Discrimination in Machine Learning: From Principles to Measures and Mechanisms

Muhammad Bilal Zafar

joint work with

Isabel Valera, Manuel Gomez Rodriguez, Krishna P. Gummadi



Algorithmic decision making in practice



Algorithms being used to assist or replace human-decision making

- Algorithmic decision making used in several domains
 - Banking: Loan approval
 - Recruiting: Filtering and ranking applicants
 - Judiciary: Bail decisions
- Learn from historical training data

An algorithm can predict human behavior better than humans

An algorithm can predict human behavior better than humans

TheUpshot

ROBO RECRUITING

Can an Algorithm Hire Better Than a Human?

An algorithm can predict human behavior better than humans

TheUpshotROBO RECRUITING Can an Algorithm Hire Better Than a Human?

Even Imperfect Algorithms Can Improve the Criminal Justice System

Racism is Poisoning Online Ad Delivery, Says Harvard Professor

Racism is Poisoning Online Ad Delivery, Says Harvard Professor

Googling Politics

How the Google issue guide on candidates is biased.

Racism is Poisoning Online Ad Delivery, Says Harvard Professor

Googling Politics

How the Google issue guide on candidates is biased.

Machine Bias

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

What constitutes discrimination?

[...] wrongful imposition of relative disadvantage on persons based on their membership in a salient social group e.g., gender, race

[Altman'16]

What constitutes discrimination?

[...] wrongful imposition of relative disadvantage on persons based on their membership in a salient social group e.g., gender, race

[Altman'16]

Challenge # 1: Measures

- How do we measure wrongful relative disadvantage?
 - Fuzzy notion, not readily measurable
 - Existing measures may be insufficient

What constitutes discrimination?

[...] wrongful imposition of relative disadvantage on persons based on their membership in a salient social group e.g., gender, race

[Altman'16]

Challenge # 2: Mechanisms

- How to incorporate nondiscrimination into algorithmic decision making?
- Algorithmic decision system training
 - Optimize prediction accuracy
 - Enabling efficient training (crucial for large training datasets)

Part 1

Computational Measures for Nondiscrimination [WWW'17]

Recap: What constitutes discrimination?

[...] wrongful imposition of relative disadvantage on persons based on their membership in a salient social group e.g., gender, race

[Altman'16]

Toy example: University admission



Gender, race, etc

Relative disadvantage measure 1: Disparate treatment

- Achieve parity (or equality) in treatment
 - Decisions should not change with change in sensitive feature



Relative disadvantage measure 1: Disparate treatment

- Achieve parity (or equality) in treatment
 - Decisions should not change with change in sensitive feature



Relative disadvantage measure 1: Disparate treatment

- Achieve parity (or equality) in treatment
 - Decisions should not change with change in sensitive feature



Disparate treatment

Disparate treatment: A measure of direct discrimination



• Wrongful relative disadvantage: Basing decisions on gender, race, etc.

• Most intuitive notion of discrimination

• Focuses on direct or intentional discrimination

Relative disadvantage measure 2: Disparate impact

- Achieve parity (or equality) in impact
 - Positive outcome rates should be same for all groups

Relative disadvantage measure 2: Disparate impact

- Achieve parity (or equality) in impact
 - Positive outcome rates should be same for all groups





Relative disadvantage measure 2: Disparate impact

- Achieve parity (or equality) in impact
 - Positive outcome rates should be same for all groups



Disparate impact: A measure of indirect discrimination

$$\mathsf{P}(\hat{y} = 1 \mid \mathbf{Q}) = \mathsf{P}(\hat{y} = 1 \mid \mathbf{Q})$$

 $\hat{y} = 1$ is the desired outcome

- Wrongful relative disadvantage: Disparity in positive outcome rates
- Historical biases in training data (perpetuated biases)

Toy example: University admission



Disparate impact: A measure of indirect discrimination

$$\mathsf{P}(\hat{y} = 1 \mid \mathbf{Q}) = \mathsf{P}(\hat{y} = 1 \mid \mathbf{Q})$$

 $\hat{y} = 1$ is the desired outcome

- Wrongful relative disadvantage: Disparity in positive outcome rates
- Historical biases in training data (perpetuated biases)
- Correlations between sensitive and non-sensitive features
 - Correlation between race and zip code
 - Helps where disparate treatment might fail

Until now: 2 existing measures of discrimination

- Disparate treatment: Targets direct discrimination
 - Requires: $P(\hat{y} | \mathbf{x}, z) = P(\hat{y} | \mathbf{x})$
- Disparate impact: Targets indirect discrimination, when biased historical labels
 - Requires: $P(\hat{y} = 1 | \mathbf{Q}) = P(\hat{y} = 1 | \mathbf{Q})$

Until now: 2 existing measures of discrimination

- Disparate treatment: Targets direct discrimination
 - Requires: $P(\hat{y} | \mathbf{x}, z) = P(\hat{y} | \mathbf{x})$
- Disparate impact: Targets indirect discrimination, when biased historical labels
 - Requires: $P(\hat{y} = 1 | \mathbf{Q}) = P(\hat{y} = 1 | \mathbf{Q})$

Unbiased historical labels (Ground truth)

Ground truth → No wrongful relative disadvantage?

Example: Credit risk assessment



No wrongful relative disadvantage?

A fictitious dataset: Predict who will return the loan





Nmin $\sum L(\mathbf{x}_i, y_i, \mathbf{w})$ $\overline{i=1}$







Nmin $\sum L(\mathbf{x}_i, y_i, \mathbf{w})$ i=1

Women: Few errors

Feature 1



Disparate mistreatment: Different error rates

Relative disadvantage measure 3: Disparate mistreatment

- Achieve parity (or equality) in error rates
 - Error rates for all groups should be the same
- Wrongful relative disadvantage: Disparity in error rates
- Specially useful when ground truth labels are available

Formalizing disparate mistreatment

Equal error rates:

$$\mathsf{P}(\mathsf{y} \neq \hat{\mathsf{y}} \mid \mathbf{\sigma}) = \mathsf{P}(\mathsf{y} \neq \hat{\mathsf{y}} \mid \mathbf{Q})$$
Equality of overall error rates is not enough



Equality of overall error rates is not enough



Similar error rate for men and women

No disparate mistreatment?

Equality of overall error rates is not enough



Similar error rate for men and women

Errors for women are false positives (wrongly granted loan)

Errors for men are false negatives (wrongly denied loan)

Disproportionate advantage to women

Error rates should be considered separately

Formalizing disparate mistreatment

Equal error rates:

$$\mathsf{P}(\mathsf{y} \neq \hat{\mathsf{y}} \mid \mathbf{\sigma}) = \mathsf{P}(\mathsf{y} \neq \hat{\mathsf{y}} \mid \mathbf{Q})$$

Equal false negative rates: $P(y \neq \hat{y} \mid y = +1, \sigma) = P(y \neq \hat{y} \mid y = +1, Q)$

Equal false positive rates:

$$P(y \neq \hat{y} \mid y = -1, \sigma) = P(y \neq \hat{y} \mid y = -1, Q)$$

Formalizing disparate mistreatment

Equal error rates:

$$\mathsf{P}(\mathsf{y} \neq \hat{\mathsf{y}} \mid \mathbf{\sigma}) = \mathsf{P}(\mathsf{y} \neq \hat{\mathsf{y}} \mid \mathbf{Q})$$

Equal false negative rates: $P(y \neq \hat{y} \mid y = +1, \sigma) = P(y \neq \hat{y} \mid y = +1, Q)$

Equal false positive rates: $P(y \neq \hat{y} \mid y = -1, \sigma) = P(y \neq \hat{y} \mid y = -1, Q)$

Similar criteria for false omission and false discovery rates

Summary: 3 Measures of discrimination

- Disparate treatment: Targets direct discrimination
 - Requires: $P(\hat{y} | \mathbf{x}, z) = P(\hat{y} | \mathbf{x})$
- Disparate impact: Targets indirect discrimination, when biased historical labels
 - Requires: $P(\hat{y} = 1 | \mathbf{Q}) = P(\hat{y} = 1 | \mathbf{Q})$
- Disparate mistreatment: Targets indirect discrimination, when ground truth available
 - Requires: $P(y \neq \hat{y} | \sigma) = P(y \neq \hat{y} | \varphi)$
 - Also for other misclassification rates

Part 2

Mechanisms for Nondiscriminatory Machine Learning

[AISTATS'17; WWW'17]

Mechanisms for nondiscrimination

• Disparate treatment: $P(\hat{y} | \mathbf{x}, z) = P(\hat{y} | \mathbf{x})$

• Disparate impact: $P(\hat{y} = 1 | \mathbf{Q}) = P(\hat{y} = 1 | \mathbf{Q})$

• Disparate mistreatment: $P(y \neq \hat{y} | \sigma) = P(y \neq \hat{y} | Q)$

Toy example: University admission



Learn a decision boundary (w) in the feature space, separating the two classes

Learning the optimal boundary

- Learning → Minimize loss on the historical data
- Convex boundary-based loss functions

Logistic loss
$$\sum_{i=1}^{N} \log(1 + e^{-y_i d_{\mathbf{w}}(\mathbf{x}_i)})$$

SVM loss $||\mathbf{w}||^2 + C \sum_{i=1}^{N} \max(0, 1 - y_i d_{\mathbf{w}}(\mathbf{x}_i))$

Classification free of disparate treatment



 $\mathsf{P}(\hat{\mathsf{y}} \mid \mathsf{X}, \mathsf{z}) = \mathsf{P}(\hat{\mathsf{y}} \mid \mathsf{X})$

Do not use the sensitive feature

Mechanisms for nondiscrimination

- Disparate treatment: $P(\hat{y} | \mathbf{x}, z) = P(\hat{y} | \mathbf{x})$
 - Just drop z

Mechanisms for nondiscrimination

- Disparate treatment: $P(\hat{y} | \mathbf{x}, z) = P(\hat{y} | \mathbf{x})$
 - Just drop z
- Disparate impact: $P(\hat{y} = 1 | \mathbf{Q}) = P(\hat{y} = 1 | \mathbf{Q})$

Key Idea: Learn under constraints

min
$$\sum_{i=1}^{N} L(\mathbf{x}_i, y_i, \mathbf{w})$$

Key Idea: Learn under constraints

s.t.
$$-\varepsilon \le P(y_{pred} = 1 | z = 0) - P(y_{pred} = 1 | z = 1) \le \varepsilon$$

- Minimizing loss → Optimizing accuracy
- Adding constraints → Nondiscrimination goals

Key Idea: Learn under constraints

min
$$\sum_{i=1}^{N} L(\mathbf{x}_i, y_i, \mathbf{w})$$
 $z = 0$ (Men)
 $z = 1$ (Women)

s.t.
$$-\varepsilon \le P(y_{pred} = 1 | z = 0) - P(y_{pred} = 1 | z = 1) \le \varepsilon$$

Important insight

Tradeoff between accuracy & nondiscrimination

Key Idea: Learn under constraints

min
$$\sum_{i=1}^{N} L(\mathbf{x}_i, y_i, \mathbf{w})$$
 $z = 0$ (Men)
 $z = 1$ (Women)

s.t.
$$-\varepsilon \le P(y_{pred} = 1 | z = 0) - P(y_{pred} = 1 | z = 1) \le \varepsilon$$

- Non-convex for many well-known classifiers (logistic regression, SVM)
- Hard to compute efficiently

Goal: $-\varepsilon \le P(y_{pred} = 1 | z = 0) - P(y_{pred} = 1 | z = 1) \le \varepsilon$

Goal:
$$-\varepsilon \le P(y_{pred} = 1 | z = 0) - P(y_{pred} = 1 | z = 1) \le \varepsilon$$

$$\frac{1}{N} \sum_{i=1}^{N} (z_i - \bar{z}) d_{\mathbf{w}}(\mathbf{x}_i)$$

Key Idea: Limit the covariance between sensitive feature value and distance from decision boundary

Goal:
$$-\varepsilon \le P(y_{pred} = 1 | z = 0) - P(y_{pred} = 1 | z = 1) \le \varepsilon$$



Key Idea: Limit the covariance between sensitive feature value and distance from decision boundary

Goal:
$$-\varepsilon \le P(y_{pred} = 1 | z = 0) - P(y_{pred} = 1 | z = 1) \le \varepsilon$$



Key Idea: Limit the covariance between sensitive feature value and distance from decision boundary

Goal:
$$-\varepsilon \le P(y_{pred} = 1 | z = 0) - P(y_{pred} = 1 | z = 1) \le \varepsilon$$



In other words: Limit the difference in average strength of acceptance or rejection between sensitive feature groups

Key Idea: Learn under constraints



(Optimal) Classifier without any constraints: $\mathbf{w}^* = \operatorname{argmin} \sum_{i=1}^{N} L(\mathbf{x}_i, y_i, \mathbf{w})$

(Optimal) Classifier without any constraints: $\mathbf{w}^* = \operatorname{argmin} \sum_{i=1}^{N} L(\mathbf{x}_i, y_i, \mathbf{w})$

$$\left|\frac{1}{N}\sum_{i=1}^{N}(z_i-\bar{z})d_{\mathbf{w}}(\mathbf{x}_i)\right.$$

min

(Optimal) Classifier without any constraints: $\mathbf{w}^* = \operatorname{argmin} \sum_{i=1}^{N} L(\mathbf{x}_i, y_i, \mathbf{w})$

min

$$\left|\frac{1}{N}\sum_{i=1}^{N}(z_i-\bar{z})d_{\mathbf{w}}(\mathbf{x}_i)\right|$$

s.t.
$$\sum_{i=1}^{N} L(\mathbf{x}_i, y_i, \mathbf{w}) \le (1+\gamma) \sum_{i=1}^{N} L(\mathbf{x}_i, y_i, \mathbf{w}^*)$$

 $\gamma = 0 \rightarrow$ No additional loss $\gamma > 0 \rightarrow$ Allow additional loss

Allows for fine-grained control over loss in accuracy

(Optimal) Classifier without any constraints: $\mathbf{w}^* = \operatorname{argmin} \sum_{i=1}^{N} L(\mathbf{x}_i, y_i, \mathbf{w})$

min
$$\left| \frac{1}{N} \sum_{i=1}^{N} (z_i - \bar{z}) d_{\mathbf{w}}(\mathbf{x}_i) \right|$$

s.t.
$$L(\mathbf{x}_j, y_j, \mathbf{w}) \le (1 + \gamma_j) L(\mathbf{x}_j, y_j, \mathbf{w}^*) \quad \forall j \in \{1, 2, ..., N\}$$

Set γ_j for individual points $\gamma_j = 0 \rightarrow$ No additional loss $\gamma_j > 0 \rightarrow$ Allow additional loss

Preserving accuracy for <u>specific</u> points

Extending constraints to non-linear models

$$\max \sum_{i=1}^{N} \alpha_{i} - \frac{1}{2} \sum_{i,j=1}^{N} \alpha_{i} \alpha_{j} y_{i} y_{j} k(\mathbf{x}_{i}, \mathbf{x}_{j})$$

s.t. $0 \le \alpha_{i} \le C \quad \forall i$
$$\sum_{i=1}^{N} \alpha_{i} y_{i} = 0$$

$$d_{\alpha}(\mathbf{x}) = \sum_{i=1}^{N} \alpha_i y_i k(\mathbf{x}, \mathbf{x}_i)$$

Extending constraints to non-linear models

N

i=1

Linear on α

constraints too

Mechanisms for nondiscrimination

- Disparate treatment: $P(\hat{y} | \mathbf{x}, z) = P(\hat{y} | \mathbf{x})$
 - Just drop z
- Disparate impact: $P(\hat{y} = 1 | \mathbf{Q}) = P(\hat{y} = 1 | \mathbf{Q})$
 - Convex constraints

Mechanisms for nondiscrimination

- Disparate treatment: $P(\hat{y} | \mathbf{x}, z) = P(\hat{y} | \mathbf{x})$
 - Just drop z
- Disparate impact: $P(\hat{y} = 1 | \mathbf{Q}) = P(\hat{y} = 1 | \mathbf{Q})$
 - Convex constraints
- Disparate mistreatment: $P(y \neq \hat{y} | \sigma) = P(y \neq \hat{y} | Q)$

Classification free of disparate mistreatment

Classification free of disparate mistreatment

Key Idea: Learn under constraints

$$\min \sum_{i=1}^{N} L(\mathbf{x}_{i}, y_{i}, \mathbf{w})$$

s.t. $-\varepsilon \leq \mathsf{P}(\mathsf{y} \neq \hat{\mathsf{y}} \mid \mathbf{\sigma}) - \mathsf{P}(\mathsf{y} \neq \hat{\mathsf{y}} \mid \mathbf{Q}) \leq \varepsilon$

$$z = 0$$
 (Men)
 $z = 1$ (Women)

- Non-convex for many well-known classifiers (logistic regression, SVM)
- Hard to compute efficiently

Disparate mistreatment constraints



Disparate mistreatment constraints



 $g_{\mathbf{w}}(\mathbf{x}, y) = \min(0, yd_{\mathbf{w}}(\mathbf{x}))$
Disparate mistreatment constraints



 $g_{\mathbf{w}}(\mathbf{x}, y) = \min(0, yd_{\mathbf{w}}(\mathbf{x}))$

Disparate mistreatment constraints



Key Idea: Limit the covariance between sensitive feature value and <u>misclassification</u> distance from decision boundary

Disparate mistreatment constraints

Key Idea: Learn under constraints

$$\min \sum_{i=1}^{N} L(\mathbf{x}_{i}, y_{i}, \mathbf{w})$$

$$g_{\mathbf{w}}(\mathbf{x}, y) = \min(0, yd_{\mathbf{w}}(\mathbf{x}))$$

$$g_{\mathbf{w}}(\mathbf{x}, y) = \min(0, yd_{\mathbf{w}}(\mathbf{x}))$$

$$g_{\mathbf{w}}(\mathbf{x}, y) = \min(0, yd_{\mathbf{w}}(\mathbf{x}))$$

Disciplined Convex-Concave Program (DCCP) (can be solved efficiently) [Shen, Diamond, Gu, Boyd, 2016]

Mechanisms for nondiscrimination

- Disparate treatment: $P(\hat{y} | \mathbf{x}, z) = P(\hat{y} | \mathbf{x})$
 - Just drop z
- Disparate impact: $P(\hat{y} = 1 | \mathbf{Q}) = P(\hat{y} = 1 | \mathbf{Q})$
 - Convex constraints
- Disparate mistreatment: $P(y \neq \hat{y} | \sigma) = P(y \neq \hat{y} | Q)$
 - Convex-concave constraints

Mechanisms for nondiscrimination

- Disparate treatment: $P(\hat{y} | \mathbf{x}, z) = P(\hat{y} | \mathbf{x})$
 - Just drop z
- Disparate impact: $P(\hat{y} = 1 | \mathbf{Q}) = P(\hat{y} = 1 | \mathbf{Q})$
 - Convex constraints
- Disparate mistreatment: $P(y \neq \hat{y} | \sigma) = P(y \neq \hat{y} | \varphi)$
 - Convex-concave constraints

Case study: Crime risk prediction

Crime risk prediction: Background

- Person arrested on suspicion of a crime
 - Appear in front of a judge before the trial
 - Judge decides whether to give bail or not
 - Predict whether the person will recidivate
 - Recidivism: Commit a crime within two years

- COMPAS tool
 - Machine learning model to advise the judge
 - Predicts a recidivism probability

ProPublica COMPAS dataset

- ProPublica gathered COMPAS assessments
 - Broward Country, FL for 2013-14
 - Features: arrest charge, #prior offenses, age,...
 - Class label: 2-year recidivism
- Train a classifier on this dataset

Case study: Crime risk prediction

Can traditional classifiers lead to disparate mistreatment?

• Can our approach help avoid disparate mistreatment?

Disparate mistreatment in risk prediction

Trained a logistic regression classifier to predict recidivism

Positive class: Will recidivate

Race	FPR	FNR
Black		
White		

Disparate mistreatment in risk prediction

Trained a logistic regression classifier to predict recidivism

Positive class: Will recidivate

Race	FPR	FNR
Black	34%	
White	15%	

• False positive: Non-recid. person wrongly classified as recid.

Disparate mistreatment in risk prediction

Trained a logistic regression classifier to predict recidivism

Positive class: Will recidivate

Race	FPR	FNR
Black	34%	32%
White	15%	55%

- False positive: Non-recid. person wrongly classified as recid.
- False negative: Recid. person wrongly classified as non-recid.

Case study: Crime risk prediction

- Can traditional classifiers lead to disparate mistreatment?
 - Disparity in both FPR and FNR!

• Can our approach help avoid disparate mistreatment?

Removing disparate mistreatment



Removing disparate mistreatment

Introduce FPR and FNR constraints



Case study: Crime risk prediction

- Can traditional classifiers lead to disparate mistreatment?
 - Disparity in both FPR and FNR!

- Can our approach help avoid disparate mistreatment?
 - Yes!
 - Experiments with several synthetic datasets

Summary: Discrimination in machine leanring

- Part 1: Defined three measures of discrimination
 - Disparate treatment / impact / mistreatment
- Part 2: Designed mechanisms for each of them
 - Proposed tractable and efficient proxies

Summary: Discrimination in machine leanring

- Part 1: Defined three measures of discrimination
 - Disparate treatment / impact / mistreatment
- Part 2: Designed mechanisms for each of them
 - Proposed tractable and efficient proxies
- Part 3: Refined measures and designed mechanisms
 - Inspired by envy-freeness and Nash bargaining solution

Summary: Discrimination in machine leanring

- Part 1: Defined three measures of discrimination
 - Disparate treatment / impact / mistreatment
- Part 2: Designed mechanisms for each of them
 - Proposed tractable and efficient proxies
- Part 3: Refined measures and designed mechanisms
 - Inspired by envy-freeness and Nash bargaining solution

Code publicly available at

https://github.com/mbilalzafar/fair-classification