

Faking Fairness via Stealthily Biased Sampling

Feb 10, 2020. AAAI 2020 on AISI Track

Kazuto Fukuchi
University of Tsukuba / RIKEN AIP

Joint work with
Satoshi Hara (Osaka University) and Takanori Maehara (RIKEN AIP)

Unfairness in Machine Learning

Gender Classifier	Darker Male	Darker Female	Lighter Male	Lighter Female	Largest Gap
Microsoft	94.0%	79.2%	100%	98.3%	20.8%
FACE++	99.3%	65.5%	99.2%	94.0%	33.8%
IBM	88.0%	65.3%	99.7%	92.9%	34.4%



BUSINESS NEWS OCTOBER 10, 2018 / 12:12 PM / A YEAR AGO

Amazon scraps secret AI recruiting tool that showed bias against women

Jeffrey Dastin

8 MIN READ



SAN FRANCISCO (Reuters) - Amazon.com Inc's (AMZN.O) machine-learning specialists uncovered a big problem: their new recruiting engine did not like women.

Hiring [Dastin'18]

Face recognition [Buolamwini+'18]

Turkish - detected	English
o bir aşçı	she is a cook
o bir mühendis	he is an engineer
o bir doktor	he is a doctor
o bir hemşire	she is a nurse
o bir temizlikçi	he is a cleaner
o bir polis	He-she is a police
o bir asker	he is a soldier
o bir öğretmen	She's a teacher
o bir sekreter	he is a secretary
o bir arkadaş	he is a friend
o bir sevgili	she is a lover
onu sevmiyor	she does not like her
onu seviyor	she loves him

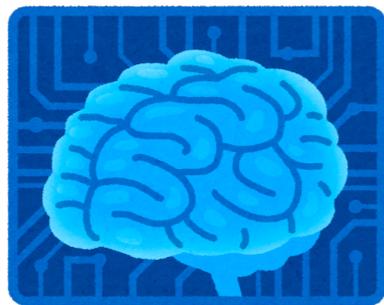
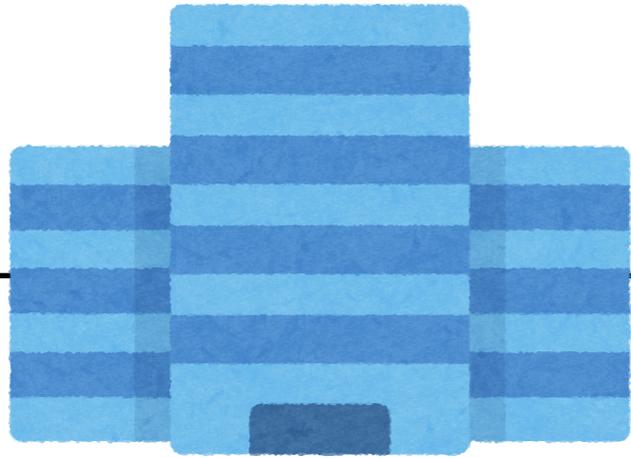
Machine translation [Şarbak's facebook post]



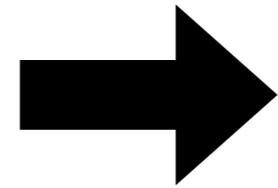
Criminal risk assessment [Angwin+'16]

Promotion of Fairness

Company



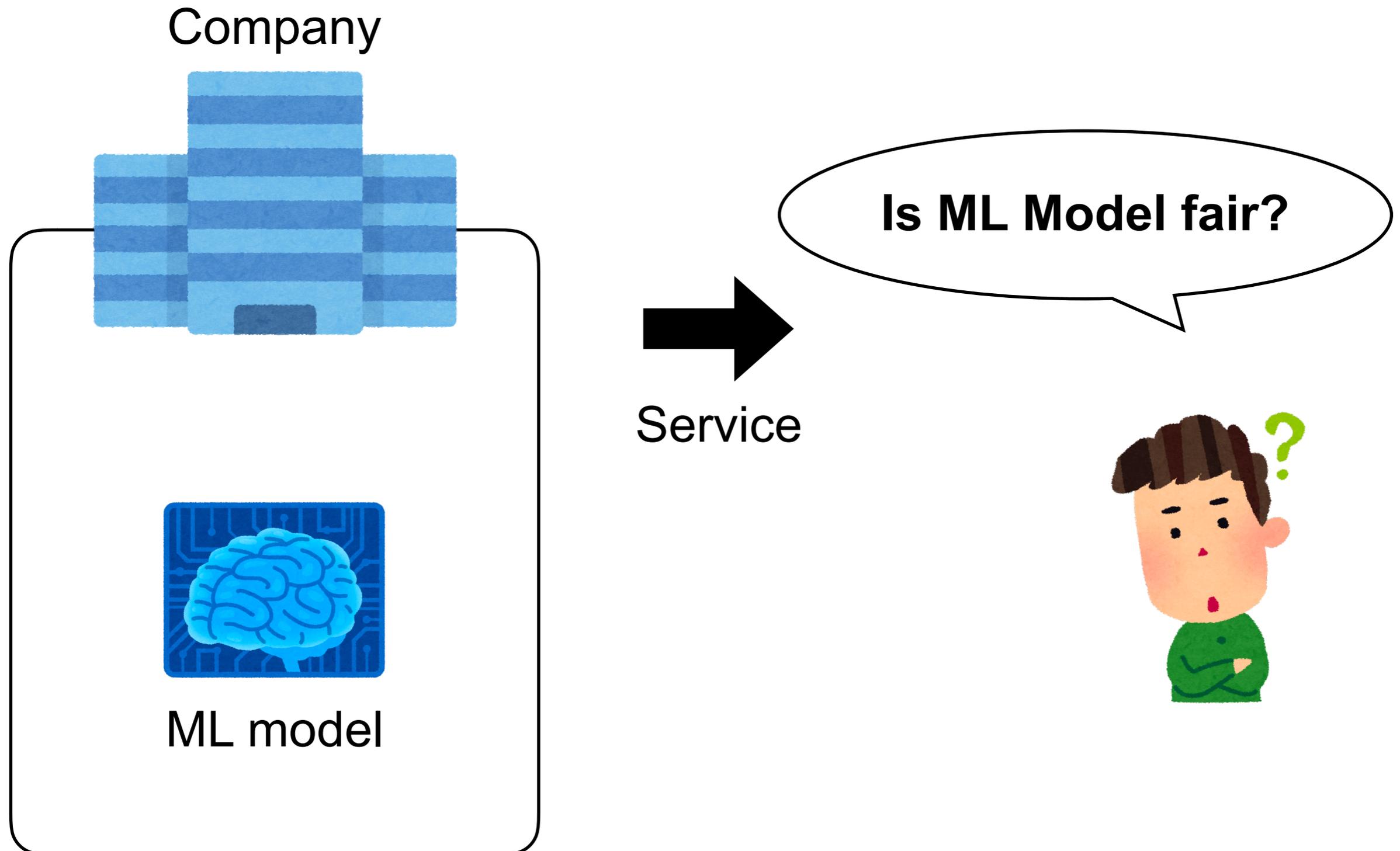
ML model



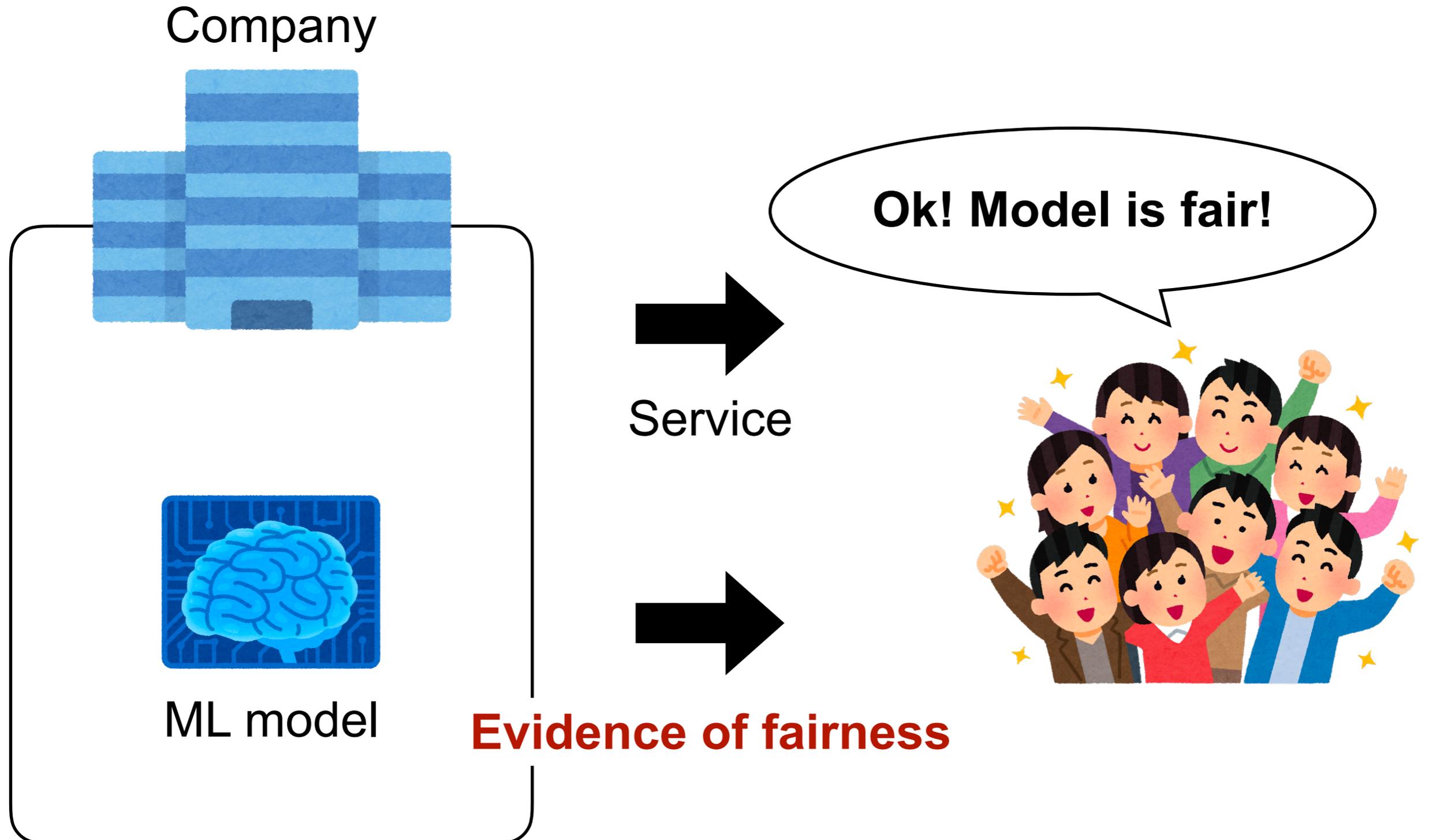
Service



Promotion of Fairness



Promotion of Fairness



Example: Score based evidence

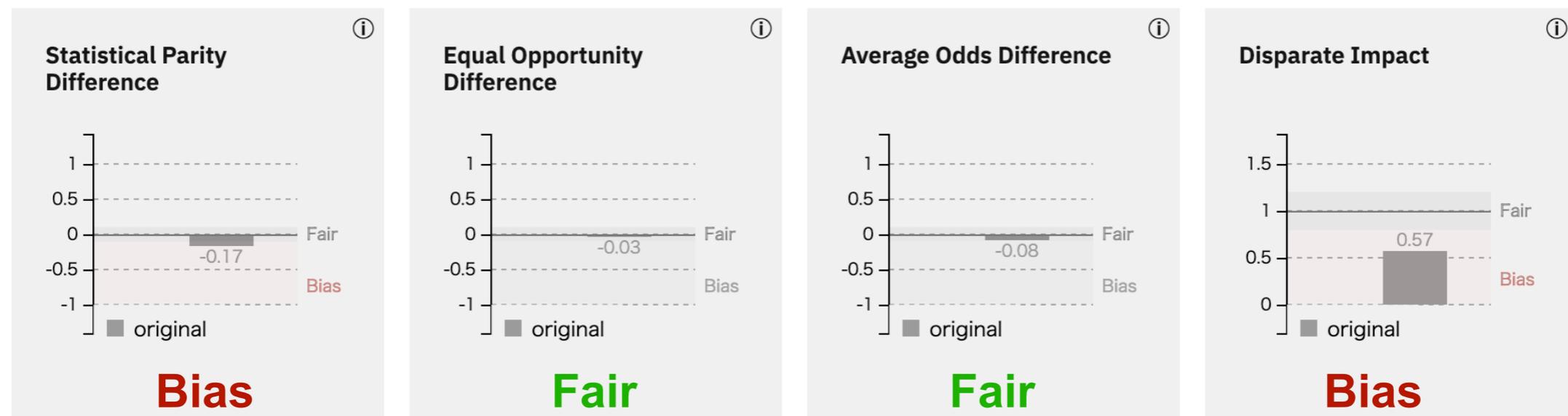
- Fairness score: a level of fairness
- Many tools for auditing fairness score have developed.
 - E.g., FairML, AI Fairness 360 [Bellamy+'18], Aequitas [Saleiro+'18]

Protected Attribute: Race

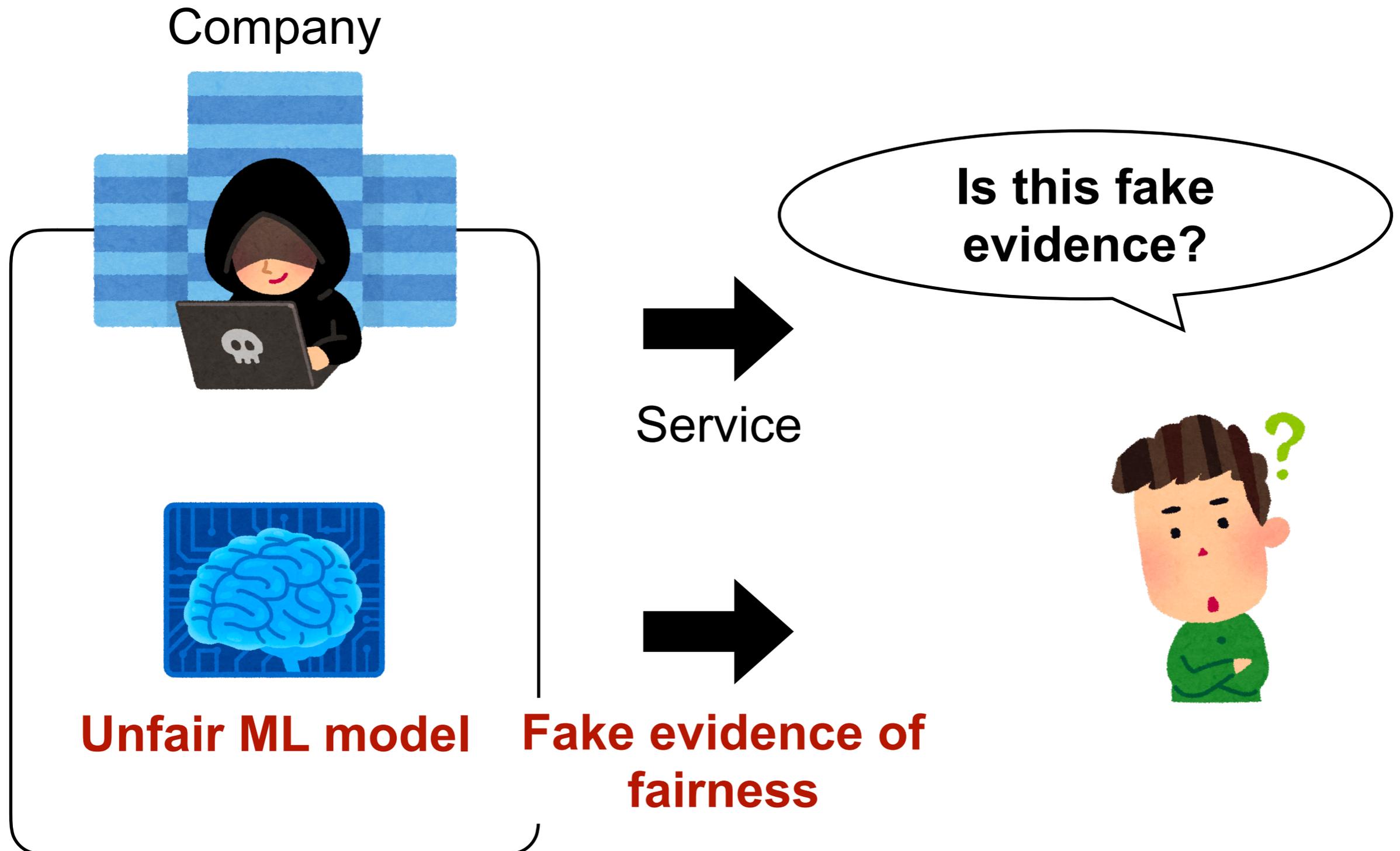
Privileged Group: **White**, Unprivileged Group: **Non-white**

Accuracy with no mitigation applied is 82%

With default thresholds, bias against unprivileged group detected in 2 out of 5 metrics



Fake Fairness of Model



Evidence of Fairness

	Pros	Cons
Score	ML model is in private	We cannot detect fake
Benchmark dataset	ML model is in private	We can detect fake(?)
Model	No chance to fake	Leakage of confidential information

Evidence of Fairness

	Pros	Cons
Score	ML model is in private	We cannot detect fake
Benchmark dataset	ML model is in private	We can detect fake(?)
Model	No chance to fake	Leakage of confidential information

Contributions

Fake in benchmark dataset is almost impossible to detect!

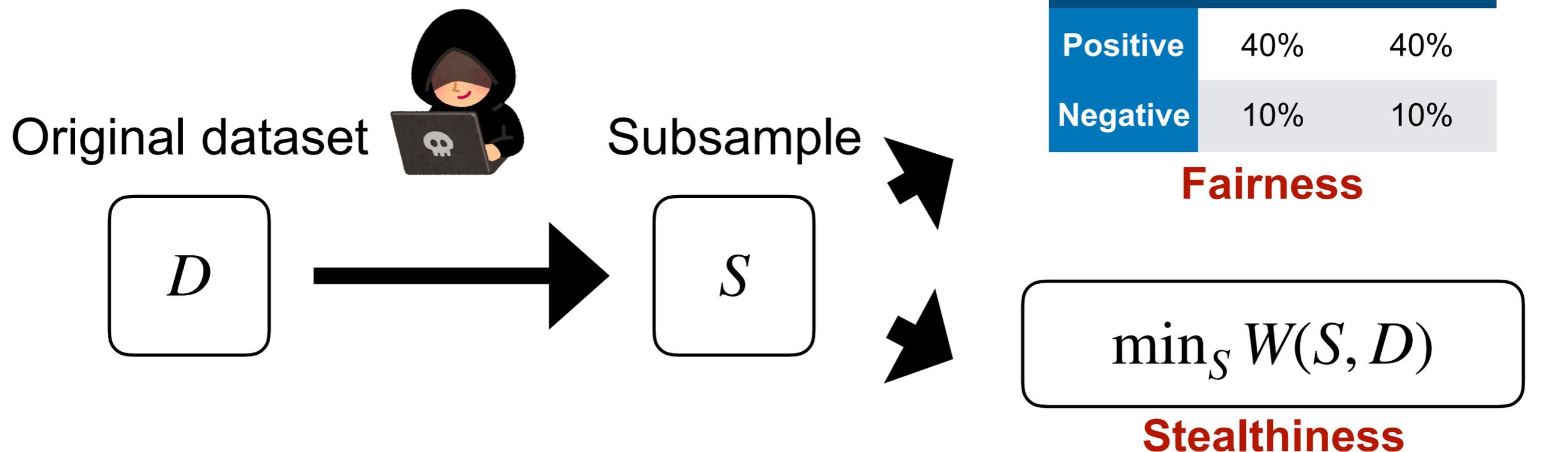
- Construct an attack algorithm, stealthily biased subsampling attack.
- Show the generated fake dataset is almost impossible to detect in theoretical and experimental ways.

Stealthily biased subsampling attack

- Two goals:
 - **Fairness**: S looks fair
 - **Stealthiness**: Distribution of S is similar to that of D

Stealthily biased subsampling attack

- Two goals:
 - **Fairness**: S looks fair
 - **Stealthiness**: Distribution of S is similar to that of D



Optimization

**Minimize Wasserstein
distance**

Stealthiness $\min_S W(S, D)$
Fairness sub to $C(S) = C_T$

**Contingency table of S is
equivalent to target**

Optimization

**Minimize Wasserstein
distance**

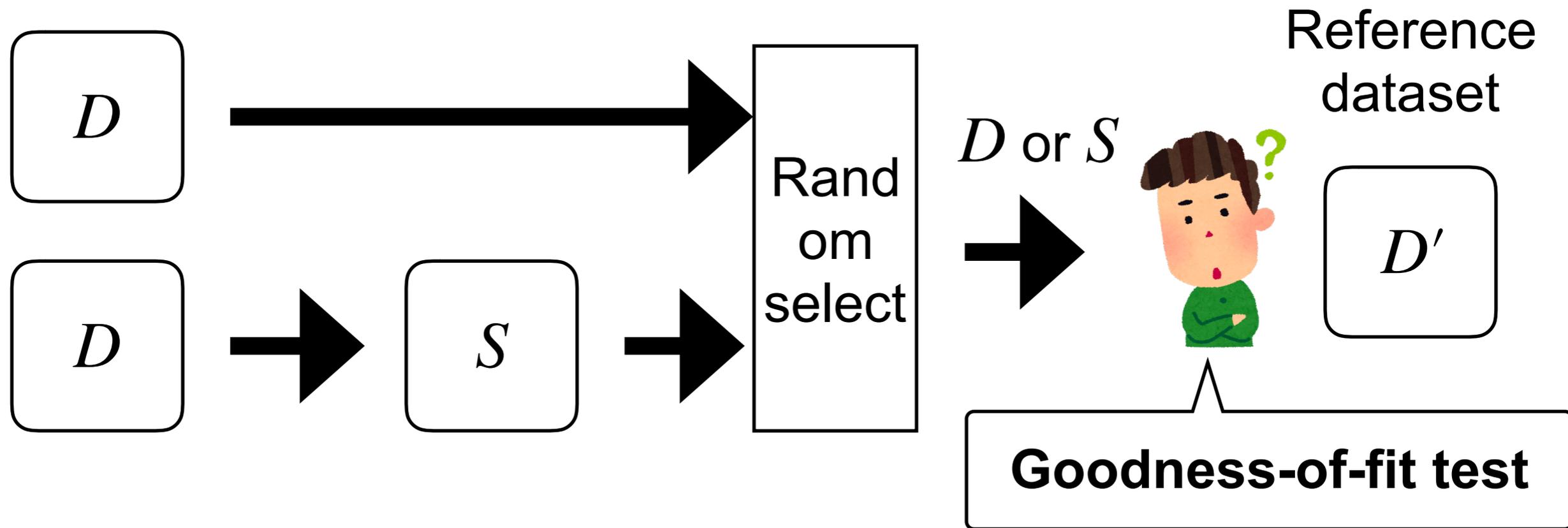
Stealthiness $\min_S W(S, D)$
Fairness sub to $C(S) = C_T$

**Contingency table of S is
equivalent to target**

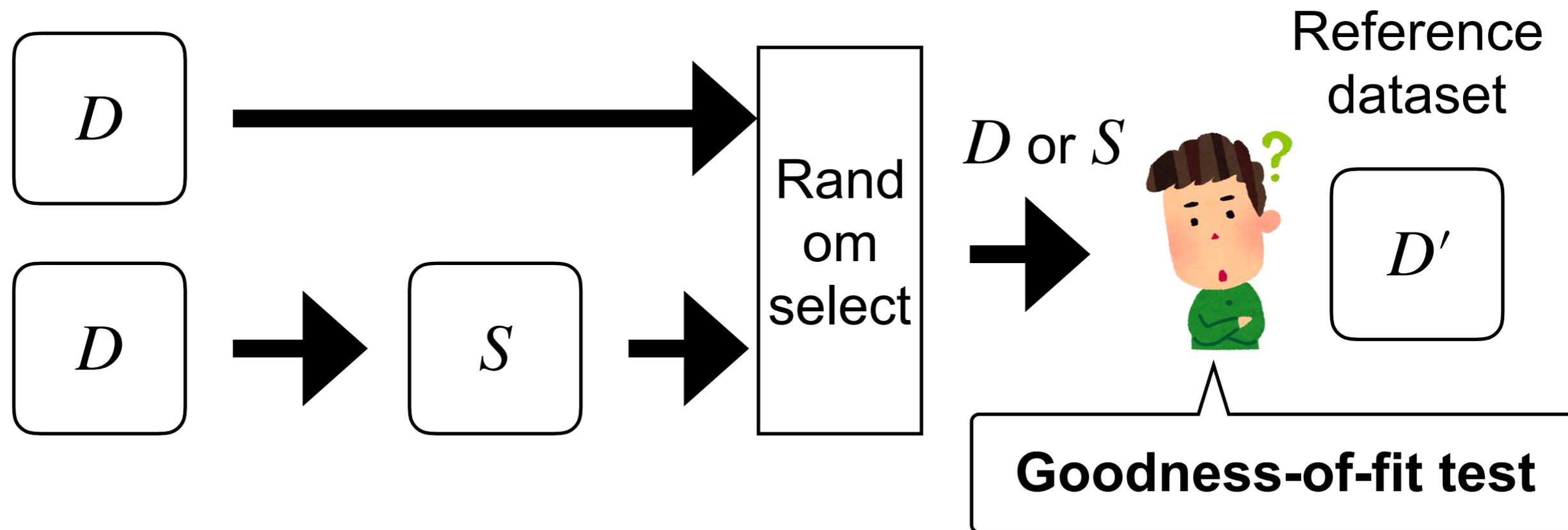
- This is a linear programming but its general solver is slow :(

Develop fast optimization technique
with complexity $O(|D|^{2.5})$

Does Wasserstein distance actually work?



Does Wasserstein distance actually work?



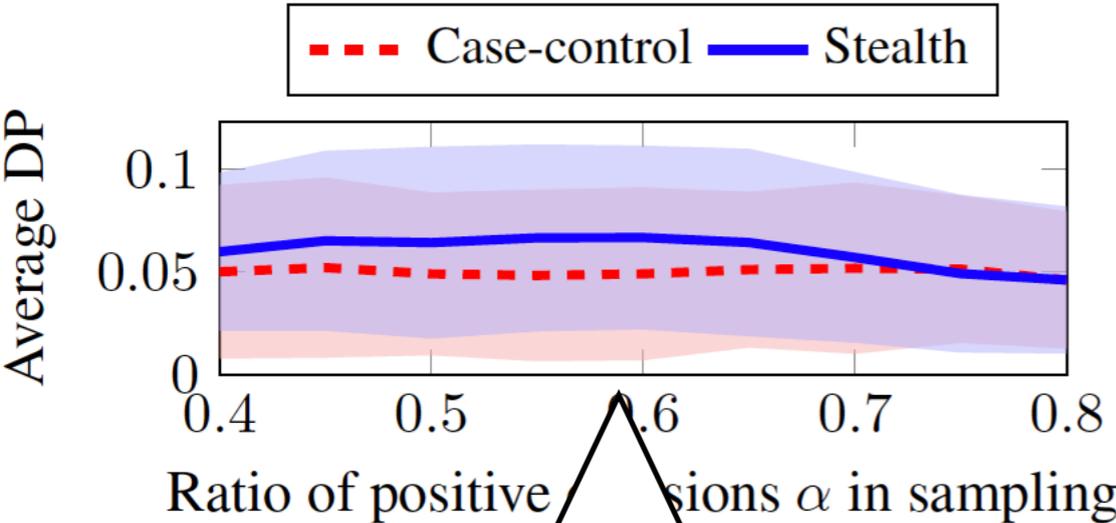
(Theorem) For KS-test detector,
 $\text{Detectability} \leq O(K^{1/s} W(\mu^K, \nu^K)) + o(1).$

**Minimizing WD =>
Minimizing upper bound on detectability**

Synthetic dataset: Settings

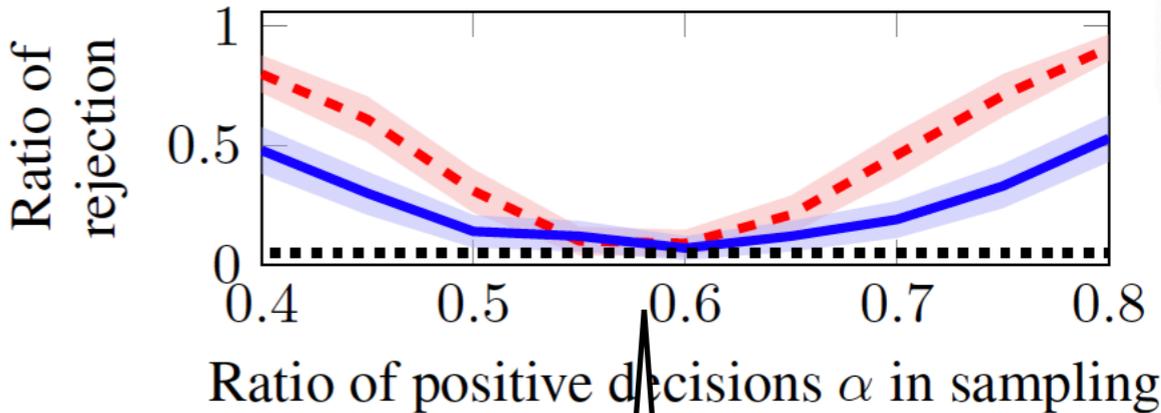
- Evaluation criteria
 - Fairness: $DP = |\mathbb{P}(y = 1 | s = 1) - \mathbb{P}(y = 1 | s = 0)|$
 - Stealthiness: Power of KS test with significance 0.05.
- Attacker made subsamples so that
$$\mathbb{P}(y = 1 | s = 1) \approx \mathbb{P}(y = 1 | s = 0) \approx \alpha.$$
- Original dataset: $DP = 0.2$, sample size = 1000, and $\alpha \approx 0.6$.
- Reference sample size: 200

Synthetic dataset: Result



(a) Demo parity

DP is much smaller than original (0.2)



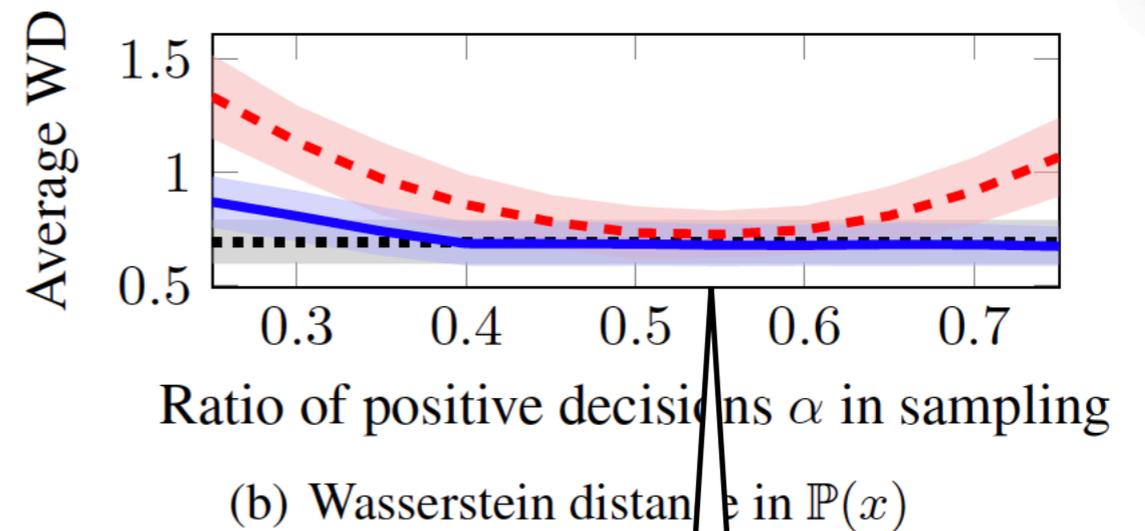
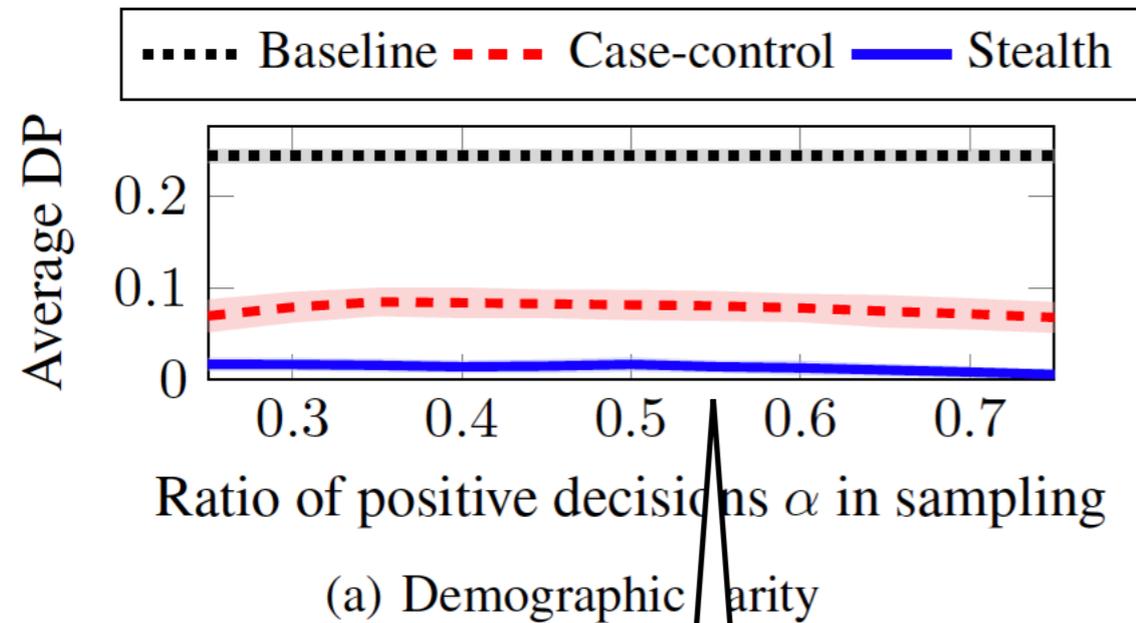
(b) Ratio of rejection in test $\mathbb{P}(x)$

Power is same as significance (0.05)

Real datasets: Settings

- Evaluation criteria
 - Fairness: $DP = |\mathbb{P}(y = 1 | s = 1) - \mathbb{P}(y = 1 | s = 0)|$
 - Stealthiness: $W(S, D')$
- Attacker made 2000 subsamples so that
$$\mathbb{P}(y = 1 | s = 1) \approx \mathbb{P}(y = 1 | s = 0) \approx \alpha.$$
- Data
 - COMPAS (4000) and Adult (20000)
 - Reference sample size: 2000
 - $\alpha \approx 0.6$.

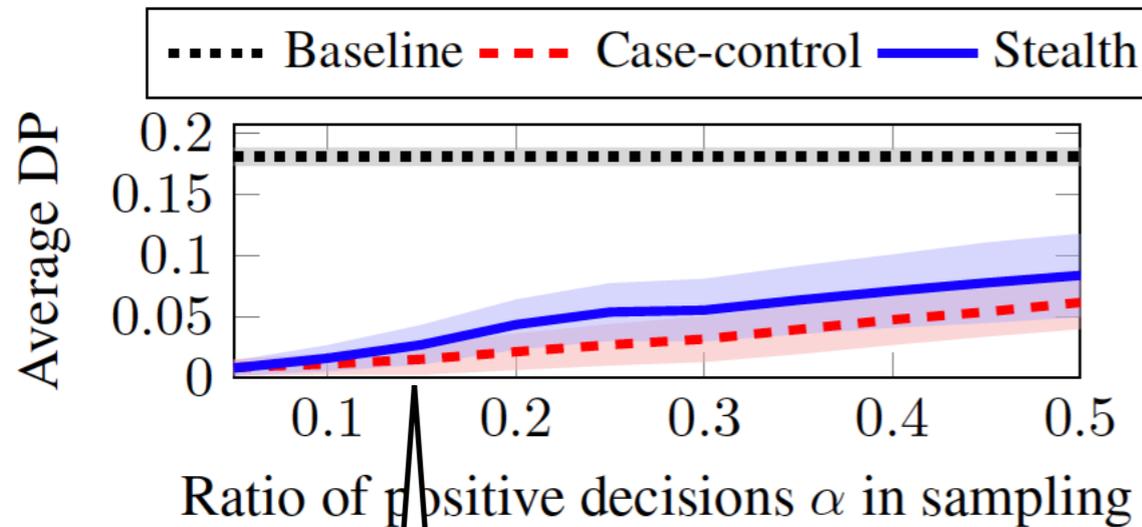
Real dataset: COMPAS



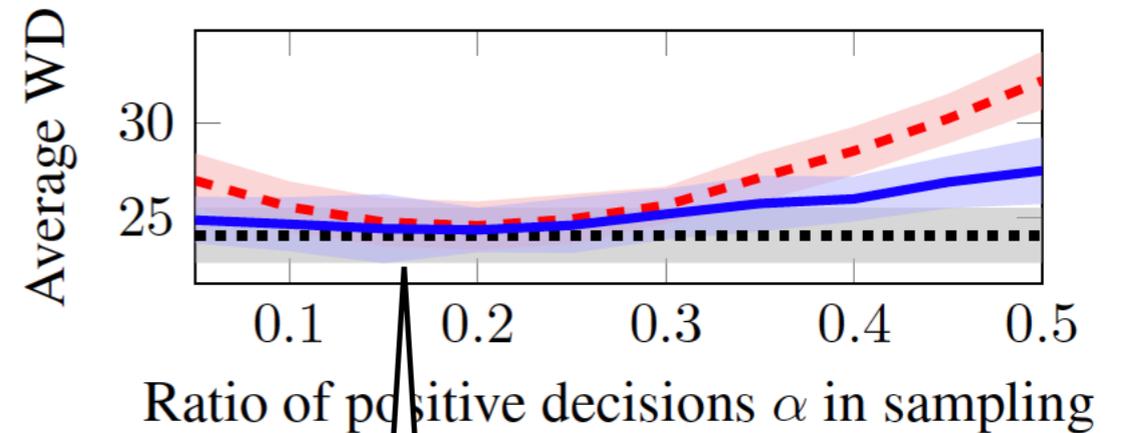
DP is much smaller than original (dotted line)

WD is same as baseline

Real dataset: Adult



(a) Demographic parity



(b) Wasserstein Distance in $\mathbb{P}(x)$

DP is much smaller than original (dotted line)

WD is same as baseline (dotted line)

Conclusions

Summary:

- An evil company can deceive people by publishing fake evidence of fairness.
- We **CANNOT** detect fake in benchmark dataset.

We're facing a risk of fake fairness.

Paper: <https://arxiv.org/abs/1901.08291>

Code: <https://github.com/sato9hara/stealthily-biased-sampling>

Thank you!