From Parity to Preference-based Notions of Fairness in Classification

Muhammad Bilal Zafar

joint work with

Isabel Valera, Manuel Gomez Rodriguez, Krishna P. Gummadi, Adrian Weller





Existing notions of fairness in classification

- Inspired by anti-discrimination legislation
 - Rooted in social sciences ideas of parity (or equality)

Existing notions of fairness in classification

- Inspired by anti-discrimination legislation
 - Rooted in social sciences ideas of parity (or equality)
- Parity in treatment (avoid disparate treatment)
 - Outcome should not depend on the sensitive group membership (e.g., men, women)

Existing notions of fairness in classification

- Inspired by anti-discrimination legislation
 - Rooted in social sciences ideas of parity (or equality)
- Parity in treatment (avoid disparate treatment)
 - Outcome should not depend on the sensitive group membership (e.g., men, women)
- Parity in impact (avoid *disparate impact*)
 - Similar fraction of beneficial outcomes received by different groups (e.g., men, women)
 - Beneficial outcome: Granted Ioan, etc.

Parity can be a very stringent criterion

- Impossibility results
 - Certain parities cannot be achieved together
 - [Chouldechova, Big Data'17; Kleinberg et al., ITCS'17]

Parity can be a very stringent criterion

- Impossibility results
 - Certain parities cannot be achieved together
 - [Chouldechova, Big Data'17; Kleinberg et al., ITCS'17]
- High cost of parity fairness
 - Inherent tradeoff between accuracy and fairness
 - [Corbett-Davies et al., KDD'17]

Parity can be a very stringent criterion

- Impossibility results
 - Certain parities cannot be achieved together
 - [Chouldechova, Big Data'17; Kleinberg et al., ITCS'17]
- High cost of parity fairness
 - Inherent tradeoff between accuracy and fairness
 - [Corbett-Davies et al., KDD'17]
- All groups may end up losing
 - Benefits for all groups decrease
 - [Zafar et al., AISTATS'17]

From parity to preference

• Explore preference-based notions of fairness

- Inspired by fair division in economics
 - Parity treatment to preferred treatment
 - Parity impact to preferred_impact
- Outcomes may not follow parity
 - Groups involved prefer their respective outcomes
 - Gains in accuracy over parity fairness

This talk

Defining preference-based fairness notions

• Formalizing preference-based notions

- Training preferentially-fair classifiers
 - Mechanism design
 - Case study: NYPD SQF dataset

This talk

• **Defining preference-based fairness notions**

Formalizing preference-based notions

- Training preferentially-fair classifiers
 - Mechanism design
 - Case study: NYPD SQF dataset

Some preliminaries

- Groups based on sensitive feature
 - Gender (men, women)
 - Race (African-American, Hispanic, White, ...)
- Classification outcomes
 - Beneficial (+ve class) and not beneficial (-ve class)
- Group preference
 - Each group collectively prefers classifier with most beneficial outcomes

Some preliminaries

- Groups based on sensitive feature
 - Gender (men, women)
 - Race (African-American, Hispanic, White, ...)
- Classification outcomes
 - Beneficial (+ve class) and not beneficial (-ve class)
- Group preference
 - Each group collectively prefers classifier with most beneficial outcomes

1

- Classifier 1 (10% beneficial outcomes)
- Classifier 2 (30% beneficial outcomes)
- Classifier 3 (15% beneficial outcomes)

Parity treatment to preferred treatment

Parity treatment (Existing notion)

- Anti-discrimination notion of parity
- Parity: Individuals in G₁ get same decisions by posing as G₂
- Changing sensitive feature
 should not change benefits
- Cannot use groupconditional classifiers

An example of parity treatment



An example of parity treatment

Parity treatment Acc: 0.83 (500/600)



Parity treatment to preferred treatment

Parity treatment (Existing notion)

- Anti-discrimination notion of parity
- Parity: Individuals in G₁ get same decisions by posing as G₂
- Changing sensitive feature should not change benefits
- Cannot use groupconditional classifiers

Preferred treatment (New notion)

- Fair division notion of group envy-freeness
- Envy-freeness: G₁ collectively will not prefer to be posing as G₂
- Changing sensitive feature
 should not increase benefits
- Can use group-conditional classifiers

Parity treatment vs. preferred treatment

Parity treatment Acc: 0.83 (500/600) **Preferred treatment Acc: 1.00** (600/600)



Parity treatment vs. preferred treatment

Men using men's classifier 33% (100/300) get benefits

Men using women's classifier 0% get benefits **Preferred treatment Acc: 1.00** (600/600)



Parity treatment vs. preferred treatment

Men using men's classifier 33% (100/300) get benefits

Men using women's classifier 0% get benefits

Women using women's classifier 66% (200/300) get benefits

Women using men's classifier 0% get benefits **Preferred treatment Acc: 1.00** (600/600)



No incentive to change group Avoids reverse discrimination (or lowering the bar) claims

Preferred treatment: A relaxation of parity treatment

- Parity treatment: Changing sensitive feature does not change group benefits
- Preferred treatment: Changing sensitive feature does not increase group benefits



Each solution satisfying parity treatment also satisfies preferred treatment (Room for more accurate solutions)

Parity impact to preferred impact

Parity impact (Existing notion)

- Anti-discrimination notion of parity
- Parity: 20% from G1 accepted, 20% from G2 accepted
- All groups gets equal fraction of beneficial outcomes

An example of parity impact

Parity impact

Acc: 0.72 Benefit: 22% (M), 22%(W)



Parity impact to preferred impact

Parity impact (Existing notion)

- Anti-discrimination notion of parity
- Parity: 20% from G1 accepted, 20% from G2 accepted
- All groups gets equal fraction of beneficial outcomes

Preferred impact (New notion)

- Fair division notion of bargaining solution
- Bargaining solution: 25% from G1 accepted, 40% from G2 accepted (or, revert to parity)
- All group gets at least as much benefits as parity impact

Parity impact vs. preferred impact

Parity impact

Acc: 0.72 Benefit: 22% (M), 22%(W)

Preferred impact

Acc: 1.00 Benefit: 33% (M), 67%(W)



Both groups have incentive to move to preferred impact

Preferred impact: A relaxation of parity impact

- Parity impact: All groups should get similar fraction of benefits
- Preferred impact: All groups should get at least as much benefits as parity impact



Each solution satisfying parity impact also satisfies preferred impact (Room for more accurate solutions)

This talk

Defining preference-based fairness notions

• Formalizing preference-based notions

- Training preferentially-fair classifiers
 - Mechanism design
 - Case study: NYPD SQF dataset

Given

- Sensitive feature group z
- Classifier θ with outcomes $\hat{y} \in [-1, 1]$
- $B_z(\theta) = P(\hat{y} = 1 | z, \theta)$

- Given
 - Sensitive feature group z
 - Classifier θ with outcomes $\hat{y} \in [-1, 1]$
 - $B_z(\theta) = P(\hat{y} = 1 | z, \theta)$

Benefits for men by **θ**

$$\mathsf{B}_{\sigma}(\boldsymbol{\theta}) = \mathsf{P}(\hat{\mathsf{y}} = 1 \mid \boldsymbol{\sigma}, \boldsymbol{\theta})$$

Benefits for women by $\boldsymbol{\theta}$

$$\mathsf{B}_{\mathsf{Q}}(\boldsymbol{\theta}) = \mathsf{P}(\hat{\mathsf{y}} = 1 \mid \mathbf{Q}, \, \boldsymbol{\theta})$$

Preferred treatment

$B_{\sigma}(\mathbf{\theta}_{\sigma})$	2	B♂(θ ♀)
Bç(θ ç)	\geq	Bç(θ ♂)

- Train separate classifiers for men (θ_{σ}) and women (θ_{φ})
- Benefits for men with their classifier more than benefits for men with women's classifier
- Benefits for women with their classifier more than benefits for women with men's classifier

Preferred treatment

$$\begin{array}{ll} \mathsf{B}_{\sigma}(\boldsymbol{\theta}_{\sigma}) & \geq & \mathsf{B}_{\sigma}(\boldsymbol{\theta}_{\varphi}) \\ \mathsf{B}_{\varphi}(\boldsymbol{\theta}_{\varphi}) & \geq & \mathsf{B}_{\varphi}(\boldsymbol{\theta}_{\sigma}) \end{array}$$

Preferred impact

Given parity impact classifiers θ'_{σ} and θ'_{φ}

$$\begin{array}{lll} \mathsf{B}_{\sigma}(\boldsymbol{\theta}_{\sigma}) & \geq & \mathsf{B}_{\sigma}(\boldsymbol{\theta}'_{\sigma}) & \overset{\mathsf{Benefit}}{\underset{\mathsf{parity imparity imparity imparity imparity (constrained}{})} \\ \mathsf{B}_{\varsigma}(\boldsymbol{\theta}_{\varsigma}) & \geq & \mathsf{B}_{\varsigma}(\boldsymbol{\theta}'_{\varsigma}) & \overset{\mathsf{Benefit}}{\underset{\mathsf{parity imparity imparity imparity (constrained}{})} \end{array}$$

Benefits from Darity impact classifier (constants)

This talk

Defining preference-based fairness notions

Formalizing preference-based notions

- Training preferentially-fair classifiers
 - Mechanism design
 - Case study: NYPD SQF dataset

min
$$\sum_{\mathcal{O}} L(\mathbf{x}, y, \boldsymbol{\theta}_{\mathcal{O}}) + \sum_{\mathcal{Q}} L(\mathbf{x}, y, \boldsymbol{\theta}_{\mathcal{Q}})$$

$$\min \sum_{\substack{\mathcal{O}}} L(\mathbf{x}, y, \boldsymbol{\theta}_{\mathcal{O}}) + \sum_{\substack{\mathcal{Q}}} L(\mathbf{x}, y, \boldsymbol{\theta}_{\mathcal{Q}}) \\ \text{s.t} \quad B_{\mathcal{O}}(\boldsymbol{\theta}_{\mathcal{O}}) \geq B_{\mathcal{O}}(\boldsymbol{\theta}_{\mathcal{Q}}) \quad \begin{array}{c} \mathsf{Add} \\ \mathsf{preferred} \\ B_{\mathcal{Q}}(\boldsymbol{\theta}_{\mathcal{Q}}) \geq B_{\mathcal{Q}}(\boldsymbol{\theta}_{\mathcal{O}}) \\ \end{array}$$

$$\begin{split} \min & \sum_{\mathcal{O}} L(\mathbf{x}, y, \boldsymbol{\theta}_{\mathcal{O}}) + \sum_{\mathcal{Q}} L(\mathbf{x}, y, \boldsymbol{\theta}_{\mathcal{Q}}) \\ \text{s.t.} & P(\hat{y} = 1 | \mathcal{O}, \boldsymbol{\theta}_{\mathcal{O}}) \geq P(\hat{y} = 1 | \mathcal{O}, \boldsymbol{\theta}_{\mathcal{Q}}) \\ & P(\hat{y} = 1 | \mathcal{Q}, \boldsymbol{\theta}_{\mathcal{Q}}) \geq P(\hat{y} = 1 | \mathcal{Q}, \boldsymbol{\theta}_{\mathcal{O}}) \end{split}$$

$$\begin{split} \min & \sum_{\mathcal{O}} L(\mathbf{x}, y, \boldsymbol{\theta}_{\mathcal{O}}) + \sum_{\mathcal{Q}} L(\mathbf{x}, y, \boldsymbol{\theta}_{\mathcal{Q}}) \\ \text{s.t.} & P(\hat{y} = 1 | \mathcal{O}, \boldsymbol{\theta}_{\mathcal{O}}) \geq P(\hat{y} = 1 | \mathcal{O}, \boldsymbol{\theta}_{\mathcal{Q}}) \\ & P(\hat{y} = 1 | \mathcal{Q}, \boldsymbol{\theta}_{\mathcal{Q}}) \geq P(\hat{y} = 1 | \mathcal{Q}, \boldsymbol{\theta}_{\mathcal{O}}) \end{split}$$

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

Idea: Learn under constraints

$$\begin{split} \min & \sum_{\mathcal{O}} L(\mathbf{x}, y, \boldsymbol{\theta}_{\mathcal{O}}) + \sum_{\mathcal{Q}} L(\mathbf{x}, y, \boldsymbol{\theta}_{\mathcal{Q}}) \\ \text{s.t.} & P(\hat{y} = 1 | \mathcal{O}, \boldsymbol{\theta}_{\mathcal{O}}) \geq P(\hat{y} = 1 | \mathcal{O}, \boldsymbol{\theta}_{\mathcal{Q}}) \\ & P(\hat{y} = 1 | \mathcal{Q}, \boldsymbol{\theta}_{\mathcal{Q}}) \geq P(\hat{y} = 1 | \mathcal{Q}, \boldsymbol{\theta}_{\mathcal{Q}}) \end{split}$$

Approximate positive class probability using ramp function $\max(0, \boldsymbol{\theta}^T \mathbf{x})$

Non-zero when $\hat{y}=1$, zero otherwise

Idea: Learn under constraints

$$\min \sum_{\mathcal{O}^{T}} L(\mathbf{x}, y, \boldsymbol{\theta}_{\mathcal{O}^{T}}) + \sum_{\mathcal{Q}} L(\mathbf{x}, y, \boldsymbol{\theta}_{\mathcal{Q}})$$
s.t
$$\sum_{\mathcal{O}^{T}} \max(0, \boldsymbol{\theta}_{\mathcal{O}^{T}}^{T} \mathbf{x}) \ge \sum_{\mathcal{O}^{T}} \max(0, \boldsymbol{\theta}_{\mathcal{Q}}^{T} \mathbf{x})$$

$$\sum_{\mathcal{O}^{T}} \max(0, \boldsymbol{\theta}_{\mathcal{Q}}^{T} \mathbf{x}) \ge \sum_{\mathcal{O}^{T}} \max(0, \boldsymbol{\theta}_{\mathcal{O}^{T}}^{T} \mathbf{x})$$

$$\sum_{\mathcal{Q}} \max(0, \boldsymbol{\theta}_{\mathcal{Q}}^{T} \mathbf{x}) \ge \sum_{\mathcal{Q}} \max(0, \boldsymbol{\theta}_{\mathcal{O}^{T}}^{T} \mathbf{x})$$

$$P(\hat{y} = 1 | \mathcal{Q}, \boldsymbol{\theta}_{\mathcal{Q}}) P(\hat{y} = 1 | \mathcal{Q}, \boldsymbol{\theta}_{\mathcal{O}^{T}})$$

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

Training preferred impact classifier

Training preferred impact classifier

θ'σ and θ'φ
parity impact
classifiers

Training preferred impact classifier

Idea: Learn under constraints

$$\min \sum_{\sigma'} L(\mathbf{x}, y, \boldsymbol{\theta}_{\sigma'}) + \sum_{\varphi} L(\mathbf{x}, y, \boldsymbol{\theta}_{\varphi})$$
s.t
$$\sum_{\sigma'} \max(0, \boldsymbol{\theta}_{\sigma'}^T \mathbf{x}) \ge \sum_{\sigma'} \max(0, \boldsymbol{\theta}_{\sigma'}^T \mathbf{x}) \quad \boldsymbol{\theta'}_{\sigma'} \text{ and } \boldsymbol{\theta'}_{\varphi}$$

$$\sum_{\sigma'} \max(0, \boldsymbol{\theta}_{\varphi}^T \mathbf{x}) \ge \sum_{\sigma'} \max(0, \boldsymbol{\theta}_{\varphi}^{'T} \mathbf{x}) \quad \boldsymbol{\theta'}_{\sigma'} \text{ and } \boldsymbol{\theta'}_{\varphi}$$

$$parity \text{ impact classifiers}$$

$$\sum_{\varphi} Convex \quad (Constant) \quad \mathbf{f}$$

$$P(\hat{y} = 1 | \varphi, \boldsymbol{\theta}_{\varphi}) \quad P(\hat{y} = 1 | \varphi, \boldsymbol{\theta'}_{\varphi})$$

Disciplined Convex-Concave Program (DCCP)

This talk

Defining preference-based fairness notions

• Formalizing preference-based notions

- Training preferentially-fair classifiers
 - Mechanism design
 - Case study: NYPD SQF dataset

Are the preferential fairness constraints effective?

 Does preferential fairness lead to accuracy gains over parity fairness?

- Two sensitive feature groups
 - African-Americans and Whites

- Two classes
 - Pedestrian in possession of a weapon (non-beneficial)
 - No weapon (beneficial)

Toy classification task

	Uncons.		
Af. Americans	67% (34%)		
Whites	12% (22%)		
Accuracy	0.74		

Dissimilar group benefits Incentive to change group

	Uncons.	Parity fairness
Af. Americans	67% (34%)	50% (50%)
Whites	12% (22%)	52% (52%)
Accuracy	0.74	0.61

Similar group benefits No incentive to change group Large drop in accuracy

Uncons.		Parity fairness	Preferred fairness
Af. Americans	67% (34%)	50% (50%)	78% (76%)
Whites	12% (22%)	52% (52%)	52% (33%)
Accuracy	0.74	0.61	0.68

Group benefits more than parity No incentive to change group Modest drop in accuracy

- Are the preferential fairness constraints effective?
 - Yes. Each group prefers their outcome.

- Does preferential fairness lead to accuracy gains over parity fairness?
 - Yes. Smaller loss in accuracy.

Conclusion

- Preference-based notions of fairness
 - Each group prefers to pose as itself
 - Each group prefers the outcomes over parity impact
- Preference-based notions can lead to more accurate solutions

Paper at https://tinyurl.com/preference-based-fairness