

Practical Individual Fairness

Mikhail Yurochkin



Machine Bias

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

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica

May 23, 2016



2019 study of 13.2 million mortgage and refinancing applications

BLACK & LATINX

61% rejection rates

+5.3 basis points
(bps) on mortgage
interest rates for
fintech lending

+7.9 bps on overall
mortgage interest
rates

EVERYONE ELSE

48%

→ \$756M
annual race premium

Bartlett, et al. "Consumer-Lending Discrimination in the FinTech Era." NBER 2019.

ECONOMIC VIEW

Biased Algorithms Are Easier to Fix Than Biased People

Racial discrimination by algorithms or by people is harmful — but that's where the similarities end.



Tim Cook

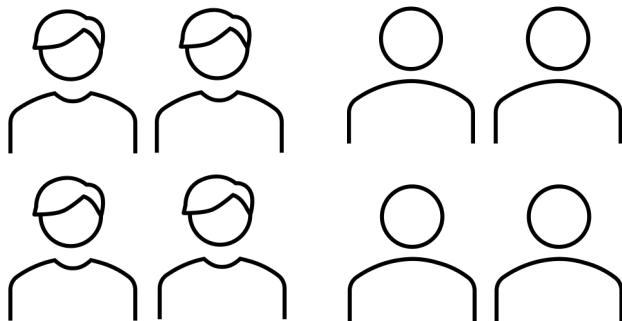
Roadmap

- AI is prone to biases
- What is a fair algorithm
- Distributional Individual Fairness (DIF)
- Enforcing DIF
- Subgroup Fairness
- Learning the Fair Metric

What is a fair algorithm?

Group Fairness:

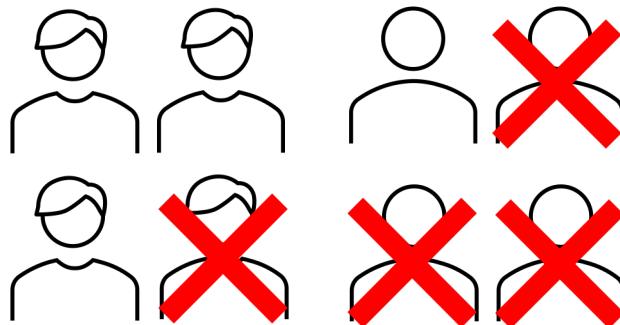
Algorithm is equally good on groups of individuals



What is a fair algorithm?

Group Fairness:

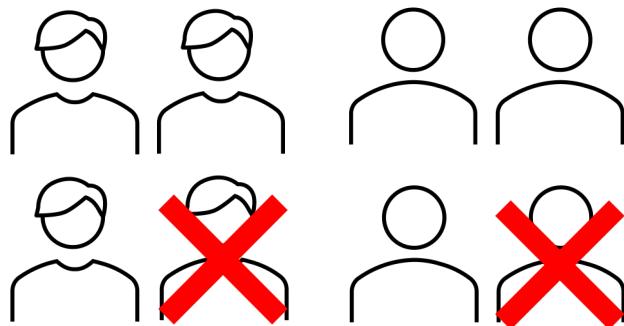
Algorithm is equally good on groups of individuals



What is a fair algorithm?

Group Fairness:

Algorithm is equally good on groups of individuals



What is a fair algorithm?

Group Fairness:

Algorithm is equally good on groups of individuals

Y – true label

A – protected attribute

\hat{Y} – prediction

Statistical Parity: $\hat{Y} \perp\!\!\!\perp A$

Equalized Odds: $\hat{Y}|Y \perp\!\!\!\perp A$

Most prior work is on group fairness

What is a fair algorithm?

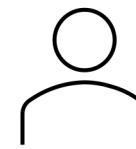
Individual Fairness:

Algorithm treats similar individuals similarly

Ryan earns 50k and lives in a predominantly white zip-code



Tyree earns 50k and lives in a predominantly black zip-code



What is a fair algorithm?

Individual Fairness:

Algorithm treats similar individuals similarly

Ryan earns 50k and lives in a predominantly white zip-code



Tyree earns 50k and lives in a predominantly black zip-code



What is a fair algorithm?

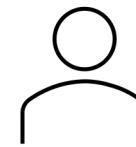
Individual Fairness:

Algorithm treats similar individuals similarly

Ryan earns 50k and lives in a predominantly white zip-code



Tyree earns 50k and lives in a predominantly black zip-code



What is a fair algorithm?

Individual Fairness (Dwork et al. 2012):

$$d_{\mathcal{Y}}(h(x_1), h(x_2)) \lesssim d_{\mathcal{X}}(x_1, x_2) \text{ for all } x_1, x_2 \in \mathcal{X}$$

- ML model is a map $h : \mathcal{X} \rightarrow \mathcal{Y}$
- **fair metric** $d_{\mathcal{X}}$ measures similarity between inputs
- $d_{\mathcal{Y}}$ measures similarity between outputs

Why Individual Fairness?

Example: Sentiment analysis – classify words as positive or negative

Positive: *admire, adorable, joy, lucky, talented, ...* 

Negative: *aggressive, distrust, nasty, radical, ...* 

Why Individual Fairness?

Example: Sentiment analysis – classify words as positive or negative

Positive: *admire, adorable, joy, lucky, talented, ...* 

Negative: *aggressive, distrust, nasty, radical, ...* 

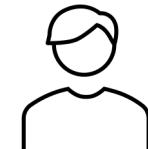
Deep Learning + Word Embeddings -> **95%** test accuracy.



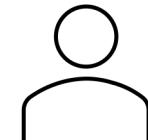
Why Individual Fairness?

What is a sentiment of a name?

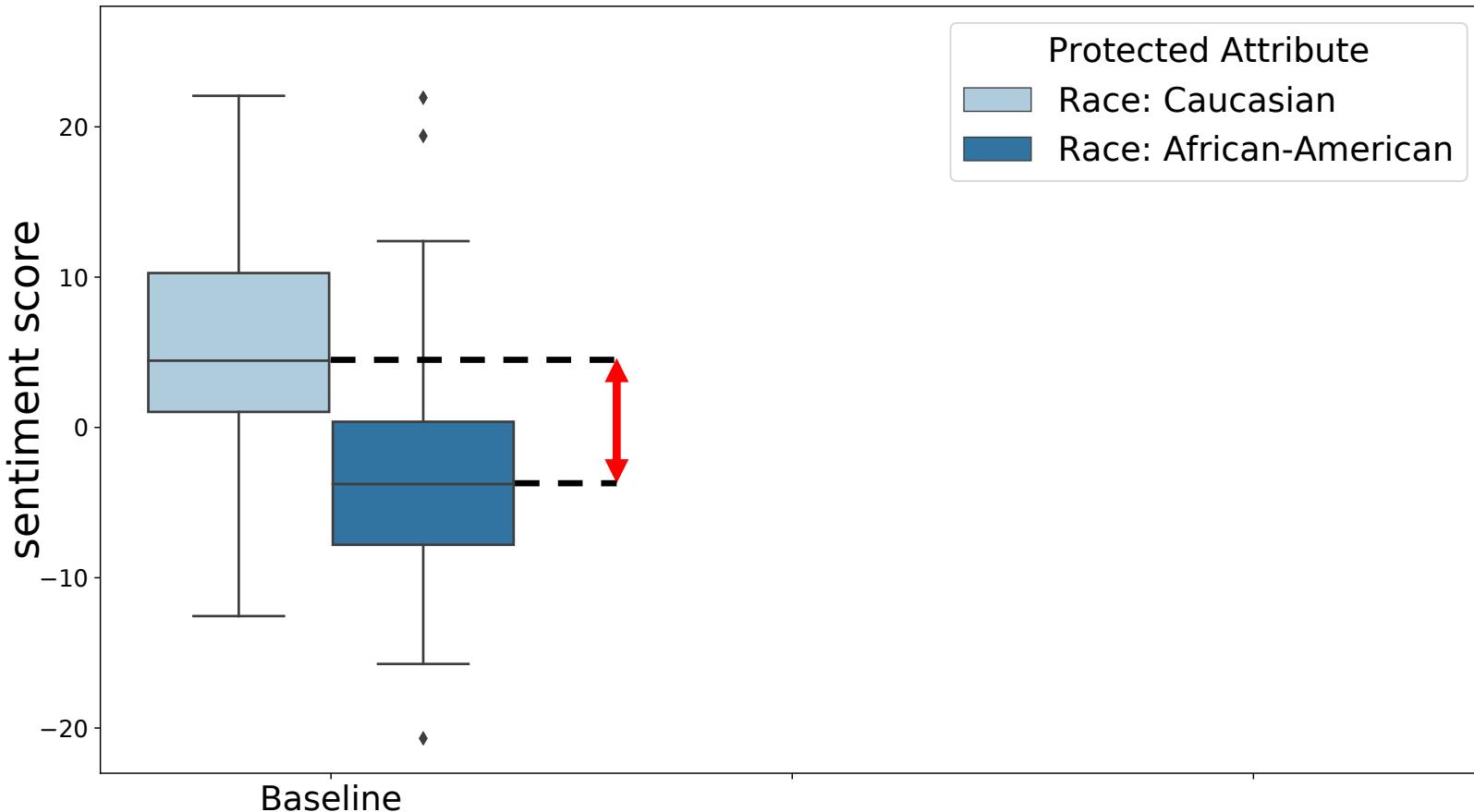
European-American names: *Adam, Ryan, Paul, ... , Courtney, Meredith, Megan, ...*



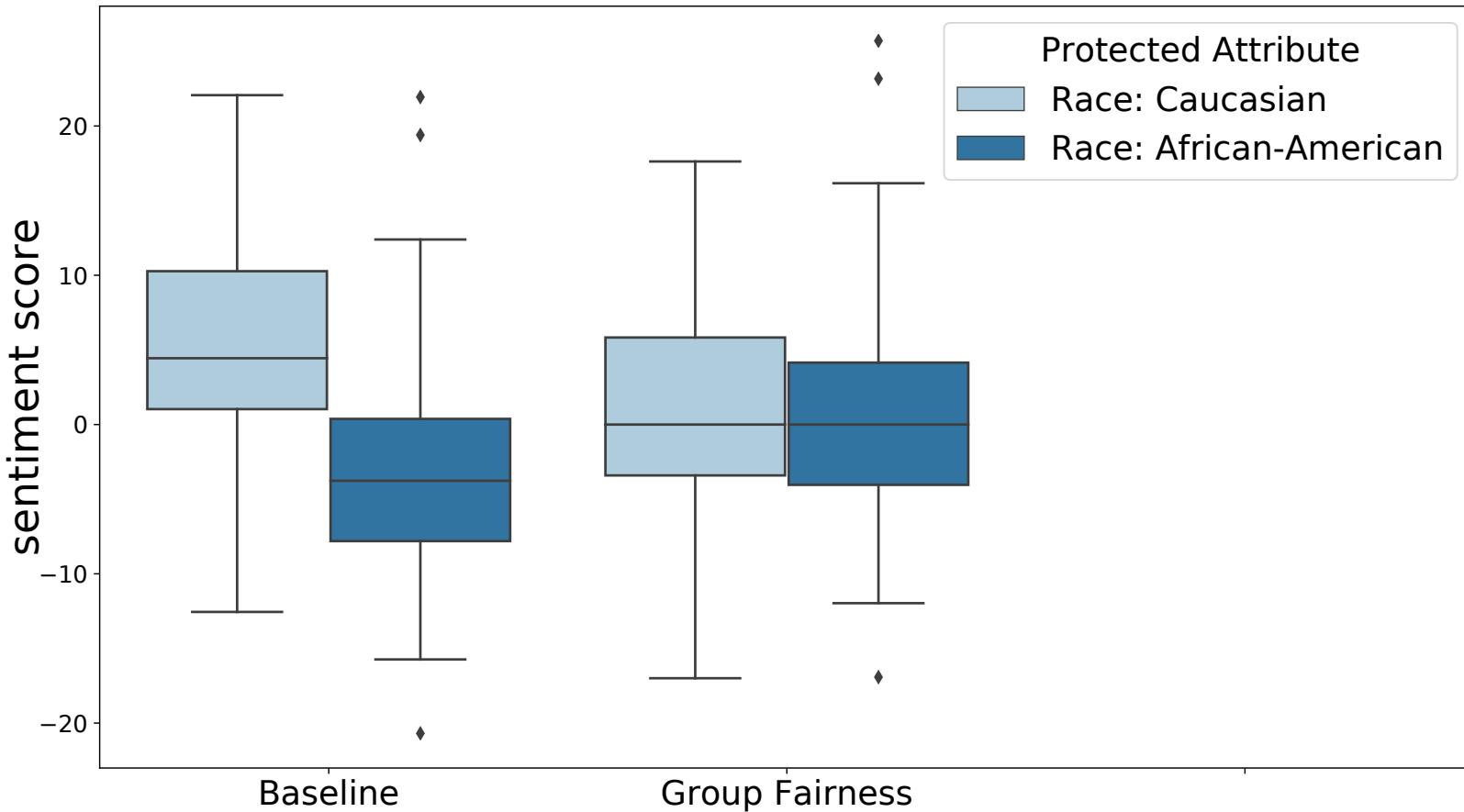
African-American names: *Alonzo, Leroy, Tyree, ... , Shereen, Sharise, Tawanda, ...*



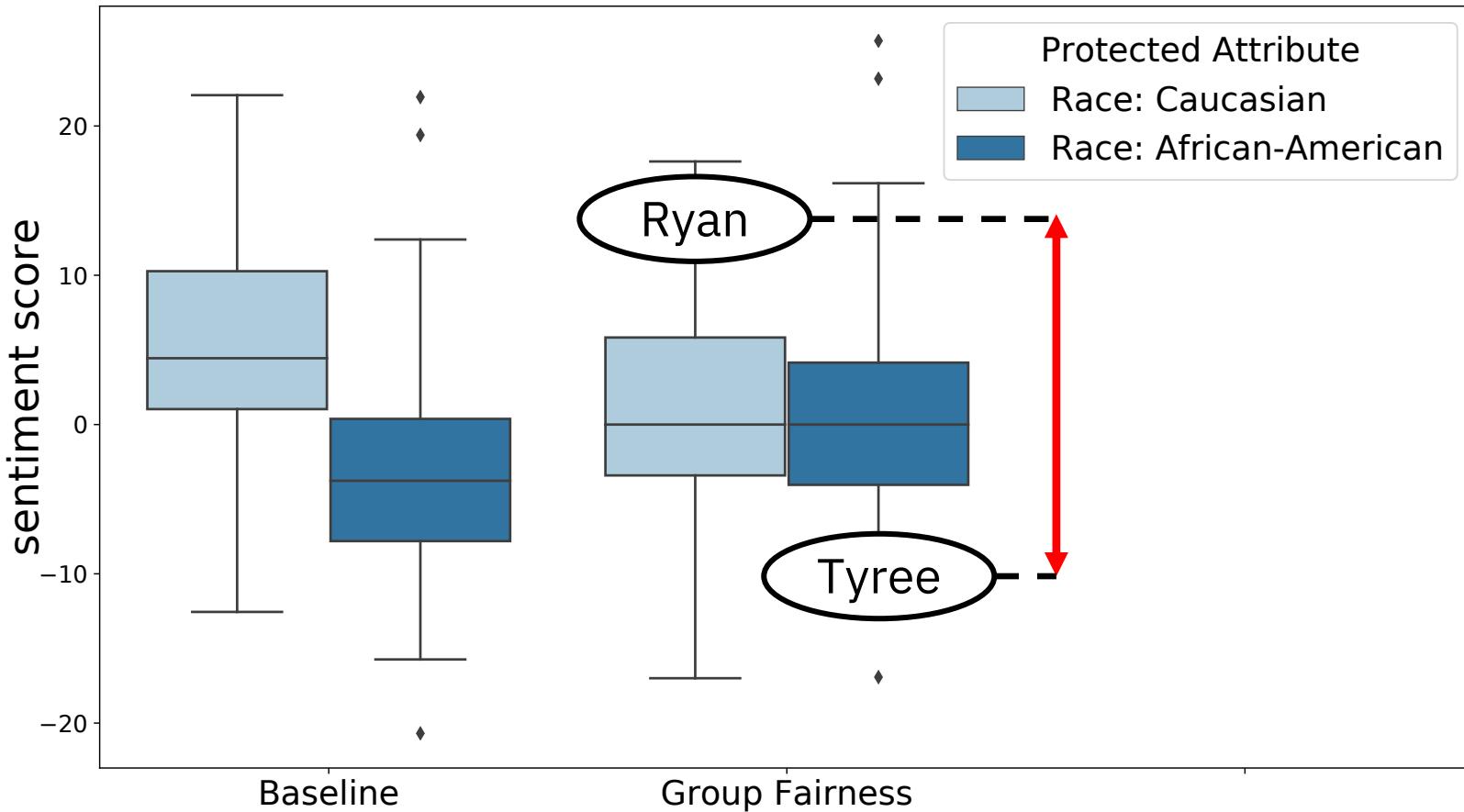
Why Individual Fairness?



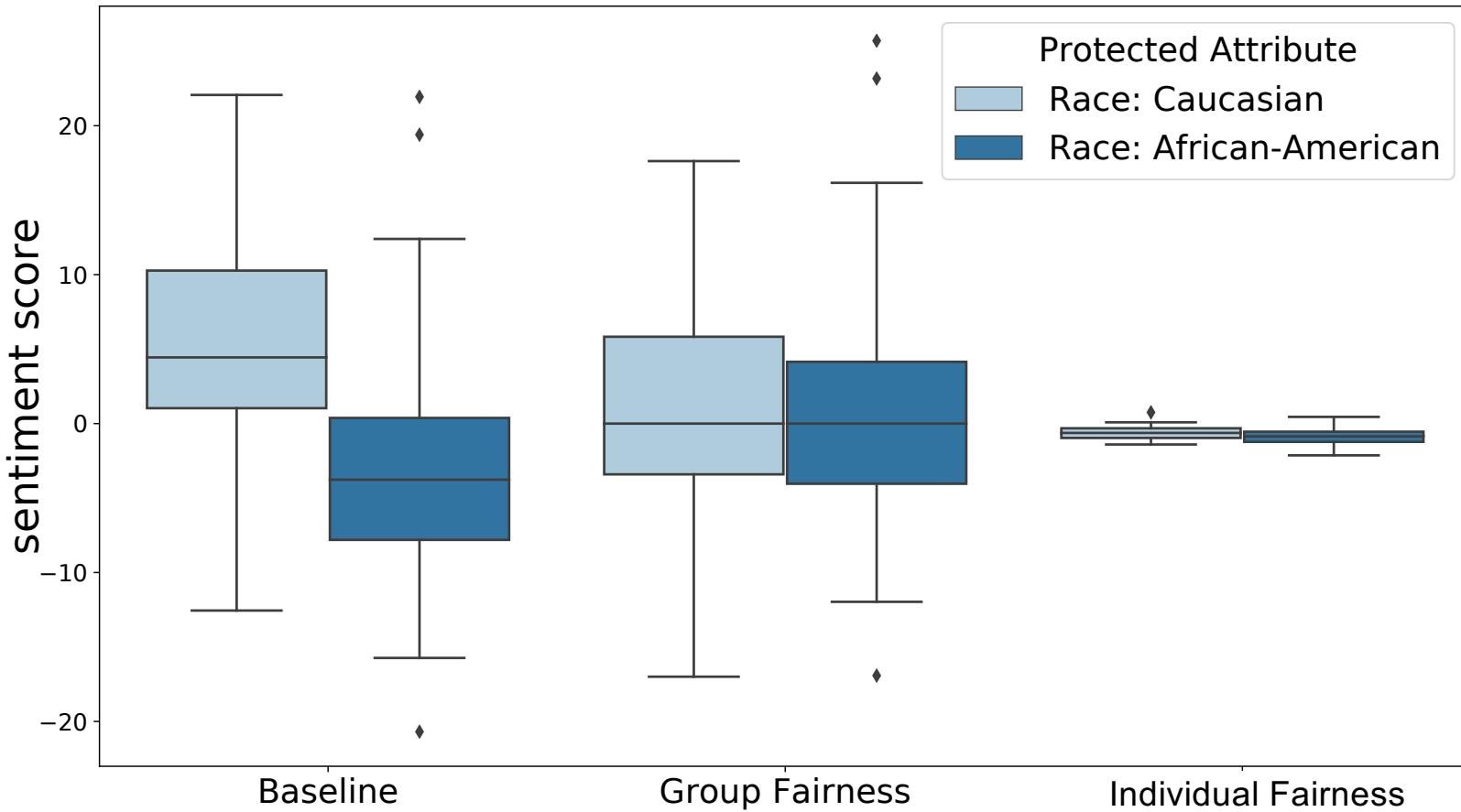
Why Individual Fairness?



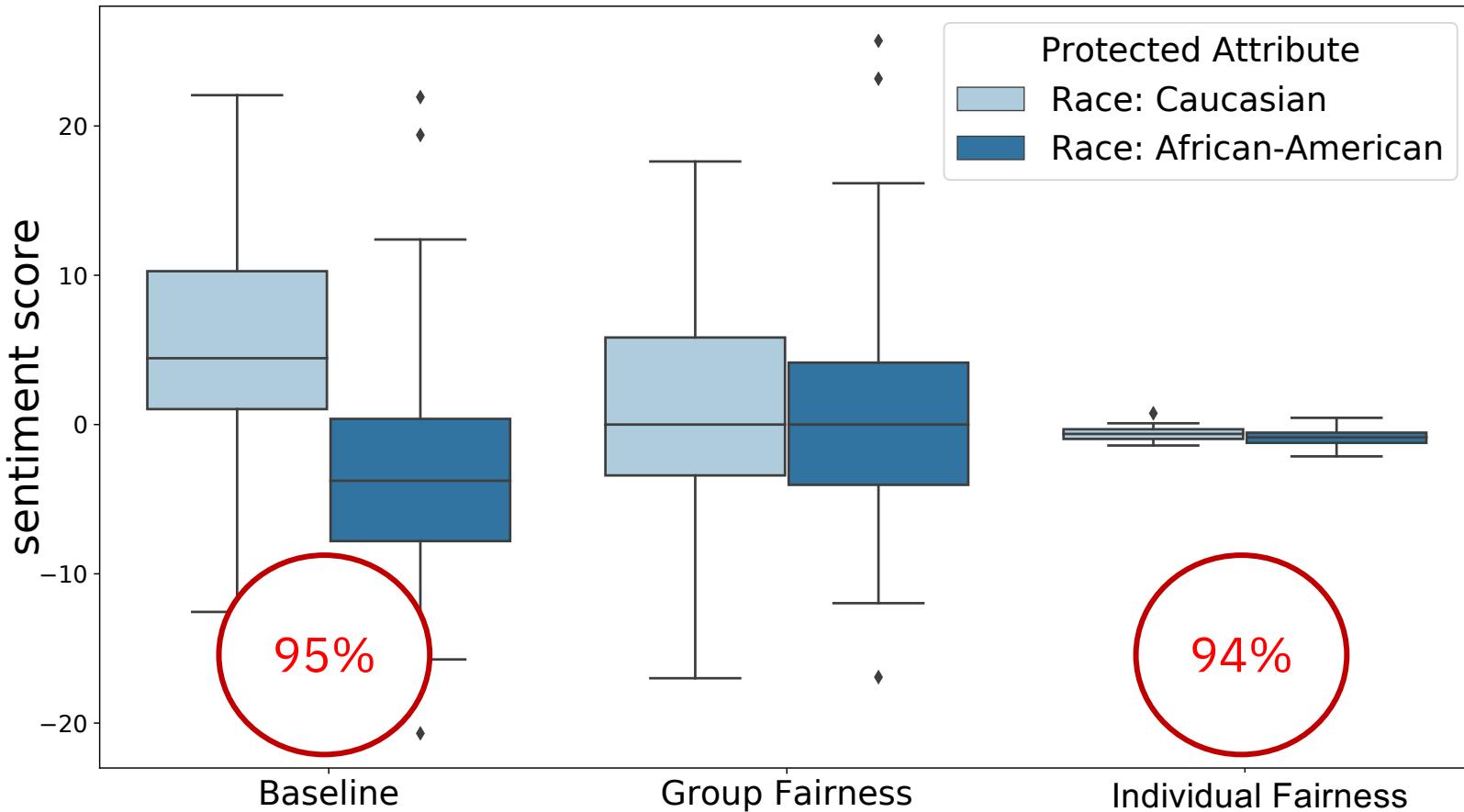
Why Individual Fairness?



Why Individual Fairness?



Why Individual Fairness?



Roadmap

- AI is prone to biases
- What is a fair algorithm
- Distributional Individual Fairness (DIF)
- Enforcing DIF
- Subgroup Fairness
- Learning the Fair Metric

Assessing Individual Fairness

$$\text{IF}(h, x_i) \triangleq \left\{ \begin{array}{ll} \max_{x'_i \in \mathcal{X}} & d_{\mathcal{Y}}(h(x_i), h(x'_i)) \\ \text{subject to} & d_{\mathcal{X}}(x_i, x'_i) \leq \epsilon \end{array} \right\}$$

- $d_{\mathcal{X}}$: fair metric
- $d_{\mathcal{Y}}$: metric on outputs
- ϵ : small tolerance parameter

Distributional Individual Fairness (DIF)

$$\text{DIF}(h) \triangleq \left\{ \begin{array}{ll} \sup_{T:\mathcal{X} \rightarrow \mathcal{X}} & \mathbb{E}_{P_X} [d_{\mathcal{Y}}(h(x), h(T(x)))] \\ \text{subject to} & \mathbb{E}_{P_X} [d_{\mathcal{X}}(x, T(x))] \leq \epsilon. \end{array} \right\}$$

- P_X (marginal) distribution of inputs
- Optimal T maps x_i to x'_i
- Constraint enforces $d_{\mathcal{X}}(x_i, x'_i) \leq \epsilon$ **on average**

Comparing DIF and IF

Individual fairness map is feasible, but may not be optimal:

$$T_{\text{IF}}(x_i) \triangleq \arg \max_{d_{\mathcal{X}}(x_i, x'_i) \leq \epsilon} d_{\mathcal{Y}}(h(x_i), h(x'_i)).$$

DIF may transport some points by more than ϵ .

IF does not imply DIF or vice a versa!

Comparing DIF and IF

Individual fairness map is feasible, but may not be optimal:

$$T_{\text{IF}}(x_i) \triangleq \arg \max_{d_{\mathcal{X}}(x_i, x'_i) \leq \epsilon} d_{\mathcal{Y}}(h(x_i), h(x'_i)).$$

Let $\text{DIF}(h) < \delta$, then

$$P_X(d_{\mathcal{Y}}(h(x), h(T_{\text{IF}}(x))) \geq \tau) \leq \frac{\delta}{\tau} \text{ for any } \tau > 0.$$

DIF implies IF with high probability

DIF in Social Science

$$\text{DIF}(h) \triangleq \left\{ \begin{array}{ll} \sup_{T: \mathcal{X} \rightarrow \mathcal{X}} & \mathbb{E}_{P_X} [d_{\mathcal{Y}}(h(x), h(T(x)))] \\ \text{subject to} & \mathbb{E}_{P_X} [d_{\mathcal{X}}(x, T(x))] \leq \epsilon. \end{array} \right\}$$

- P_X (marginal) distribution of inputs
- Optimal T maps x_i to x'_i
- Large $\text{DIF}(h)$ implies unfairness

DIF in Social Science

Bertrand & Mullainathan (2004) studied racial bias in the US labor market.

- P_X (marginal) distribution of inputs
- Optimal T maps x_i to x'_i
- Large $\text{DIF}(h)$ implies unfairness

DIF in Social Science

Bertrand & Mullainathan (2004) studied racial bias in the US labor market.

- The investigators responded to job ads in Boston and Chicago newspapers with fictitious resumes.
- Optimal T maps x_i to x'_i
- Large $\text{DIF}(h)$ implies unfairness

DIF in Social Science

Bertrand & Mullainathan (2004) studied racial bias in the US labor market.

- The investigators responded to job ads in Boston and Chicago newspapers with fictitious resumes.
- They randomly assigned African-American or white sounding names to the resumes.
- Large $\text{DIF}(h)$ implies unfairness

DIF in Social Science

Bertrand & Mullainathan (2004) studied racial bias in the US labor market.

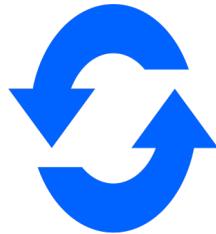
- The investigators responded to job ads in Boston and Chicago newspapers with fictitious resumes.
- They randomly assigned African-American or white sounding names to the resumes.
- The investigators concluded there is discrimination against African-Americans: the resumes assigned white names received 50% more callbacks for interviews.

Roadmap

- AI is prone to biases
- What is a fair algorithm
- Distributional Individual Fairness (DIF)
- Enforcing DIF
- Subgroup Fairness
- Learning the Fair Metric

Training DIF models

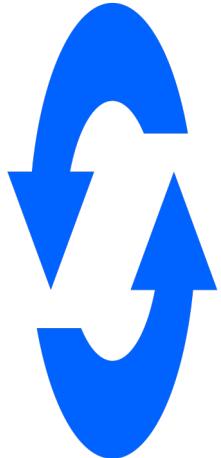
Usual AI: train algorithm accurate on the available data



- Observe data
- Update parameters to minimize prediction error
- Repeat

Training DIF models

Individually Fair AI: train algorithm accurate on the available data **and** all possible similar data



- Observe data
- Generate similar data where algorithm performs differently to evaluate DIF
- Update parameters to minimize prediction error **and** DIF
- Repeat

Training DIF models

$$\min_{h \in \mathcal{H}} L(h) + \rho \text{DIF}(h)$$

$$L(h) \triangleq \mathbb{E}[\ell(y, h(x))]$$

- \mathcal{H} : model class (e.g. neural nets with a certain architecture)
- $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ is a loss function
- ρ : regularization parameter

Training DIF models

$$\text{DIF}(h) = \inf_{\lambda > 0} \{ \lambda \epsilon + \mathbb{E}_{P_X} [r_\lambda(h, x)] \}$$

$$r_\lambda(h, x) \triangleq \sup_{x' \in \mathcal{X}} \{ d_{\mathcal{Y}}(h(x), h(x')) - \lambda d_{\mathcal{X}}(x, x') \}$$

DIF is infinite-dimensional, but its dual is univariate

Training DIF models

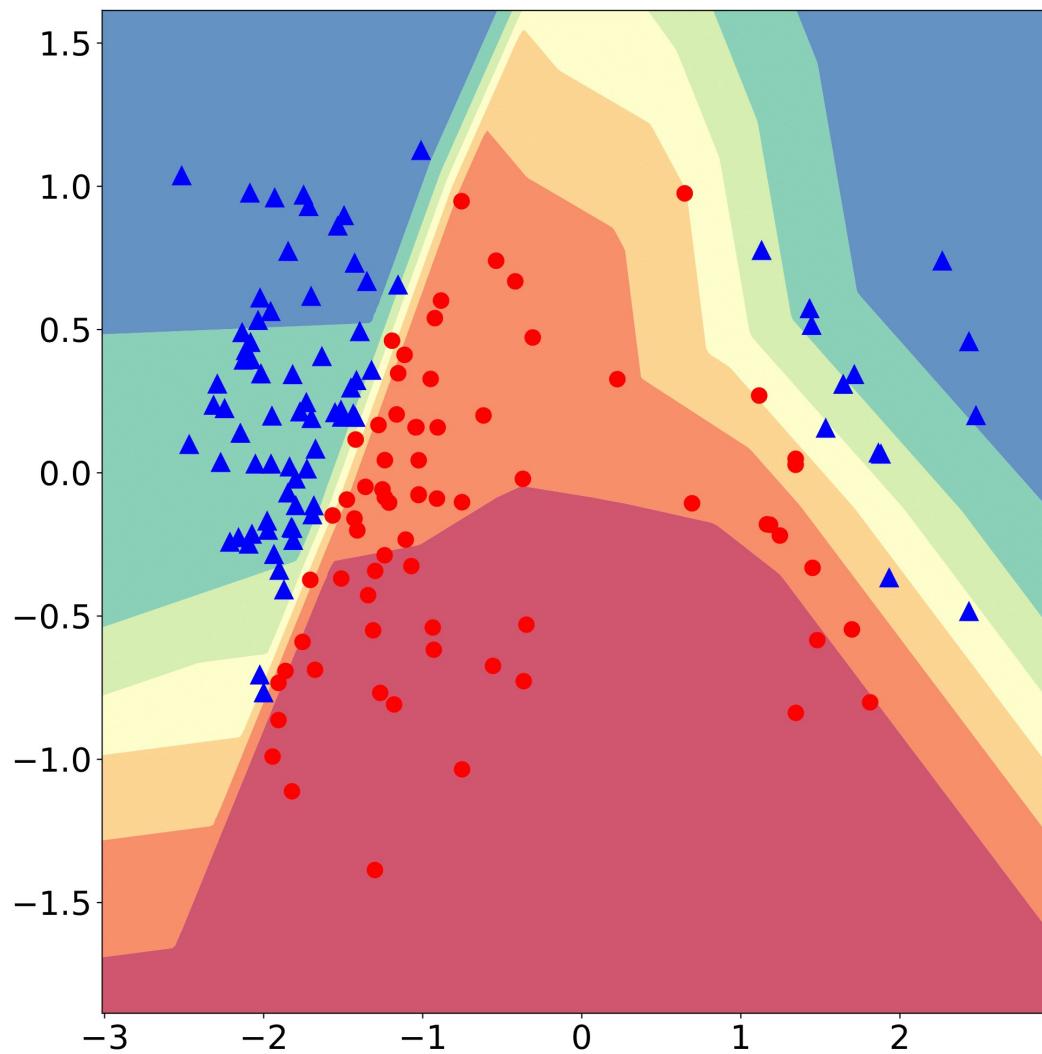
$$\min_{\theta \in \Theta, \lambda \in \mathbb{R}_+} \mathbb{E}[f(\theta, (x, y))]$$

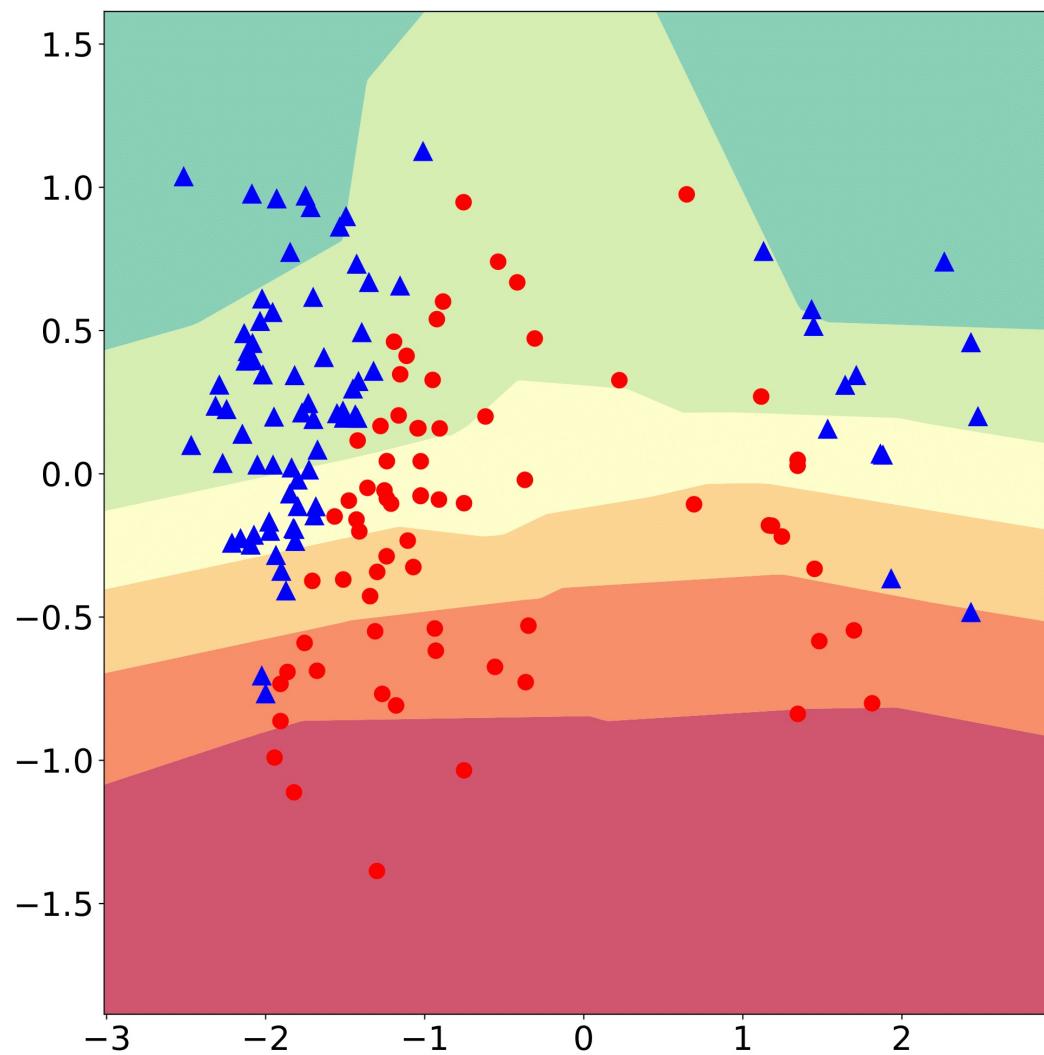
$$f(\theta, (x, y)) \triangleq \ell(y, h_\theta(x)) + \rho(\lambda \epsilon + r_\lambda(h_\theta, x))$$

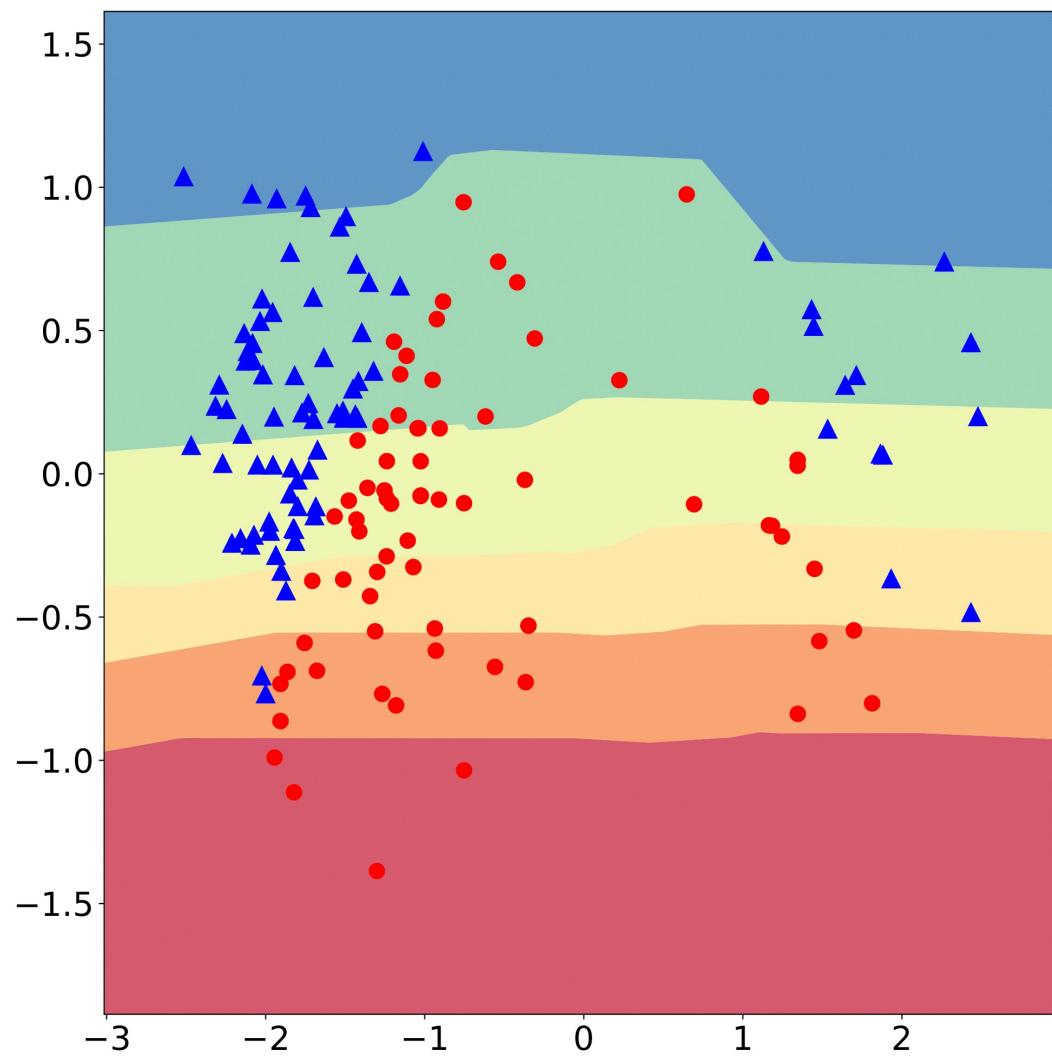
DIF is amendable to stochastic optimization

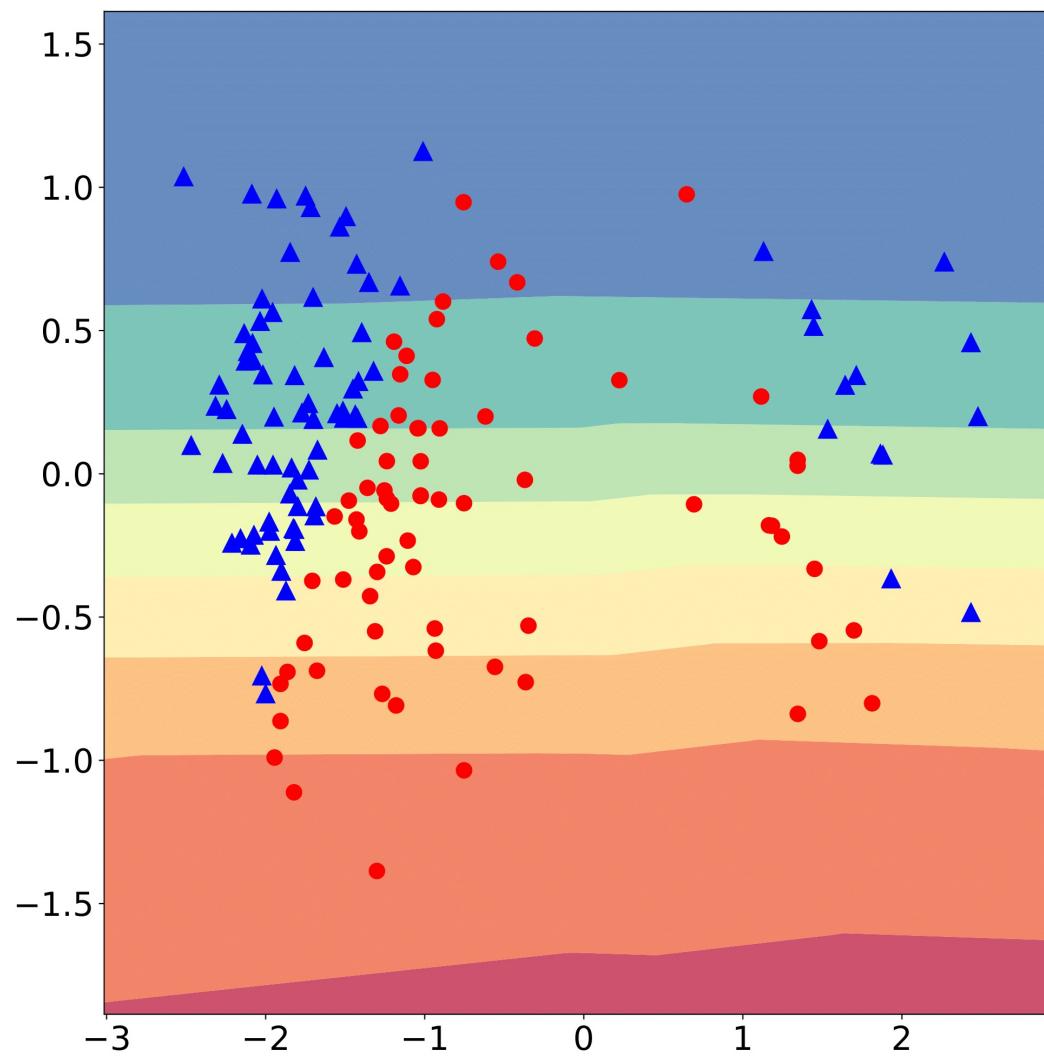
DIF Example

- Points on horizontal lines (identical y -values) are similar
- Training data is biased: $P_{Y|X}$ not constant on horizontal lines
- Use DIF regularizer to correct bias in training data



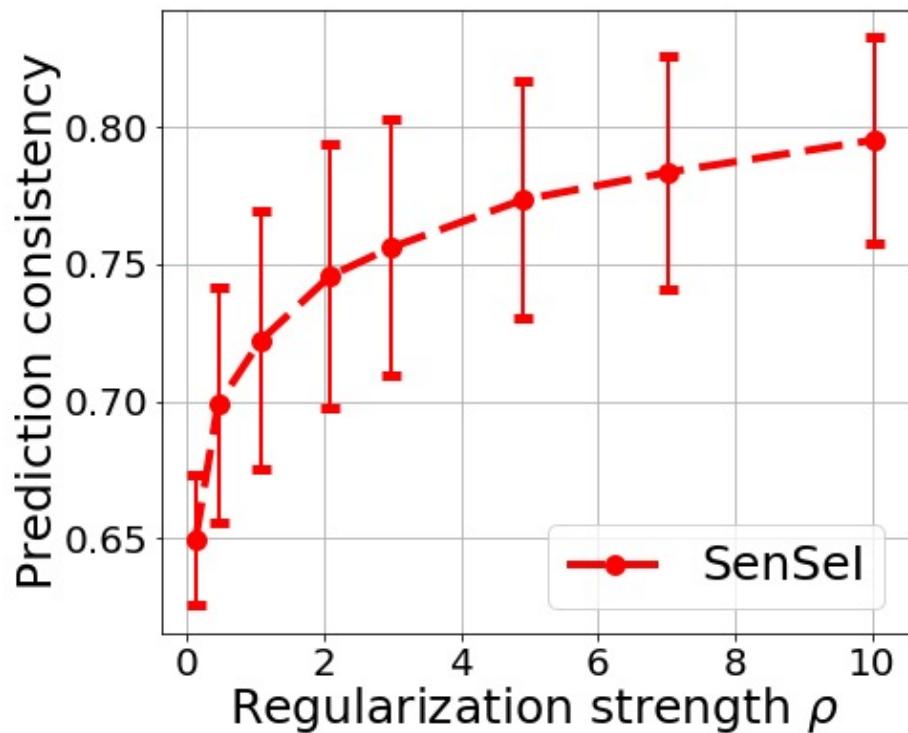
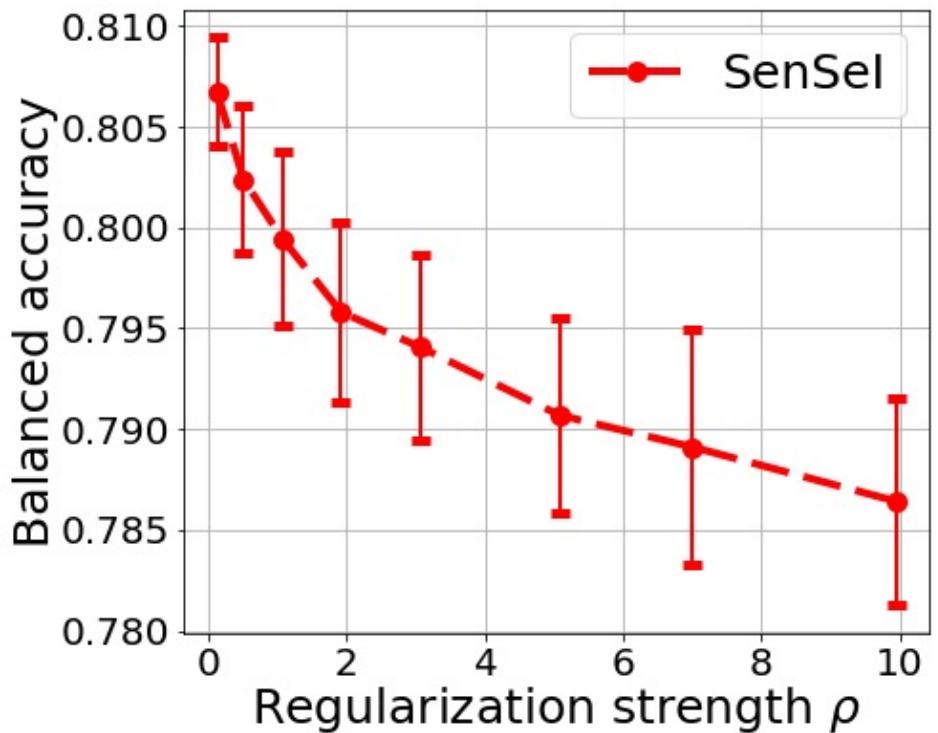






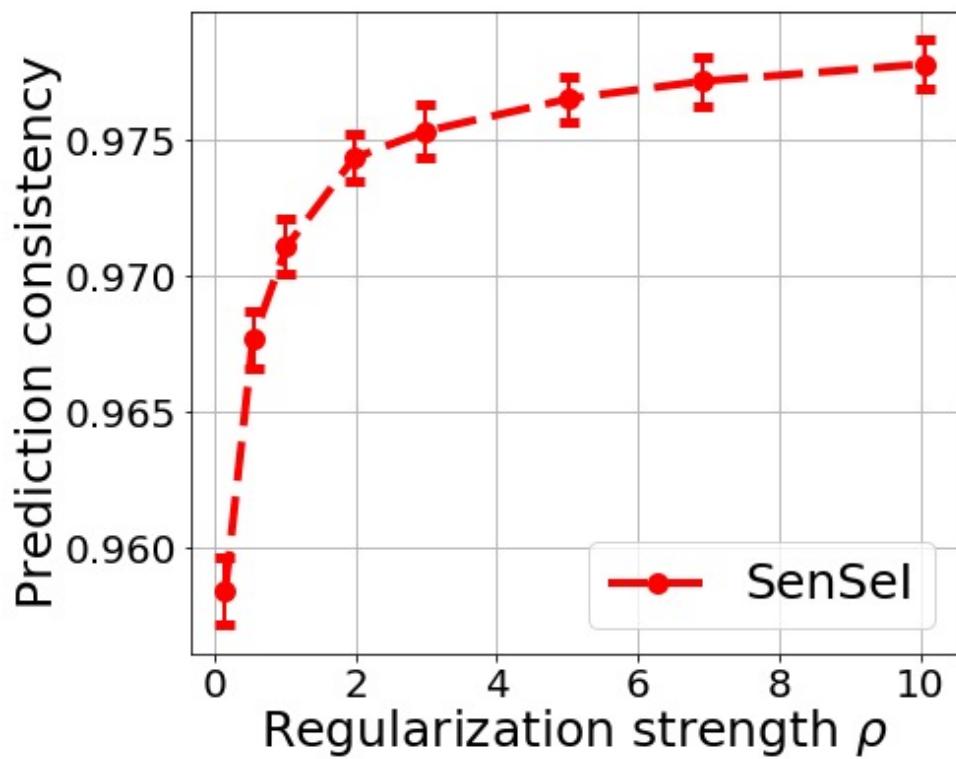
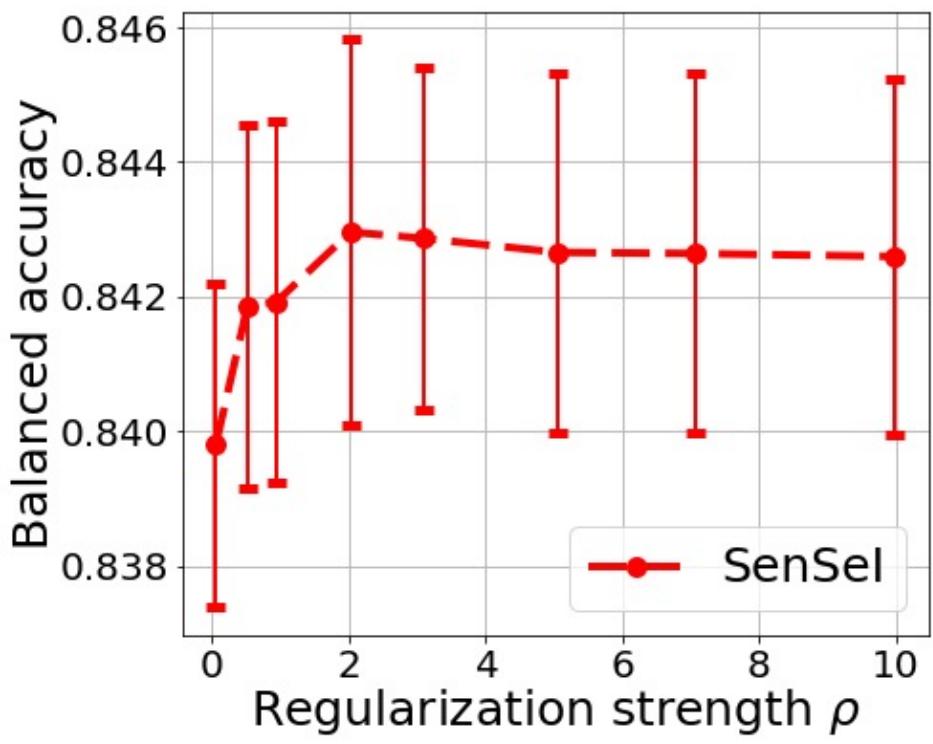
Achieving IF on toxicity classification

- Data: “Toxic Comment Classification Challenge”. 165k text comments
- Task: predict if a comment is toxic or not
- Measuring IF: does prediction change when changing identity tokens.
Are predictions for “Some people are gay” and “Some people are straight” the same?
- 50 identity tokens: prediction on all variations should be the same to satisfy IF



Achieving IF on occupation prediction

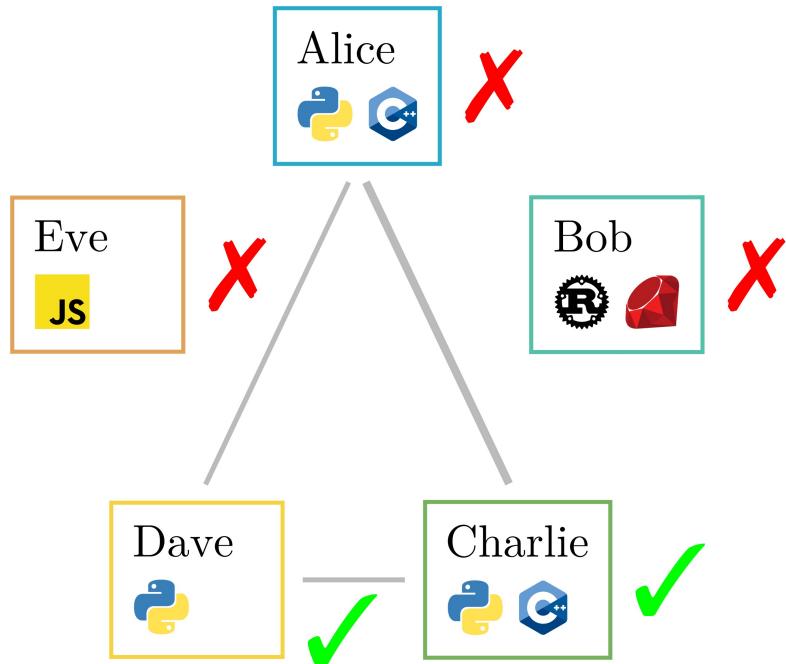
- Data: Bias in Bios (De-Arteaga et al. 2019). 400k textual bio descriptions
- Task: predict one of the 28 occupations from a bio
- Measuring IF: does prediction change when changing names and gender pronouns in a bio?



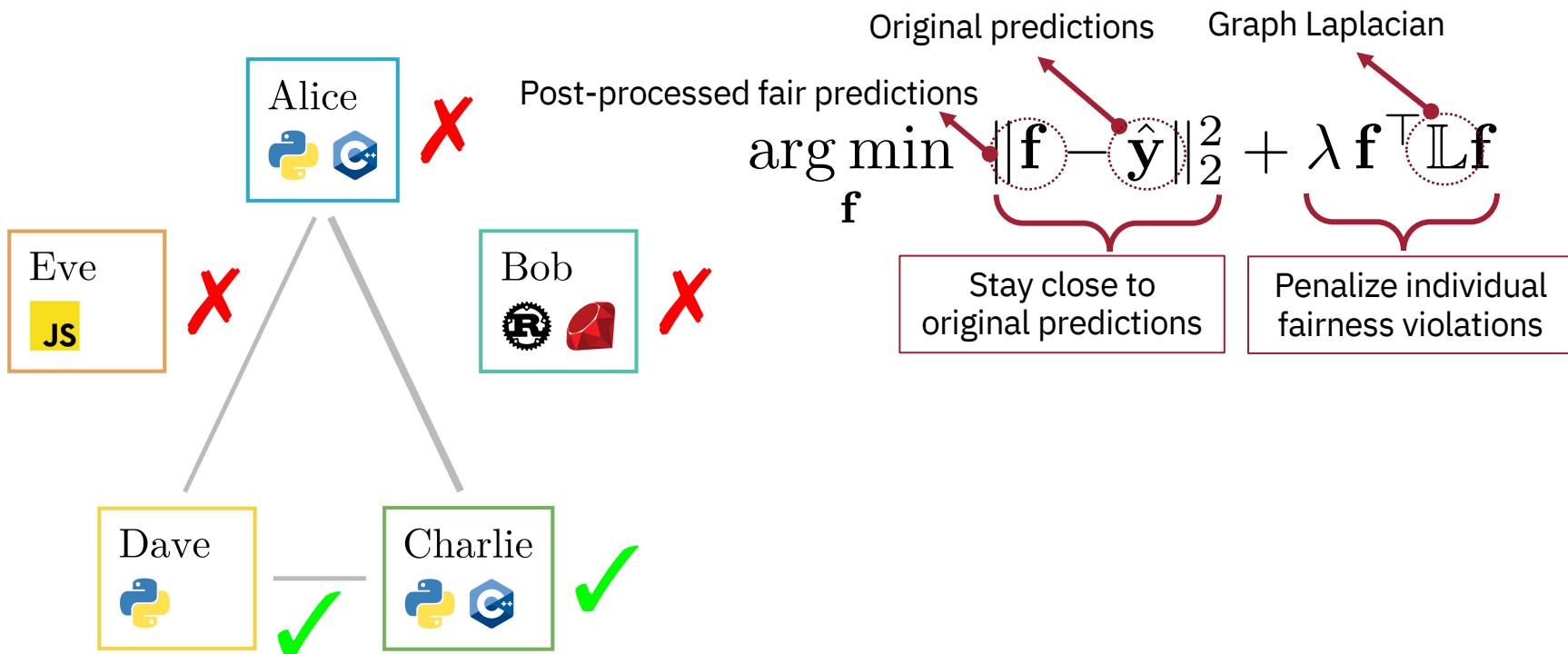
Post-processing for Individual Fairness

- Works with any trained model without re-training
- Fast and easy to implement
- Does not require knowledge of the model parameters
- Only needs access to the outputs
- Key idea: represent individuals as a graph and quantify fairness with the graph Laplacian quadratic form

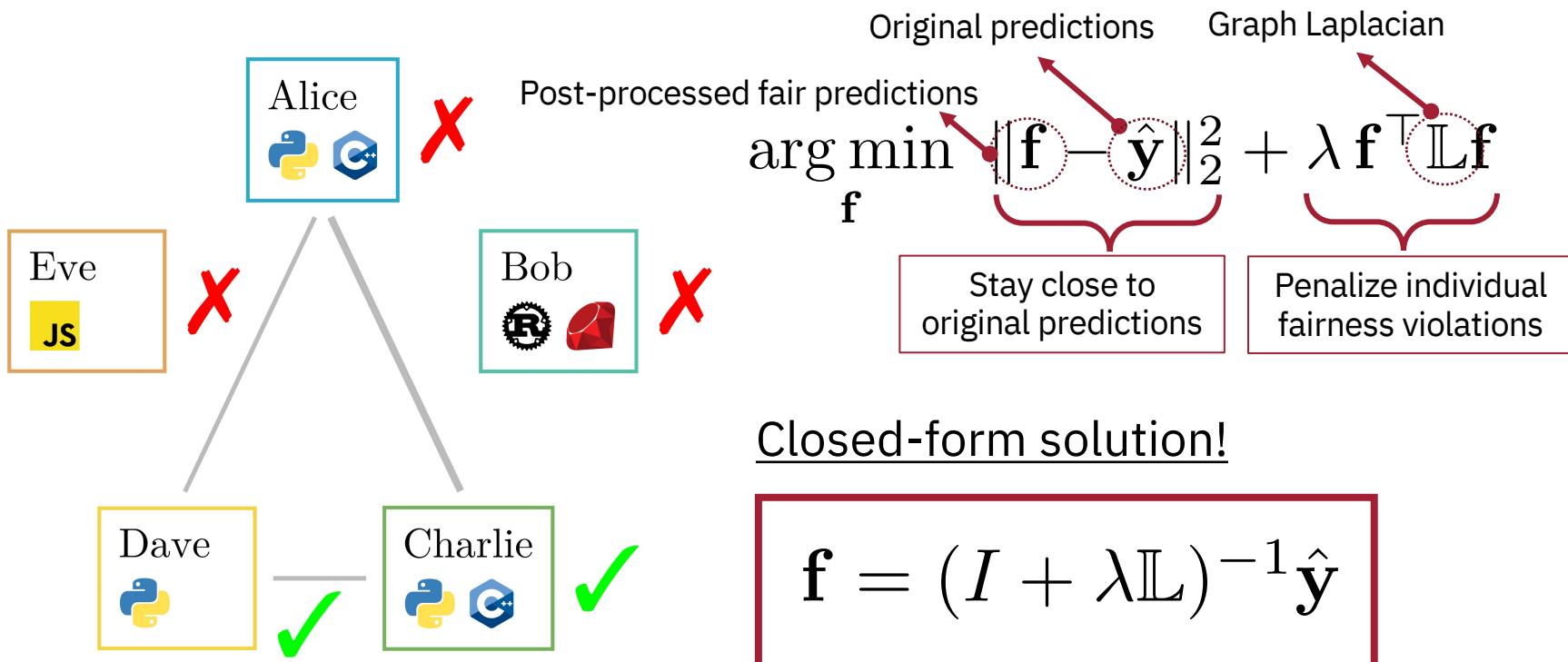
Post-processing for Individual Fairness



Post-processing for Individual Fairness



Post-processing for Individual Fairness



Roadmap

- AI is prone to biases
- What is a fair algorithm
- Distributional Individual Fairness (DIF)
- Enforcing DIF
- Subgroup Fairness
- Learning the Fair Metric

Fairness problems Today



goldman sachs women credit



Google

apple women credit



The @AppleCard is such a [REDACTED] sexist program. My wife and I filed joint tax returns, live in a community-property state, and have been married for a long time. Yet Apple's black box algorithm thinks I deserve 20x the credit limit she does. No appeals work.

1:34 PM · Nov 7, 2019 · Twitter for iPhone

9K Retweets 3.5K Quote Tweets 28K Likes

Nov 9, 2019 — A Wall Street regulator is opening a probe into Goldman Sachs Group Inc.'s credit card practices after a viral tweet from a tech entrepreneur ...

Nov 11, 2019 — Danish entrepreneur David Heinemeier Hansson says his credit limit is ... differences in Apple Card credit lines for male and female customers.

Fairness problems Today



goldman sachs women credit



Google



All

News

Images

Shopping

Videos

More

Settings

Tools

About 6,810,000 results (0.55 seconds)

www.nytimes.com › Apple-credit-card-investigation

Apple Card Investigated After Gender Discrimination Complaints

Nov 10, 2019 — The card, a partnership between Apple and **Goldman Sachs**, made its ... on the Apple Card's treatment of **female** credit applicants, he said his ...

www.washingtonpost.com › business › 2019/11/11 › a... ▾

Apple Card algorithm sparks gender bias inquiry - The ...

Nov 11, 2019 — Danish entrepreneur David Heinemeier Hansson says his credit limit is ...

Apple Card algorithm sparks gender bias allegations against **Goldman Sachs** ... differences in Apple Card credit lines for male and **female** customers.

People also search for

apple card discrimination cnn goldman apple credit card
goldman sachs political bias apple discrimination lawsuit
how many apple card users amazon scraps secret ai recruiting tool

www.engadget.com › 2019-11-12-goldman-sachs-credi... ▾

Goldman will re-check Apple Card credit scores after sexism ...

Nov 12, 2019 — ... after accusations that the bank behind the program, **Goldman Sachs**, has been discriminatory against **women** in its provision of **credit** lines.

www.bloomberg.com › news › articles › viral-tweet-ab... ▾

Goldman Sachs Probed After Viral Tweet About Apple Card ...

Nov 9, 2019 — A Wall Street regulator is opening a probe into **Goldman Sachs** Group Inc.'s **credit** card practices after a viral tweet from a tech entrepreneur ...

apple women credit



Google



All

News

Images

Videos

More

Settings

Tools

About 350,000,000 results (0.46 seconds)

www.nytimes.com › Apple-credit-card-investigation

Apple Card Investigated After Gender Discrimination Complaints

Nov 10, 2019 — A prominent software developer said on Twitter that the **credit** card was "sexist" against **women** applying for **credit**. Something curious happened when a husband and wife recently compared their **Apple Card** spending limits.

www.theverge.com › 2019/11/11 › apple-credit-card-g... ▾

Apple's credit card is being investigated for discriminating ...

Nov 11, 2019 — **Apple's** **credit** card is being investigated by financial regulators after customers accused its lending algorithms of discriminating against **women**.

www.cnn.com › 2019/11/12 › business › apple-card-gend... ▾

Apple Card is accused of gender bias. Here's how that can ...

Nov 12, 2019 — Some **Apple Card** customers say the **credit** card's issuer, **Goldman Sachs**, is giving **women** far lower **credit** limits, even if they share assets and ...

www.wired.com › story › the-apple-card-didnt-see-gen... ▾

The Apple Card Didn't 'See' Gender—and That's the Problem ...

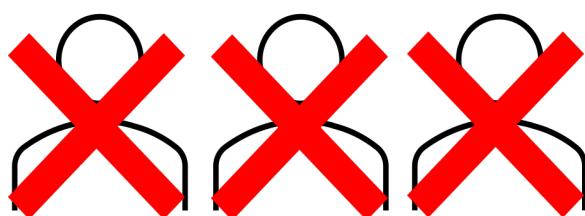
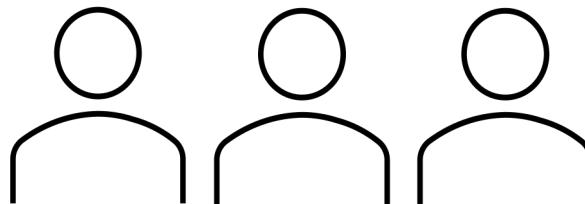
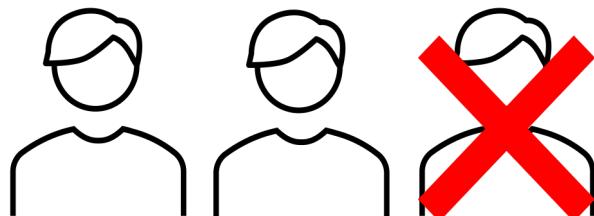
Nov 19, 2019 — The **Apple** **credit** card, launched in August, ran into major problems last ... A gender-blind algorithm could end up biased against **women** as ...

www.washingtonpost.com › business › 2019/11/11 › a... ▾

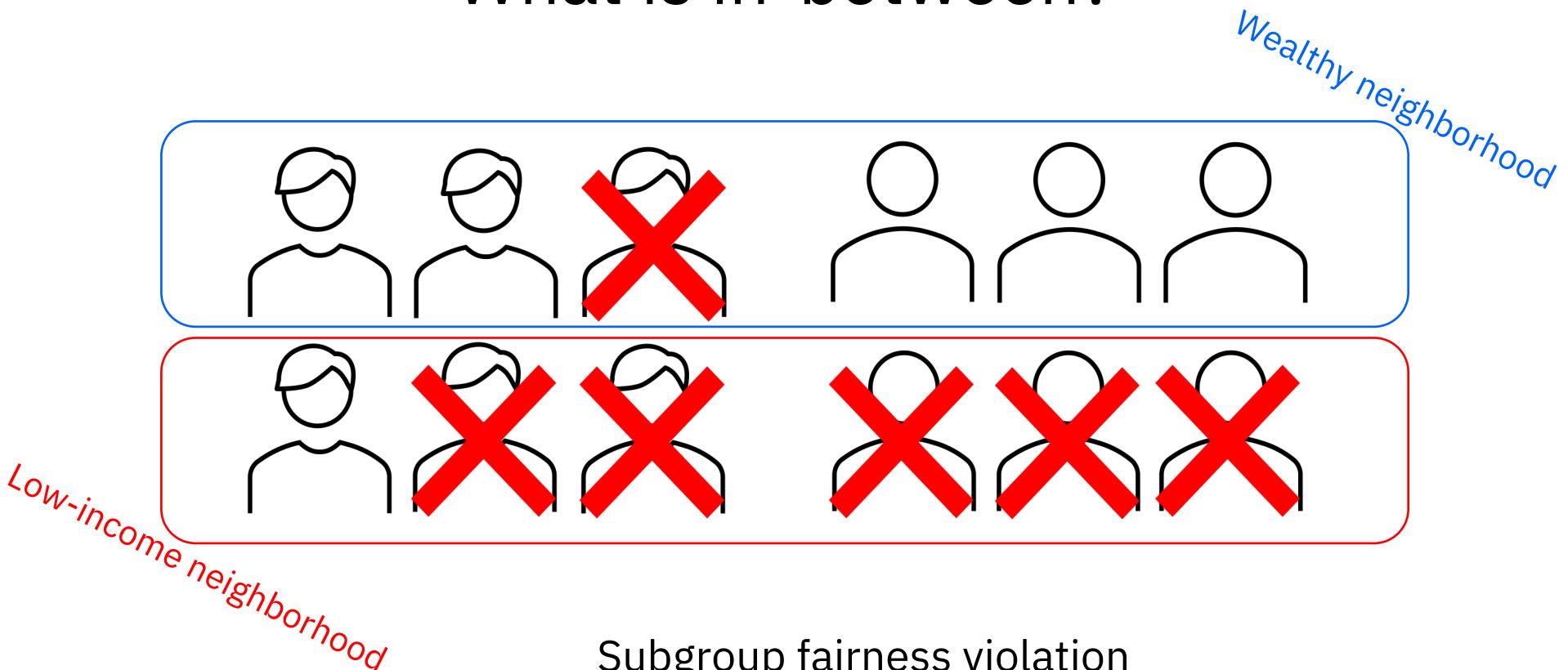
Apple Card algorithm sparks gender bias inquiry - The ...

Nov 11, 2019 — Danish entrepreneur David Heinemeier Hansson says his **credit** limit is ... differences in **Apple Card** **credit** lines for male and **female** customers.

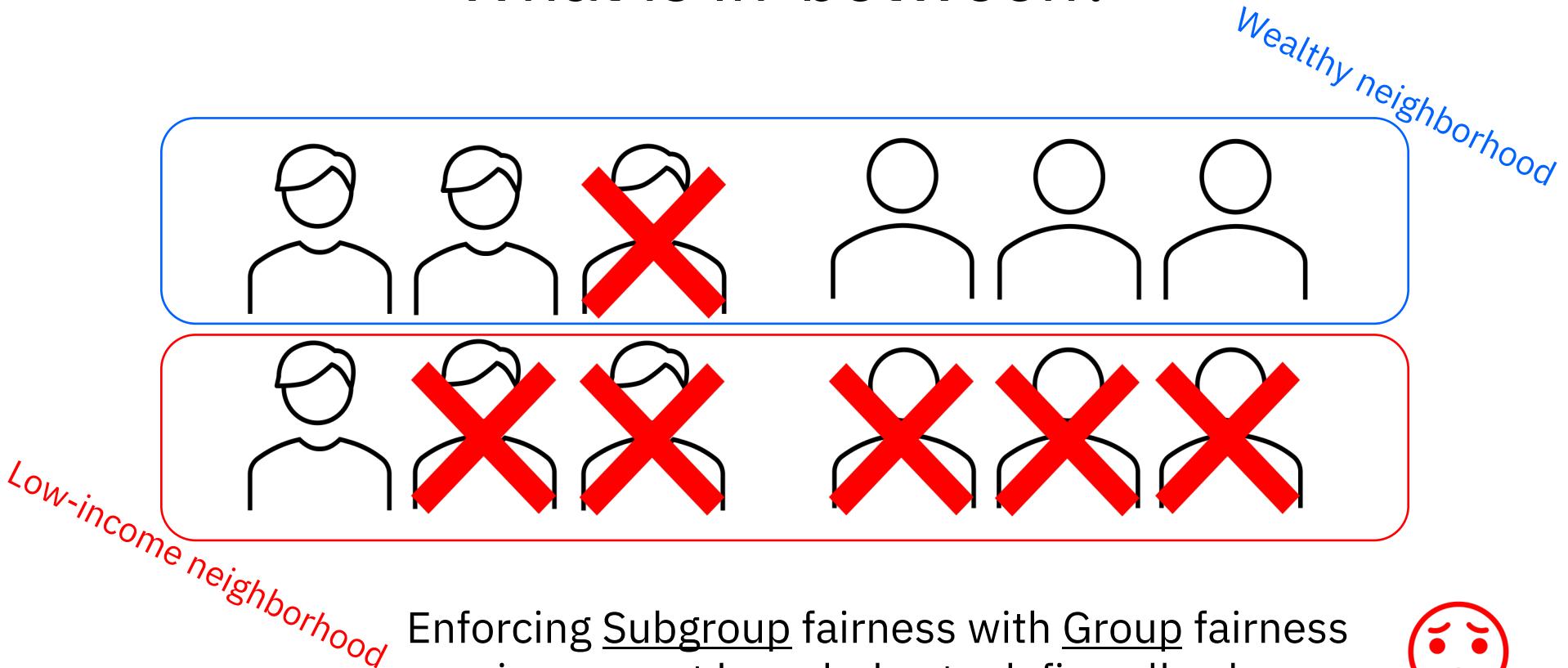
What is in-between?



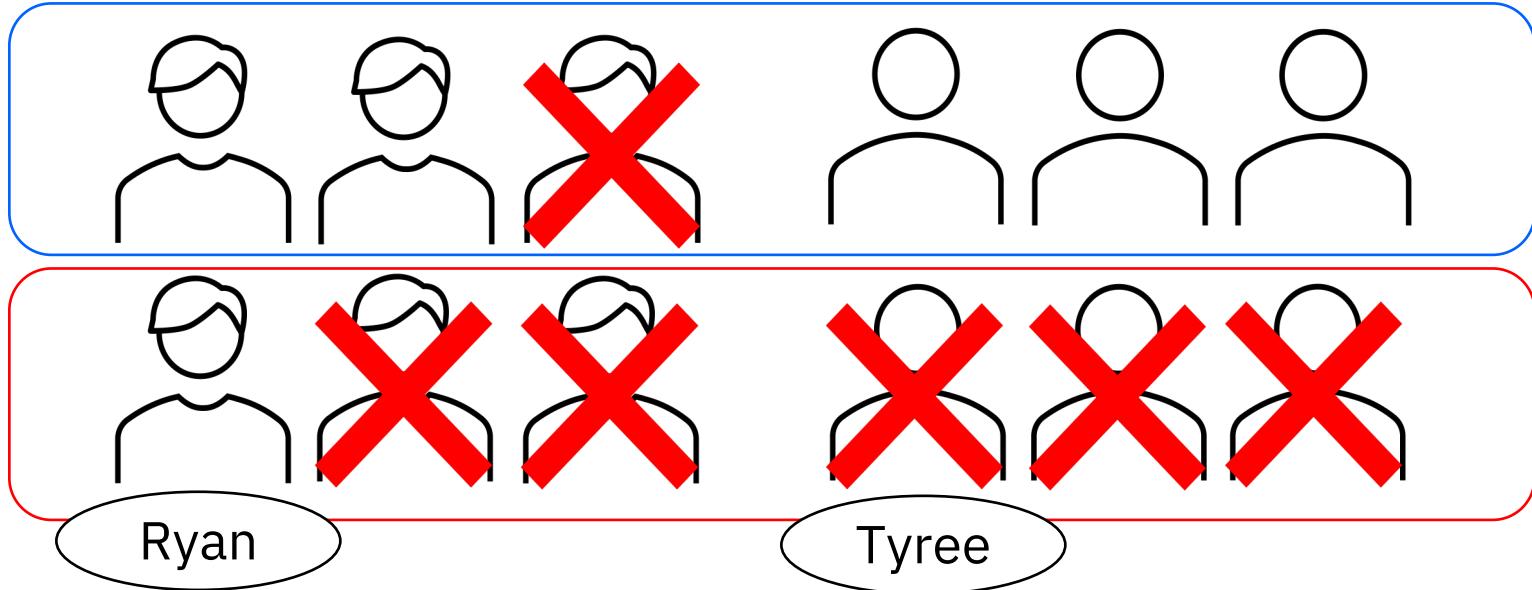
What is in-between?



What is in-between?



What is in-between?



Subgroup fairness can be enforced with Individual fairness generalizing to unforeseen subgroups



Example: Income prediction

Predict individuals' income based on Census data (ADULT dataset).

Evaluate SenSR against baseline and Adversarial Debiasing (group fairness).

	B-Acc, %	S-Con.	Gap _G ^{RMS}	Gap _R ^{RMS}	Gap _G ^{max}	Gap _R ^{max}
SenSR+EXPLORE	79.4	.966	.065	.044	.084	.059
SenSR	78.9	.934	.068	.055	.087	.067
Baseline	82.9	.848	.179	.089	.216	.105
Adv. debiasing	81.5	.807	.082	.070	.110	.078

Spousal Consistency

Spousal consistency measures counterfactual sensitivity to “relationship” status, i.e. how often the classification remains unchanged when relationship status is perturbed.

SenSR + EXPLORE is best for S-Con; the group fairness method is worse than baseline.

*SenSR is our earlier algorithm for enforcing individual fairness.

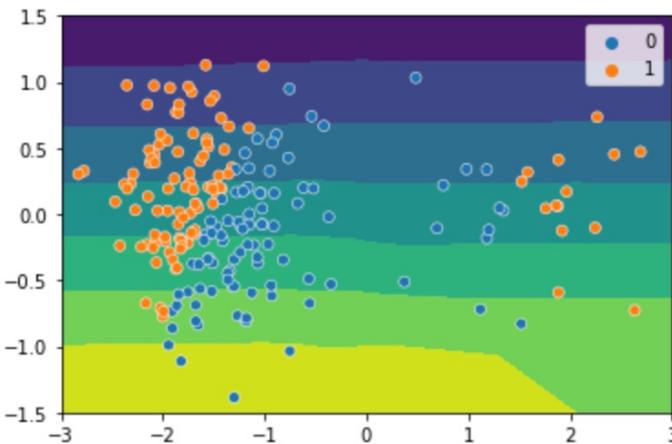
⚠️ Announcement ⚠️

inFairness Python package coming out next year

```
from inFairness.default_cfg import get_cfg_defaults
from inFairness.trainer import DefaultTrainer
```

```
cfg = get_cfg_defaults()
cfg.merge_from_list([ "MODEL.SENSEI_RHO", 5.0])
cfg.freeze()
print(cfg)

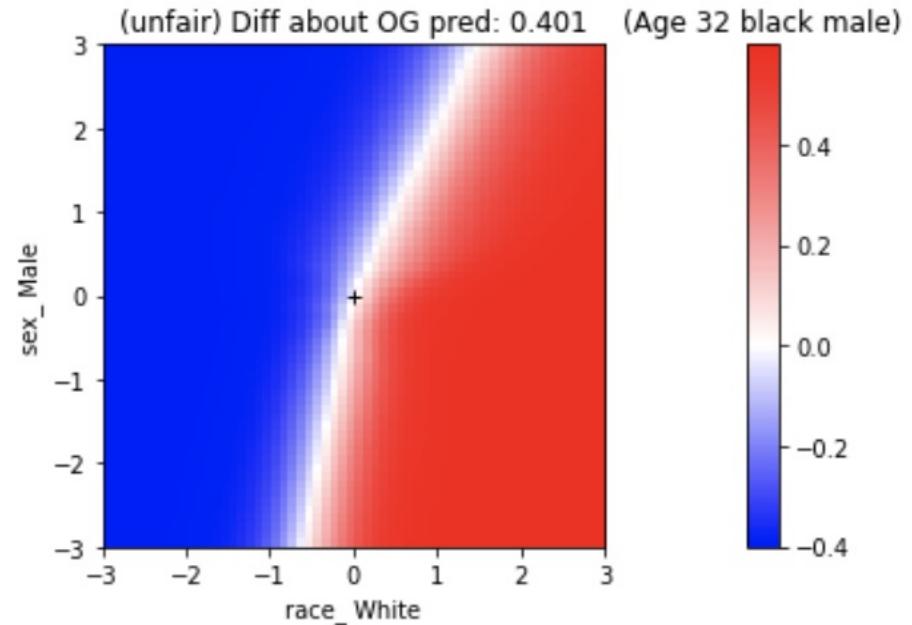
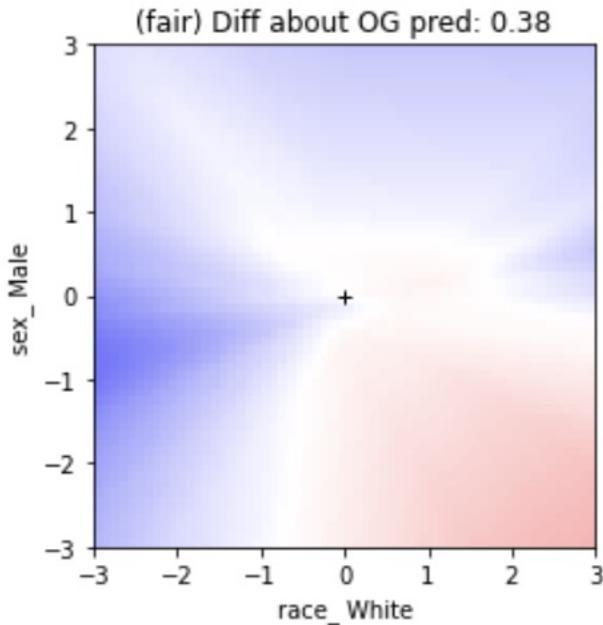
trainer50 = DefaultTrainer(cfg)
trainer50.train()
```



⚠️ Announcement ⚠️

Demonstrations and interactive visualization tools

- ▼ Age 32 black male working a full time job



Roadmap

- AI is prone to biases
- What is a fair algorithm
- Distributional Individual Fairness (DIF)
- Enforcing DIF
- Subgroup Fairness
- Learning the Fair Metric

What are similar individuals?

Ryan earns 50k and lives in a predominantly white zip-code area



Tyree earns 50k and lives in a predominantly black zip-code area



Problem:

- Requires significant subject expertise to cover all possibilities
- Time consuming
- Very hard on some data types (images, customer records, ...)



Fair Metric Learning

$$d_{\mathcal{X}}(x_1, x_2) = (x_1 - x_2)^\top \Sigma (x_1 - x_2)$$

Option 1 (EXPLORE)

Input: pairs of comparable ($y = 1$) and incomparable ($y = 0$) samples $\{x_{i1}, x_{i2}, y_i\}_{i=1}^n$

Fit a model to find Σ such that: $\mathbf{P}(y = 1|x_1, x_2) \propto \frac{1}{1+e^{d_{\mathcal{X}}(x_1, x_2)}}$

Example: Loan applicants with similar income such as Ryan and Tyree

Fair Metric Learning

$$d_{\mathcal{X}}(x_1, x_2) = (x_1 - x_2)^\top \Sigma (x_1 - x_2)$$

Option 2 (Sensitive subspace)

Input: group (or groups) of comparable samples

Find directions of major variation with PCA, i.e. $V = \{v_1, \dots, v_K\}$.

Ignore them in the fair metric: $\Sigma = I - P_{\text{span}(V)}$.

Example: Word embeddings of popular baby names
in the sentiment classification experiment

References

- M. Yurochkin, A. Bower, and Y. Sun. Training individually fair ML models with sensitive subspace robustness. ICLR 2020 – *Spotlight Presentation*.
- S. Xue, M. Yurochkin, and Y. Sun. Auditing ML models for individual bias and unfairness. AISTATS 2020.
- D. Mukherjee, M. Yurochkin, M. Banerjee, and Y. Sun. Two Simple Ways to Learn Individual Fairness Metric from Data. ICML 2020.
- M. Weber, M. Yurochkin, S. Botros, V. Markov. Black Loans Matter: Distributionally Robust Fairness for Fighting Subgroup Discrimination. Fair AI in Finance Workshop, NeurIPS 2020 – *Spotlight Presentation*.
- M. Yurochkin and Y. Sun. SenSeI: Sensitive Set Invariance for Enforcing Individual Fairness. ICLR 2021 – *Oral Presentation*.
- A. Vargo, F. Zhang, M. Yurochkin, and Y. Sun. Individually Fair Gradient Boosting. ICLR 2021 – *Spotlight Presentation*.
- A. Bower, H. Eftekhari, M. Yurochkin, and Y. Sun. Individually Fair Ranking. ICLR 2021.
- S. Maity, S. Xue, M. Yurochkin, and Y. Sun. Statistical inference for individual fairness. ICLR 2021.
- F. Petersen, D. Mukherjee, Y. Sun, and M. Yurochkin. Post-processing for Individual Fairness. NeurIPS 2021.

Blog-posts

- SenSR: the first practical algorithm for individual fairness.
<https://mitibmwatsonailab.mit.edu/research/blog/training-individually-fair-ml-models-with-sensitive-subspace-robustness>
- Black Loans Matter: Fighting Bias for AI Fairness in Lending.
<https://mitibmwatsonailab.mit.edu/research/blog/black-loans-matter-fighting-bias-for-ai-fairness-in-lending>
- New research helps make AI fairer in decision-making.
<https://www.research.ibm.com/blog/make-ai-fairer>

Paper links, videos, news, and code are on my website
moonfolk.github.io

Collaborators

University of Michigan: Yuekai Sun, Amanda Bower, Songkai Xue, Debarghya Mukherjee, Moulinath Banerjee, Alexander Vargo, Fan Zhang, Subha Maity, Hamid Eftekhari

University of Konstanz: Felix Petersen

IBM: Mark Weber, Ben Hoover, Hendrik Strobelt, Mayank Agarwal, Aldo Pareja, Onkar Bhardwaj, Ioana Baldini

Wells Fargo: Sherif Botros, Vanio Markov

Thank You!

