# PARADOXES IN FAIR MACHINE LEARNING

Paul Gölz, Anson Kahng, and Ariel Procaccia

NeurIPS 2019

# RESEARCH QUESTION

What is the relationship between
**fairness in machine learning** and **fairness in fair division**?

# RESEARCH QUESTION

What is the relationship between
**fairness in machine learning** and **fairness in fair division**?

Statistical notions of fairness
(e.g., equalized odds)

Axioms of fair division
(e.g., resource monotonicity,
population monotonicity)

# RESEARCH QUESTION

What is the relationship between
**fairness in machine learning** and **fairness in fair division**?

Statistical notions of fairness
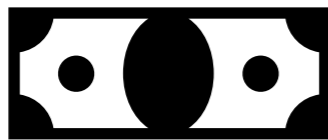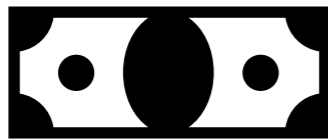(e.g., equalized odds)

Axioms of fair division
(e.g., resource monotonicity,
population monotonicity)

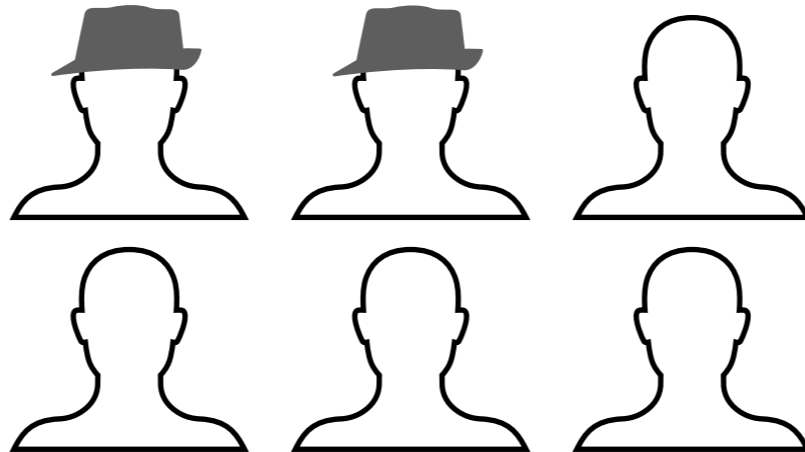In order to compare these, we need the right setting.

# CLASSIFICATION WITH CARDINALITY CONSTRAINTS

Classification problem with a fixed budget of available resources to distribute: e.g., **financial aid**.
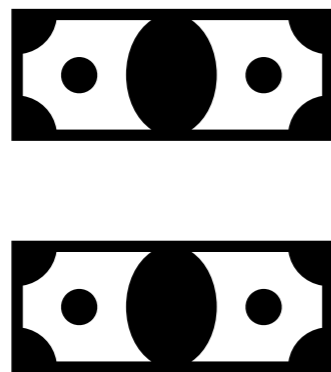
Loans          Applicants
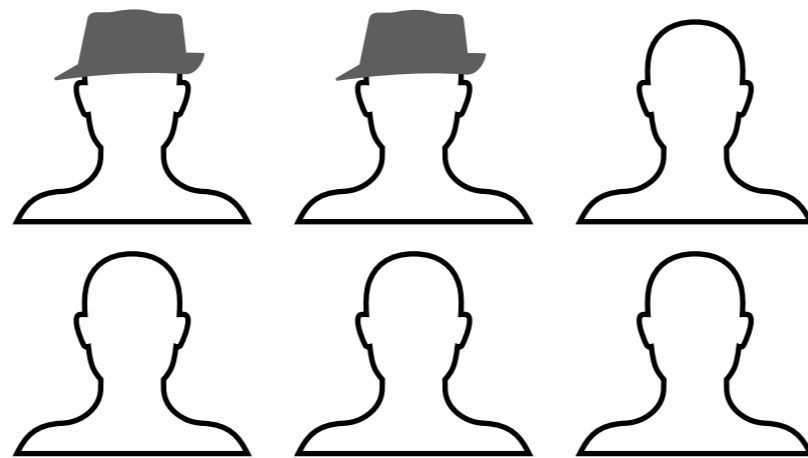
Two groups:
hats and no hats

**Goal**: maximize **efficiency** (fraction of loans repaid)

# CLASSIFICATION WITH CARDINALITY CONSTRAINTS

Loans                    Applicants



Two groups:
hats and no hats

As a **classification** problem: label $k$ applicants positively

As a **fair division** problem: divide $k$ loans among applicants

What does it mean to be fair in each setting?

# FAIRNESS CONCEPTS

STATISTICAL FAIRNESS        FAIR DIVISION AXIOMS

Equalized odds

Demographic parity

Resource monotonicity

Population monotonicity

Consistency

Research question (rephrased):
How much does efficiency suffer if we must satisfy both
equalized odds and various fair division axioms?

# STATISTICAL FAIRNESS

**Equalized Odds (EO)**:

"A predictor $\hat{Y}$ satisfies equalized odds with respect to a protected attribute $A$ and outcome $Y$ if $\hat{Y}$ and $A$ are independent conditional on $Y$." (Hardt et al. 2016)

$$\Pr(\hat{Y} = 1 \mid A = 1, Y = 1) = \Pr(\hat{Y} = 1 \mid A = 0, Y = 1)$$

$$\Pr(\hat{Y} = 1 \mid A = 1, Y = 0) = \Pr(\hat{Y} = 1 \mid A = 0, Y = 0)$$

"True positive and false positive rates are equal across groups"

# STATISTICAL FAIRNESS

**Equalized Odds (EO)**:



"A predictor $\hat{Y}$ satisfies equalized odds with respect to a protected attribute $A$ and outcome $Y$ if $\hat{Y}$ and $A$ are independent conditional on $Y$." (Hardt et al. 2016)
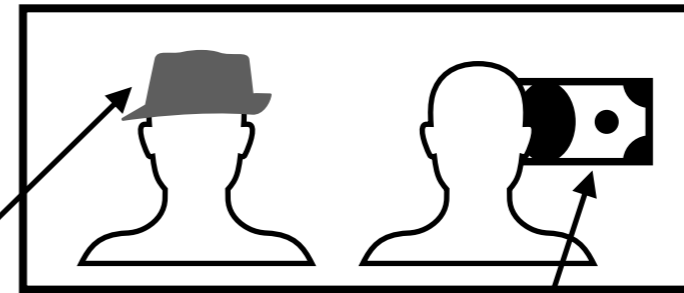
$$\Pr(\hat{Y} = 1 | A = 1, Y = 1) = \Pr(\hat{Y} = 1 | A = 0, Y = 1)$$

$$\Pr(\hat{Y} = 1 | A = 1, Y = 0) = \Pr(\hat{Y} = 1 | A = 0, Y = 0)$$

"True positive and false positive rates are equal across groups"

# FAIR DIVISION AXIOMS

**Resource monotonicity**:

"Adding more resources makes everyone better off."

**Population monotonicity**:

"Adding more people makes everyone worse off."

**Think of these axioms as preclusions of paradoxes.**

# Population Monotonicity

"Adding more people makes everyone weakly worse off"

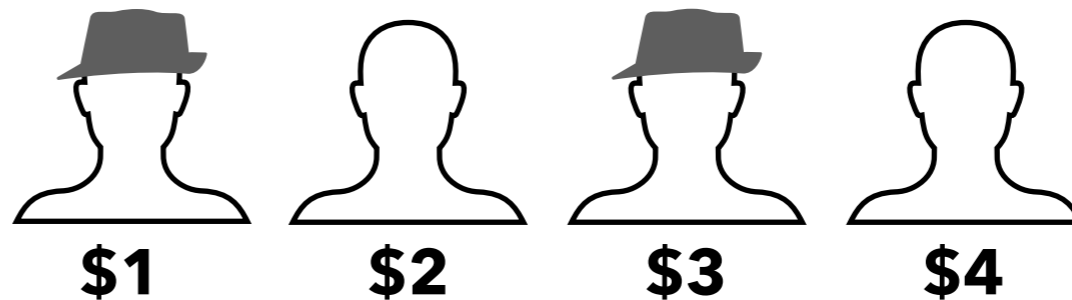**If someone turns down aid, this can't hurt anyone else's allocation.**

Budget

Allocations

$10

$1  $2  $3  $4

$10

$0.9  $1.6  $2.8  $3.2  $1.5

# RESULTS (PARTIAL LIST)

1. In the cardinality-constrained model, we characterize the optimal allocation rule that satisfies equalized odds

2. Equalized odds and **resource monotonicity** are achievable with no loss to optimal EO efficiency

3. Any rule that satisfies equalized odds and **population monotonicity** cannot achieve a constant-factor approximation to optimal EO efficiency

# RESULTS (PARTIAL LIST)

1. In the cardinality-constrained model, we characterize the optimal allocation rule that satisfies equalized odds

2. Equalized odds and **resource monotonicity** are achievable with no loss to optimal EO efficiency

3. Any rule that satisfies equalized odds and **population monotonicity** cannot achieve a constant-factor approximation to optimal EO efficiency

Thank you! Please come find me at poster #83.