \quad \forall y \in \{0, 1\}.$$

- The definition is applicable to other setups, e.g., multi-class classification.
- If  $y = 1$ , this constrains equalizes true positive rates for both  $A = 0$  and  $A = 1$ .

## Better Fairness Definitions

### Definition (equalized odd)

We say that a predictor  $\hat{Y}$  satisfies equalized odds with respect to the protected attribute  $A$  and outcome  $Y$  if  $\hat{Y}$  and  $A$  are independent conditional on  $Y$ , e.g.,

$$\mathbb{P} \left\{ \hat{Y} = 1 \mid A = 0, Y = y \right\} = \mathbb{P} \left\{ \hat{Y} = 1 \mid A = 1, Y = y \right\} \quad \forall y \in \{0, 1\}.$$

- The definition is applicable to other setups, e.g., multi-class classification.
- If  $y = 1$ , this constraint equalizes true positive rates for both  $A = 0$  and  $A = 1$ .
- If  $y = 0$ , this constraint equalizes false positive rates for both  $A = 0$  and  $A = 1$ .

## Better Fairness Definitions

### Definition (equalized odd)

We say that a predictor  $\hat{Y}$  satisfies equalized odds with respect to the protected attribute  $A$  and outcome  $Y$  if  $\hat{Y}$  and  $A$  are independent conditional on  $Y$ , e.g.,

$$\mathbb{P} \left\{ \hat{Y} = 1 \mid A = 0, Y = y \right\} = \mathbb{P} \left\{ \hat{Y} = 1 \mid A = 1, Y = y \right\} \quad \forall y \in \{0, 1\}.$$

- The definition is applicable to other setups, e.g., multi-class classification.
- If  $y = 1$ , this constraint equalizes true positive rates for both  $A = 0$  and  $A = 1$ .
- If  $y = 0$ , this constraint equalizes false positive rates for both  $A = 0$  and  $A = 1$ .
- Is this enough?

# Better Fairness Definitions

## Definition (equalized odd)

We say that a predictor  $\hat{Y}$  satisfies equalized odds with respect to the protected attribute  $A$  and outcome  $Y$  if  $\hat{Y}$  and  $A$  are independent conditional on  $Y$ , e.g.,

$$\mathbb{P} \left\{ \hat{Y} = 1 \mid A = 0, Y = y \right\} = \mathbb{P} \left\{ \hat{Y} = 1 \mid A = 1, Y = y \right\} \quad \forall y \in \{0, 1\}.$$

- The definition is applicable to other setups, e.g., multi-class classification.
- If  $y = 1$ , this constraint equalizes true positive rates for both  $A = 0$  and  $A = 1$ .
- If  $y = 0$ , this constraint equalizes false positive rates for both  $A = 0$  and  $A = 1$ .
- Is this enough?
  - ✗ The accuracy is equally high for all demographics  $\rightarrow$  a model good at the majority will be penalized.



## Better Fairness Definitions

### Definition (Equal opportunity)

We say that a binary predictor  $\hat{Y}$  satisfies *equal opportunity* with respect to  $A$  and  $Y$  if

$$\mathbb{P} \left\{ \hat{Y} = 1 \mid A = 0, Y = 1 \right\} = \mathbb{P} \left\{ \hat{Y} = 1 \mid A = 1, Y = 1 \right\}.$$

- Suppose  $Y = 1$  is the “advantaged” outcome.
- Equal opportunity is weaker than equalized odd but typically allows stronger utility.

# A Score-based Predictor

A score-based predictor

$$\hat{Y} = \mathbb{1}(\hat{R} > t)$$

- We consider a real valued score  $\hat{R} \in [0, 1]$ , from which a classifier decides a label.
- e.g., a neural network  $R = f_{\text{NN}}(X)$
- Here, we suppose a pre-trained model is given and fixed; only change the threshold.

# A Score-based Predictor

## A score-based predictor

$$\hat{Y} = \mathbb{1}(\hat{R} > t)$$

- We consider a real valued score  $\hat{R} \in [0, 1]$ , from which a classifier decides a label.
- e.g., a neural network  $R = f_{\text{NN}}(X)$
- Here, we suppose a pre-trained model is given and fixed; only change the threshold.
- The equalized odds and equal opportunity definitions are characterized by true positive and false positive rates, which is controlled by the threshold, *i.e.*,

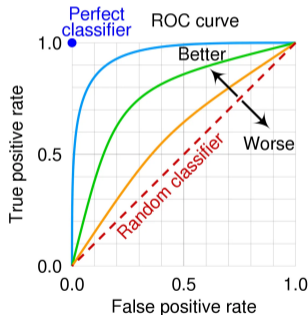
$$(\text{FP}) = \mathbb{P} \left\{ \hat{R} > t \mid A = a, Y = 0 \right\}$$

$$(\text{TP}) = \mathbb{P} \left\{ \hat{R} > t \mid A = a, Y = 1 \right\}.$$

# Receiver Operator Characteristic (ROC) Curves

## A-conditional ROC Curves

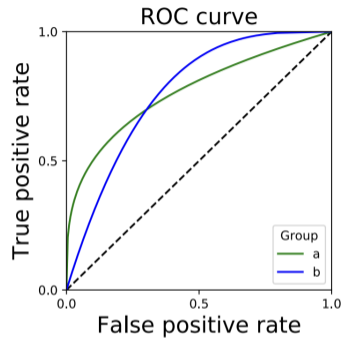
$$C_a(t) := \left( \underbrace{\mathbb{P} \left\{ \hat{R} > t \mid A = a, Y = 0 \right\}}_{\text{false positive (FP)}}, \underbrace{\mathbb{P} \left\{ \hat{R} > t \mid A = a, Y = 1 \right\}}_{\text{true positive (TP)}} \right)$$



Picture Credit: Ilyurek Kilic

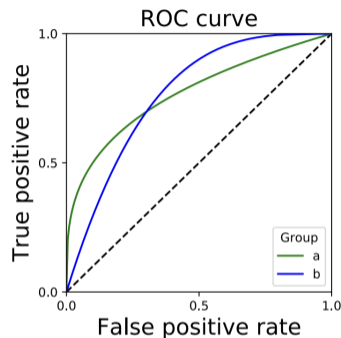
- $t \uparrow \rightarrow \text{FP} \downarrow$  and  $\text{TP} \downarrow$ .

## Algorithm for Equalized Odds



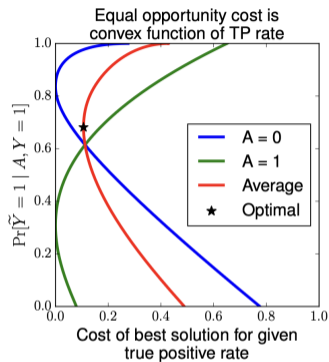
- Assume that two ROC curves are intersected, so let the intersecting points be  $(FP^*, TP^*)$
- Find  $(t_0, t_1)$  such that  $C_0(t_0) = (FP^*, TP^*)$  and  $C_1(t_1) = (FP^*, TP^*)$ .
- Our classifier is  $\hat{Y} := \mathbb{1}(\hat{R} > t_a)$

## Algorithm for Equalized Odds



- Assume that two ROC curves are intersected, so let the intersecting points be  $(FP^*, TP^*)$
- Find  $(t_0, t_1)$  such that  $C_0(t_0) = (FP^*, TP^*)$  and  $C_1(t_1) = (FP^*, TP^*)$ .
- Our classifier is  $\hat{Y} := \mathbb{1}(\hat{R} > t_a)$
- ✗ The accuracy is determined; when the accuracy is poor, no room to tune.

# Algorithm for Equal Opportunity

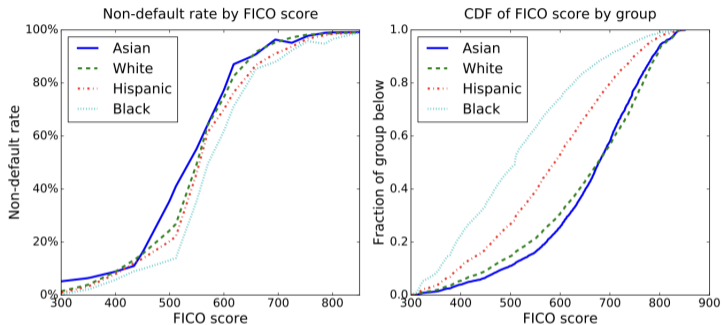


- Our classifier is  $\hat{Y} := \mathbb{1}(\hat{R} > t_a)$ .
- The algorithm solves the following constraint minimization.

$$\min_{t_0, t_1} \mathbb{E} \ell(\hat{Y}, Y) \quad \text{s.t.} \quad \text{TP}_0(\hat{Y}) = \text{TP}_1(\hat{Y})$$

►  $\ell$ : loss

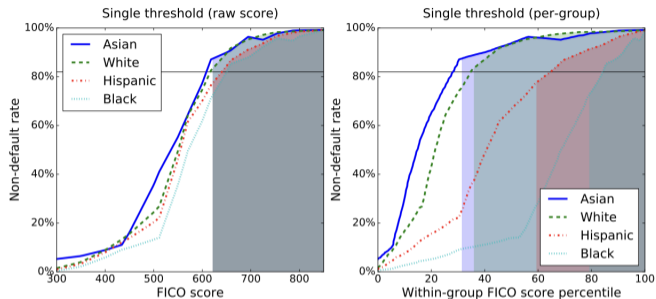
## Experiments: FICO Score (1/2)



- FICO score  $\hat{R}$ : a classifier to predict credit worthiness
- $Y = (\text{non-default})$ : failed to pay a debt
- $A$ : a race attribute (*i.e.*, Asian, white, Hispanic, black)

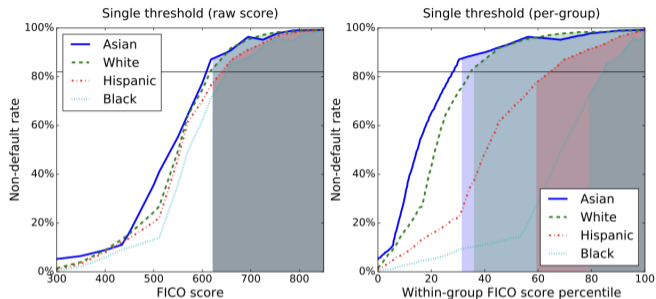


## Experiments: FICO Score (2/2)



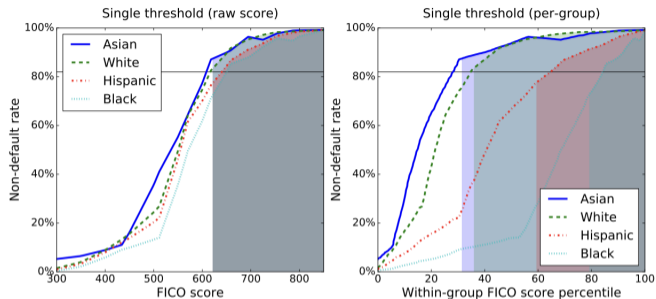
- $\hat{Y} := \mathbb{1}(\hat{R} > 620)$ : A standard classifier; is this fair classifier?
- (right  $x$  axis): rescaled, within-group score percentile
- (the fraction of the right shaded area) =  $\mathbb{P}\{\hat{Y} = 1 \mid Y = 1, A\}$

## Experiments: FICO Score (2/2)



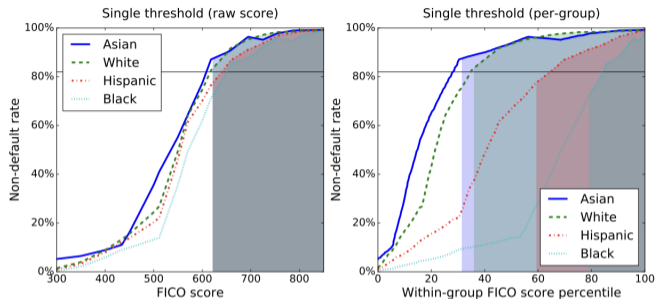
- $\hat{Y} := \mathbb{1}(\hat{R} > 620)$ : A standard classifier; is this fair classifier?
- (right  $x$  axis): rescaled, within-group score percentile
- (the fraction of the right shaded area) =  $\mathbb{P}\{\hat{Y} = 1 \mid Y = 1, A\}$
- Black non-defaulters are less likely to qualify for loans (than white or Asian ones)

## Experiments: FICO Score (2/2)



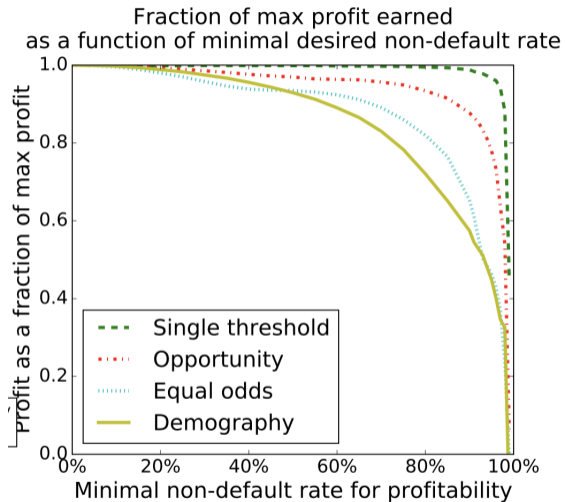
- $\hat{Y} := \mathbb{1}(\hat{R} > 620)$ : A standard classifier; is this fair classifier?
- (right  $x$  axis): rescaled, within-group score percentile
- (the fraction of the right shaded area) =  $\mathbb{P}\{\hat{Y} = 1 \mid Y = 1, A\}$
- Black non-defaulters are less likely to qualify for loans (than white or Asian ones)
- This classifier violates the fairness in equal opportunity.

## Experiments: FICO Score (2/2)



- $\hat{Y} := \mathbb{1}(\hat{R} > 620)$ : A standard classifier; is this fair classifier?
- (right  $x$  axis): rescaled, within-group score percentile
- (the fraction of the right shaded area) =  $\mathbb{P}\{\hat{Y} = 1 \mid Y = 1, A\}$
- Black non-defaulters are less likely to qualify for loans (than white or Asian ones)
- This classifier violates the fairness in equal opportunity.
- Satisfy the qualified odds?

## Experiments: Utility Performance



- Equal opportunity balances well between utility and fairness.

# Conclusion

- Fairness definitions
  - ① Demographic parity
  - ② Equalized Odds
  - ③ Equal Opportunity
- Fairness algorithms
  - ① Algorithm for Equalized Odds
  - ② Algorithm for Equal Opportunity
- There are neither “ $(\epsilon, \delta)$ -fairness” nor the proof of fairness; why?
  - ▶ Proving the fairness may be impossible without clearly understanding on domain-specific knowledge.
  - ▶ Fairness through Awareness!