DELAYED IMPACT OF FAIR MACHINE LEARNING Lydia T. Liu (UC Berkeley)





Joint work with Sarah Dean, Esther Rolf, Max Simchowitz, Moritz Hardt





Lydia T. Liu (UC Berkeley) | ICML 2018



Lydia T. Liu (UC Berkeley) | ICML 2018



Lydia T. Liu (UC Berkeley) | ICML 2018

"21 DEFINITIONS OF FAIRNESS" [Narayanan 2018] I. DEMOGRAPHIC PARITY 2. EQUALITY OF OPPORTUNITY **3. PREDICTIVE VALUE PARITY 4. GROUP CALIBRATION**



Machine learning systems are "fair"

Lydia T. Liu (UC Berkeley) | ICML 2018

"21 DEFINITIONS OF FAIRNESS" [Narayanan 2018] I. DEMOGRAPHIC PARITY 2. EQUALITY OF OPPORTUNITY **3. PREDICTIVE VALUE PARITY 4. GROUP CALIBRATION**



Protected groups are "better off"

Two groups with different score distributions (e.g. credit scores)



Lydia T. Liu (UC Berkeley) | ICML 2018



Would repay

Two groups with different score distributions (e.g. credit scores)

BLUE GROUP Count

Lydia T. Liu (UC Berkeley) | ICML 2018

ORANGE GROUP



Would repay

Two groups with different score distributions (e.g. credit scores)



Lydia T. Liu (UC Berkeley) | ICML 2018

ORANGE GROUP



Would repay

Approve loans according to **DEMOGRAPHIC PARITY**.



Lydia T. Liu (UC Berkeley) | ICML 2018



Would repay

Credit scores change with repayment (+) or default (-).



Lydia T. Liu (UC Berkeley) | ICML 2018



Would repay

Credit scores change with repayment (+) or default (-).



Lydia T. Liu (UC Berkeley) | ICML 2018







Credit scores change with repayment (+) or default (-).



Lydia T. Liu (UC Berkeley) | ICML 2018







WHAT HAPPENED?

Lydia T. Liu (UC Berkeley) | ICML 2018

WHAT HAPPENED?

Lydia T. Liu (UC Berkeley) | ICML 2018

Fairness criteria didn't seem to help the protected group, once we considered the *impact* of loans on scores.

Lydia T. Liu (UC Berkeley) | ICML 2018



OUR WORK

criteria

1. Introduce the "outcome curve", a tool for comparing the delayed impact of fairness

OUR WORK

- criteria

1. Introduce the "outcome curve", a tool for comparing the delayed impact of fairness

2. Provide a **complete characterization** of the delayed impact of 3 different fairness criteria

OUR WORK

- criteria
- 2. Provide a complete characterization of the delayed impact of 3 different fairness criteria

1. Introduce the "outcome curve", a tool for comparing the delayed impact of fairness

3. Show that fairness constraints may cause harm to groups they intended to protect

 \blacktriangleright A score R(X) is a scalar random variable that is a function of an individual's features X

► e.g. credit score is an integer from 300 to 850

- - ► e.g. credit score is an integer from 300 to 850
- > Any group of individuals has a particular **distribution** over scores:

 \blacktriangleright A score R(X) is a scalar random variable that is a function of an individual's features X

- - e.g. credit score is an integer from 300 to 850
- > Any group of individuals has a particular **distribution** over scores:



 \blacktriangleright A score R(X) is a scalar random variable that is a function of an individual's features X



- - ► e.g. credit score is an integer from 300 to 850
- > Any group of individuals has a particular **distribution** over scores:

$$\mathbb{P}\{R=r\}$$

a loan) once accepted

 \blacktriangleright A score R(X) is a scalar random variable that is a function of an individual's features X



Each score corresponds to an individual's success probability (e.g. probability of repaying



- - e.g. credit score is an integer from 300 to 850
- > Any group of individuals has a particular **distribution** over scores:

$$\mathbb{P}\{R=r\}$$

- a loan) once accepted
- Monotonicity assumption: Higher scores implies more likely to repay

 \blacktriangleright A score R(X) is a scalar random variable that is a function of an individual's features X



Each score corresponds to an individual's success probability (e.g. probability of repaying



Lydia T. Liu (UC Berkeley) | ICML 2018

Institution classifies individuals by characterize their expected utility:

► Institution classifies individuals by choosing an acceptance threshold score T to

- maximize their expected **utility**:

► Institution classifies individuals by choosing an acceptance threshold score T to

 $\mathbb{E}[\text{utility}|T] = \mathbb{E}[\text{reward from repayments}|T] - \mathbb{E}[\text{loss from defaults}|T]$

Institution classifies individuals by cho maximize their expected utility:

 $\mathbb{E}[\text{utility}|T] = \mathbb{E}[\text{reward from repayments}|T] - \mathbb{E}[\text{loss from defaults}|T]$



Lydia T. Liu (UC Berkeley) | ICML 2018

► Institution classifies individuals by choosing an acceptance threshold score T to

Institution classifies individuals by cho maximize their expected utility:

 $\mathbb{E}[\text{utility}|T] = \mathbb{E}[\text{reward from repayments}|T] - \mathbb{E}[\text{loss from defaults}|T]$



Lydia T. Liu (UC Berkeley) | ICML 2018

► Institution classifies individuals by choosing an acceptance threshold score T to

Threshold T corresponds to **acceptance rate** β for the group.

Institution classifies individuals by cho maximize their expected utility:

 $\mathbb{E}[\text{utility}|T] = \mathbb{E}[\text{reward from repayments}|T] - \mathbb{E}[\text{loss from defaults}|T]$



► Institution classifies individuals by choosing an acceptance threshold score T to

Threshold T corresponds to **acceptance rate** β for the group.

maximize their expected **utility**:

 $\mathbb{E}[\text{utility}|T] = \mathbb{E}[\text{reward from repayments}|T] - \mathbb{E}[\text{loss from defaults}|T]$



> When there are multiple groups, thresholds can be group-dependent.

Lydia T. Liu (UC Berkeley) | ICML 2018

> Institution classifies individuals by choosing an acceptance threshold score T to

Threshold *T* corresponds to **acceptance rate \beta** for the group.

- maximize their expected **utility**:



> When there are multiple groups, thresholds can be group-dependent.



Lydia T. Liu (UC Berkeley) | ICML 2018

► Institution classifies individuals by choosing an acceptance threshold score T to

 $\mathbb{E}[\text{utility}|T] = \mathbb{E}[\text{reward from repayments}|T] - \mathbb{E}[\text{loss from defaults}|T]$

Threshold *T* corresponds to **acceptance rate \beta** for the group.



Lydia T. Liu (UC Berkeley) | ICML 2018

Scores of accepted individuals change depending on their success.

Scores of accepted individuals change depending on their success.

$$R_{\text{new}} = \begin{cases} R_{\text{old}} + \\ R_{\text{old}} + \end{cases}$$

- c_+ if repaid
- *c*^{_} if defaulted

Scores of accepted individuals change depending on their success.

$$R_{\text{new}} = \begin{cases} R_{\text{old}} + \\ R_{\text{old}} + \end{cases}$$



Lydia T. Liu (UC Berkeley) | ICML 2018

- c_+ if repaid
- *c*_____if defaulted


MODEL DELAYED IMPACT ON GROUPS

Scores of accepted individuals change depending on their success.

$$R_{\text{new}} = \begin{cases} R_{\text{old}} + \\ R_{\text{old}} + \end{cases}$$



> The average change in score of each group is the **delayed impact**:

- if repaid \mathcal{C}_+
- if defaulted С_



MODEL DELAYED IMPACT ON GROUPS

Scores of accepted individuals change depending on their success.

$$R_{\text{new}} = \begin{cases} R_{\text{old}} + R_{\text{old}} + R_{\text{old}} \end{cases}$$



> The average change in score of each group is the **delayed impact**:

- if repaid \mathcal{C}_+
- if defaulted С_



- $\Delta \mu = \mathbb{E}[R_{\text{new}} R_{\text{old}}]$

Lydia T. Liu (UC Berkeley) | ICML 2018

Lemma: $\Delta \mu$ is a **concave** function of acceptance rate β under mild assumptions.



Lydia T. Liu (UC Berkeley) | ICML 2018

Lemma: $\Delta \mu$ is a **concave** function of acceptance rate β under mild assumptions.



Lydia T. Liu (UC Berkeley) | ICML 2018

Lemma: $\Delta \mu$ is a **concave** function of acceptance rate β under mild assumptions.



Lydia T. Liu (UC Berkeley) | ICML 2018

Lemma:



Lydia T. Liu (UC Berkeley) | ICML 2018

Lemma:



Lydia T. Liu (UC Berkeley) | ICML 2018

Lemma:



Lydia T. Liu (UC Berkeley) | ICML 2018

Lemma:



Lydia T. Liu (UC Berkeley) | ICML 2018

 $\Delta \mu$ is a **concave** function of acceptance rate β under mild assumptions.



← Acceptance Rate →

Lemma:



Lydia T. Liu (UC Berkeley) | ICML 2018

 $\Delta \mu$ is a **concave** function of acceptance rate β under mild assumptions.



← Acceptance Rate →

Lemma:



 $\Delta \mu$ is a **concave** function of acceptance rate β under mild assumptions.



← Acceptance Rate →

Lemma:



Lemma:



Lemma: $\Delta \mu$ is a **concave** function of acceptance rate β under mild assumptions.





← Acceptance Rate →

Lemma: $\Delta \mu$ is a **concave** function of acceptance rate β under mild assumptions.





 β_0

← Acceptance Rate →

Lemma:



Lydia T. Liu (UC Berkeley) | ICML 2018

Lydia T. Liu (UC Berkeley) | ICML 2018

Alternative to unconstrained utility maximization

- Alternative to unconstrained utility maximization
- > **Demographic Parity**: Equal Acceptance Rate

- Alternative to unconstrained utility maximization
- > **Demographic Parity**: Equal Acceptance Rate





- Alternative to unconstrained utility maximization
- > Demographic Parity: Equal Acceptance Rate





- > **Demographic Parity**: Equal Acceptance Rate



- > **Demographic Parity**: Equal Acceptance Rate



- > **Demographic Parity**: Equal Acceptance Rate



Equal Opportunity: Equal True Positive Rates [Hardt, Price and Srebro 2016]



- > **Demographic Parity**: Equal Acceptance Rate



Equal Opportunity: Equal True Positive Rates [Hardt, Price and Srebro 2016]





- > **Demographic Parity**: Equal Acceptance Rate



Equal Opportunity: Equal True Positive Rates [Hardt, Price and Srebro 2016]





- > **Demographic Parity**: Equal Acceptance Rate



Equal Opportunity: Equal True Positive Rates [Hardt, Price and Srebro 2016]





Lydia T. Liu (UC Berkeley) | ICML 2018

Theorem 1 [All outcome regimes are possible]

Equal opportunity and demographic parity may cause relative improvement, relative harm, or active harm.

Theorem 1 [All outcome regimes are possible]

relative harm, or active harm.

unconstrained utility maximization never causes active harm.

Equal opportunity and demographic parity may cause relative improvement,

Theorem 1 [All outcome regimes are possible]

relative harm, or active harm.

unconstrained utility maximization never causes active harm.



Lydia T. Liu (UC Berkeley) | ICML 2018

Equal opportunity and demographic parity may cause relative improvement,



Lydia T. Liu (UC Berkeley) | ICML 2018

Theorem 2

Demographic parity (DP) may cause active or relative harm by **over-acceptance**; equal opportunity (EO) doesn't.

Theorem 2

Demographic parity (DP) may cause active or relative harm by **over-acceptance**; equal opportunity (EO) doesn't.



Lydia T. Liu (UC Berkeley) | ICML 2018

Relative Improvement

Theorem 2

Demographic parity (DP) may cause active or relative harm by **over-acceptance**; equal opportunity (EO) doesn't.



Lydia T. Liu (UC Berkeley) | ICML 2018

Relative Improvement
Theorem 2

Demographic parity (DP) may cause active or relative harm by **over-acceptance**; equal opportunity (EO) doesn't.



Lydia T. Liu (UC Berkeley) | ICML 2018

Relative Improvement

Theorem 2

Demographic parity (DP) may cause active or relative harm by **over-acceptance**; equal opportunity (EO) doesn't.



Lydia T. Liu (UC Berkeley) | ICML 2018



Theorem 3

- Equal opportunity may cause relative harm by **under-acceptance**; demographic parity never under-accepts.
- **Relative Improvement**

Theorem 2

Demographic parity (DP) may cause active or relative harm by **over-acceptance**; equal opportunity (EO) doesn't.



Lydia T. Liu (UC Berkeley) | ICML 2018

Theorem 3

Equal opportunity may cause relative harm by **under-acceptance**; demographic parity never under-accepts.

Theorem 2

Demographic parity (DP) may cause active or relative harm by **over-acceptance**; equal opportunity (EO) doesn't.



Lydia T. Liu (UC Berkeley) | ICML 2018

Theorem 3

Equal opportunity may cause relative harm by **under-acceptance**; demographic parity never under-accepts.

Theorem 2

Demographic parity (DP) may cause active or relative harm by **over-acceptance**; equal opportunity (EO) doesn't.



Lydia T. Liu (UC Berkeley) | ICML 2018

Theorem 3

Equal opportunity may cause relative harm by **under-acceptance**; demographic parity never under-accepts.

Lydia T. Liu (UC Berkeley) | ICML 2018

- ► 300,000+ TransUnion TransRisk scores from 2003
- Scores range from 300 to 850 and are meant to predict default risk

- ► 300,000+ TransUnion TransRisk scores from 2003
- Scores range from 300 to 850 and are meant to predict default risk
- What we did

- ► 300,000+ TransUnion TransRisk scores from 2003
- Scores range from 300 to 850 and are meant to predict default risk

What we did

distributions, repayment probabilities, and relative sizes



- ► 300,000+ TransUnion TransRisk scores from 2003
- Scores range from 300 to 850 and are meant to predict default risk

What we did

- distributions, repayment probabilities, and relative sizes
- ► Model the bank's profit/loss ratio, e.g. +1:-4



- ► 300,000+ TransUnion TransRisk scores from 2003
- Scores range from 300 to 850 and are meant to predict default risk What we did
- distributions, repayment probabilities, and relative sizes
- ► Model the bank's profit/loss ratio, e.g. +1:-4
- Model the delayed impact of repayment/default on credit score, e.g. +75/-150



- ► 300,000+ TransUnion TransRisk scores from 2003
- Scores range from 300 to 850 and are meant to predict default risk What we did
- distributions, repayment probabilities, and relative sizes
- ► Model the bank's profit/loss ratio, e.g. +1:-4
- Model the delayed impact of repayment/default on credit score, e.g. +75/-150
- Compute "outcome curves" and delayed impact under different fairness criteria













1.0

















Lydia T. Liu (UC Berkeley) | ICML 2018





Lydia T. Liu (UC Berkeley) | ICML 2018

score change $\Delta\mu$

 \mathcal{N} bank utility

Lydia T. Liu (UC Berkeley) | ICML 2018



Why the large difference in delayed impact?

bank utility ${\cal U}$

Lydia T. Liu (UC Berkeley) | ICML 2018



Why the large difference in delayed impact? Maxima of outcome and utility curves under fairness criteria are **more misaligned** in the minority "black" group

bank utility \mathcal{U}^{-1}



Lydia T. Liu (UC Berkeley) | ICML 2018

Outcome curves provide a way to deviate from maximum utility while improving outcomes.

Lydia T. Liu (UC Berkeley) | ICML 2018

- Outcome curves provide a way to devolute
- ► Need for **domain-specific** models of delayed impact
 - Context-sensitive nature of fairness in machine learning

- Outcome curves provide a way to devolute outcomes.
- ► Need for **domain-specific** models of delayed impact
 - Context-sensitive nature of fairness in machine learning

FUTURE WORK

- Outcome curves provide a way to devolution
- ► Need for **domain-specific** models of delayed impact
 - Context-sensitive nature of fairness in machine learning

FUTURE WORK

► Moving beyond **binary** decisions

- outcomes.
- Need for domain-specific models of delayed impact
 - Context-sensitive nature of fairness in machine learning

FUTURE WORK

- Moving beyond binary decisions
- Moving beyond the mean score as measure of impact

- Outcome curves provide a way to devolute outcomes.
- ► Need for **domain-specific** models of delayed impact
 - Context-sensitive nature of fairness in machine learning

FUTURE WORK

- Moving beyond binary decisions
- ► Moving beyond the **mean** score as measure of impact
- Dynamics of the distributional impact of machine learning algorithms [Ensign et al. 2017; Hu and Chen 2017; Hashimoto et al. 2018]

Thank you!

DELAYED IMPACT OF FAIR MACHINE LEARNING Lydia T. Liu (UC Berkeley)



Poster today | 06:15 -- 09:00 PM | Hall B #213 Full version at arxiv.org/abs/1803.04383







Joint work with Sarah Dean, Esther Rolf, Max Simchowitz, Moritz Hardt