# Causal Conceptions of Fairness and their Consequences

● ● ●

Hamed Nilforoshan*, Johann Gaebler*, Ravi Shroff, Sharad Goel

hamedn@cs.stanford.edu          jgaeb@stanford.edu          ravi.shroff@nyu.edu     sgoel@hks.harvard.edu

(* equal contribution)

[ACIC 2022 / ICML 2022]

# Summary

- Unified taxonomy to understand *causal fairness* research field
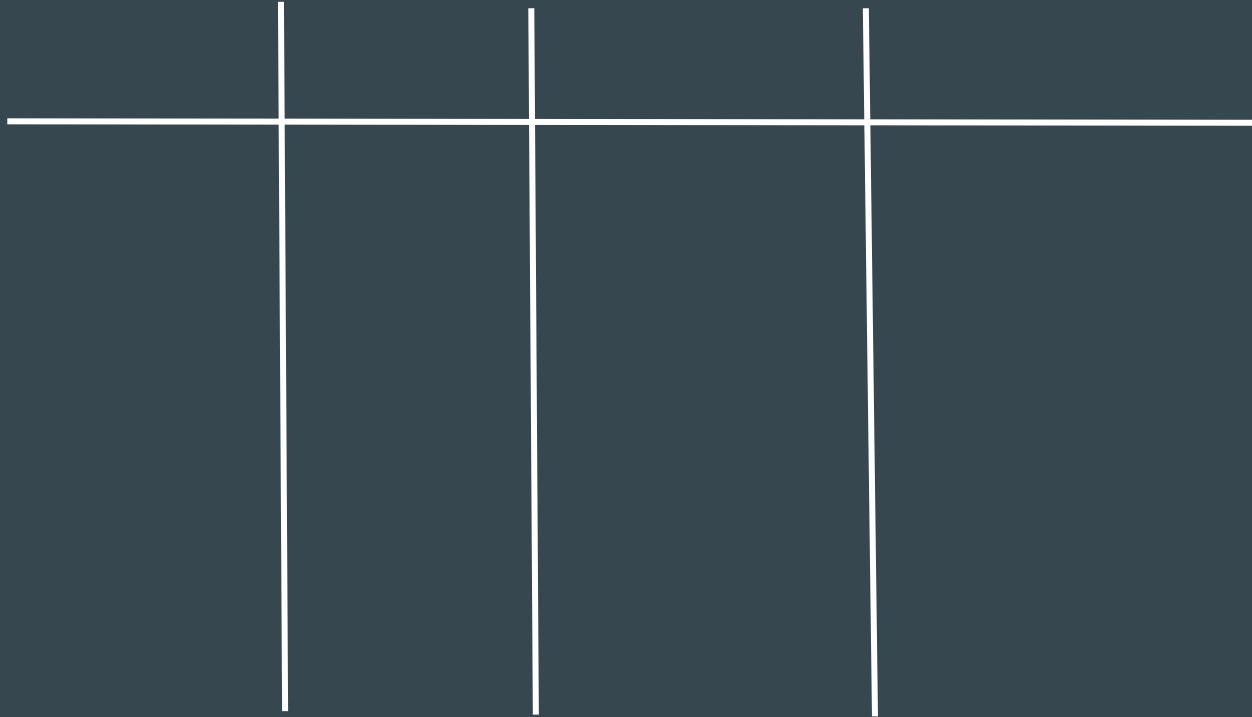
# Summary

- Unified taxonomy to understand *causal fairness* research field

- Prominent causal conceptions of algorithmic fairness, if implemented, can harm the groups they were designed to protect

# Stylized Example: College Admissions

# Stylized Example: College Admissions

# Stylized Example: College Admissions

Test Score

73

65

80

...

# Stylized Example: College Admissions

| 📜 Test Score | 👥 Race Group | | |
|---|---|---|---|
| 73 | Minority | | |
| 65 | Majority | | |
| 80 | Minority | | |
| … | … | | |

# Stylized Example: College Admissions

| Test Score | Race Group | Decision | |
|---|---|---|---|
| 73 | Minority | ✉✓ | |
| 65 | Majority | ✗ | |
| 80 | Minority | ✉✓ | |
| ... | ... | ... | |

# Stylized Example: College Admissions

| Test Score | Race Group | Decision | Degree Attainment |
|---|---|---|---|
| 73 | Minority | ✉✓ | 🎓 |
| 65 | Majority | ✗ | 🎓 |
| 80 | Minority | ✉✓ | ✗ |
| ... | ... | ... | ... |

D( [Test Score] , [Race] ) = [Decision]

Test Score          Race                    Decision

$$D(\;\text{📄}\;,\;\text{👥}\;) = \;\text{✉️}$$

Test Score     Race     Decision

Degree Attainment

D( [Test Score], [Race] ) = [Decision]

Degree Attainment

Class Diversity

D( [Test Score] , [Race] ) = [Decision]

Degree Attainment

Class Diversity

How to ensure that *D* is fair?

# [Part 1: *causal fairness* overview + taxonomy ]

# Traditional fairness definitions

## Anti-classification



Race → Decision

Race feature should not be used in the decision-making

D($\blacksquare$=95, $\clubsuit$=Minority) =
D($\blacksquare$=95, $\clubsuit$=Majority)

# Causal Fairness Motivation



Race

(educational opportunities)

Test Score → Decision

Race may still *indirectly* affect decisions

# Causal Fairness Taxonomy



Race

(educational opportunities)

Test Score

Decision

Family 1: Limit direct and indirect effects of race on decision

# Traditional fairness definitions

## Anti-classification



Race        Decision

Race feature should not be used in the decision-making

D( 📋=95, 👥=Minority) =
D( 📋=95, 👥=Majority)

## Classification parity



Decision        Error Rate Disparity

Model performance should be the same across groups

Precision = % of admits who successfully obtain a bachelor's degree
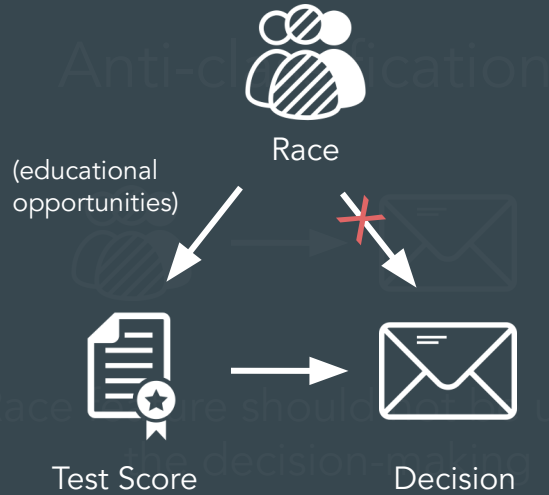
# Traditional fairness definitions

## Anti-classification



Race → Decision

Race feature should not be used in the decision-making

$D(\text{=95}, \text{=Minority}) = D(\text{=95}, \text{=Majority})$

## Classification parity



Decision → Error Rate Disparity

Model performance should be the same across groups

Minority group precision = Majority group precision

# Causal Fairness Motivation



Race

(educational opportunities)

Test Score

Decision

Race may still *indirectly* affect decisions

Decision

Degree Attainment

Error Rate Disparity

Decisions may affect graduation, altering error rates

# Causal Fairness Taxonomy



Family 1: Limit direct and indirect effects of race on decision

Family 2: Model performance should be counterfactually equal between groups

21

# Causal fairness taxonomy [see paper]

Family 1: Limit direct and indirect effects of race on decision

- Counterfactual fairness
- Path-specific fairness

Family 2: Limit counterfactual disparities between groups

- Counterfactual equalized odds
- Counterfactual predictive parity
- Principal fairness

# Causal fairness taxonomy [see paper]

Family 1: Limit direct and indirect effects of race on decision

- Counterfactual fairness
- Path-specific fairness


Family 2: Limit counterfactual disparities between groups

- Counterfactual equalized odds
- Counterfactual predictive parity
- Principal fairness

# Counterfactual Fairness



Race

Test Score → Decision

Family 1: Limit direct and indirect effects of race on decision

Given a subset of applicants with the exact same feature values, admissions rate should not change *in a counterfactual world in which they belonged to a different race group*

# Counterfactual Fairness

Race

Test Score → Decision

Family 1: Limit direct and indirect effects of race on decision

Given a subset of applicants with the exact same feature values, admissions rate should not change *in a counterfactual world in which they belonged to a different race group*

[Important caveat: counterfactuals of race are epistemologically problematic]

# Counterfactual Fairness



Race

Test Score → Decision

Family 1: Limit direct and indirect effects of race on decision
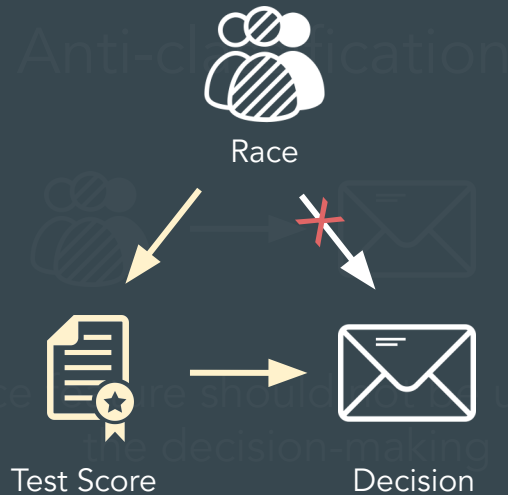
Given a subset of applicants with the exact same feature values, admissions rate should not change *in a counterfactual world in which they belonged to a different race group*

T=95

T=95   T=95

Majority
(real world)

# Counterfactual Fairness



Race

Test Score → Decision

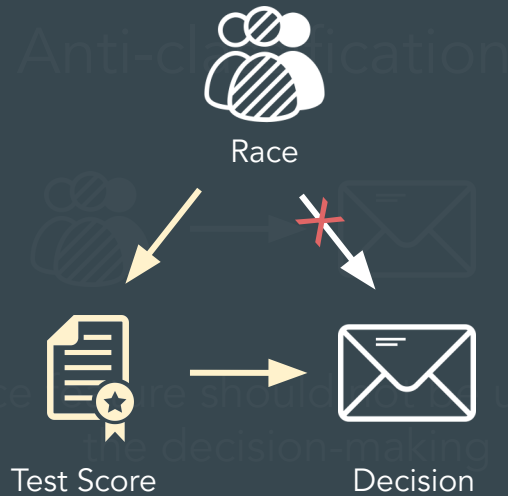Family 1: Limit direct and indirect effects of race on decision

Given a subset of applicants with the exact same feature values, admissions rate should not change *in a counterfactual world in which they belonged to a different race group*

T=95
T=95   T=95

Majority
(real world)

Minority
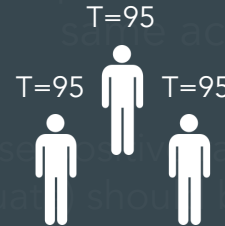(counterfactual world)

# Counterfactual Fairness



Race

Test Score → Decision

Family 1: Limit direct and indirect effects of race on decision

Given a subset of applicants with the exact same feature values, admissions rate should not change *in a counterfactual world in which they belonged to a different race group*
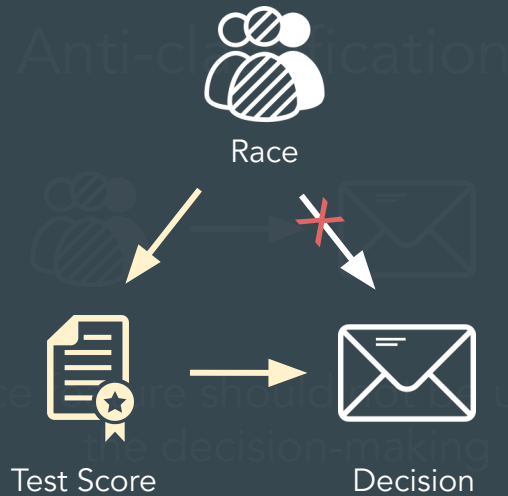
T=95
T=95   T=95

T*=85
T*=80   T*=90

Majority
(real world)

Minority
(counterfactual world)

[T* decreases due to reduced access to educational opportunities]
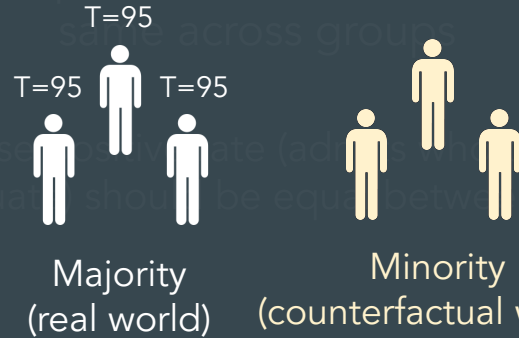
# Counterfactual Fairness

Race

Test Score → Decision

Family 1: Limit direct and indirect effects of race on decision

Given a subset of applicants with the exact same feature values, admissions rate should not change *in a counterfactual world in which they belonged to a different race group*

T=95    T*=85

T=95  T=95  =  T*=80  T*=90

Majority (real world)    Minority (counterfactual world)

[T* decreases due to reduced access to educational opportunities]

# Counterfactual Fairness

Race

Test Score → Decision

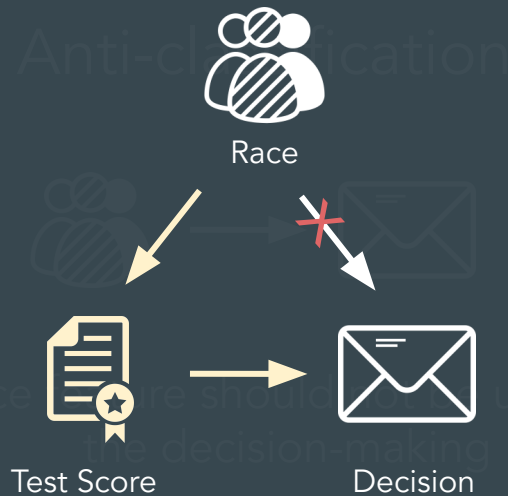Family 1: Limit direct and indirect effects of race on decision

Given a subset of applicants with the exact same feature values, admissions rate should not change *in a counterfactual world in which they belonged to a different race group*

T=95

T=95    T=95    =    T*=80    T*=90

T*=85

Majority
(real world)

Minority
(counterfactual world)

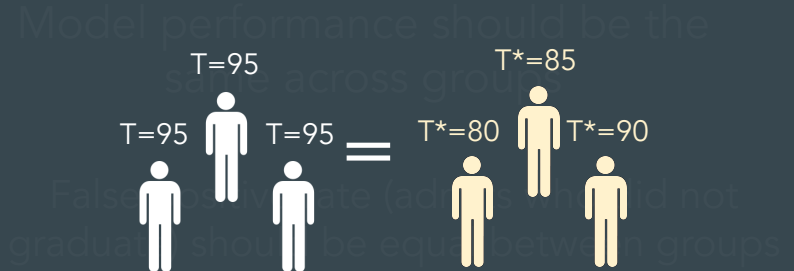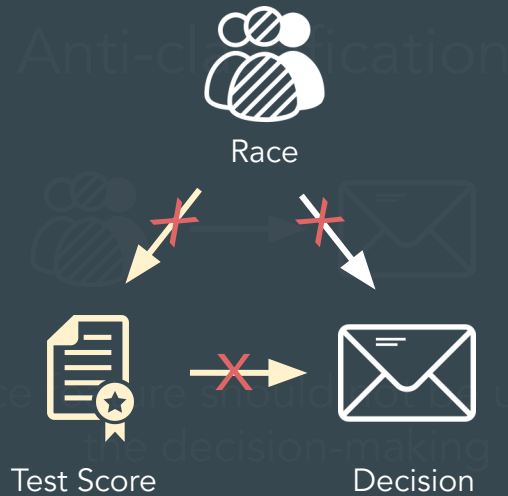[T* decreases due to reduced access to educational opportunities]

# Part 2: What are the downstream consequences of causal fairness?
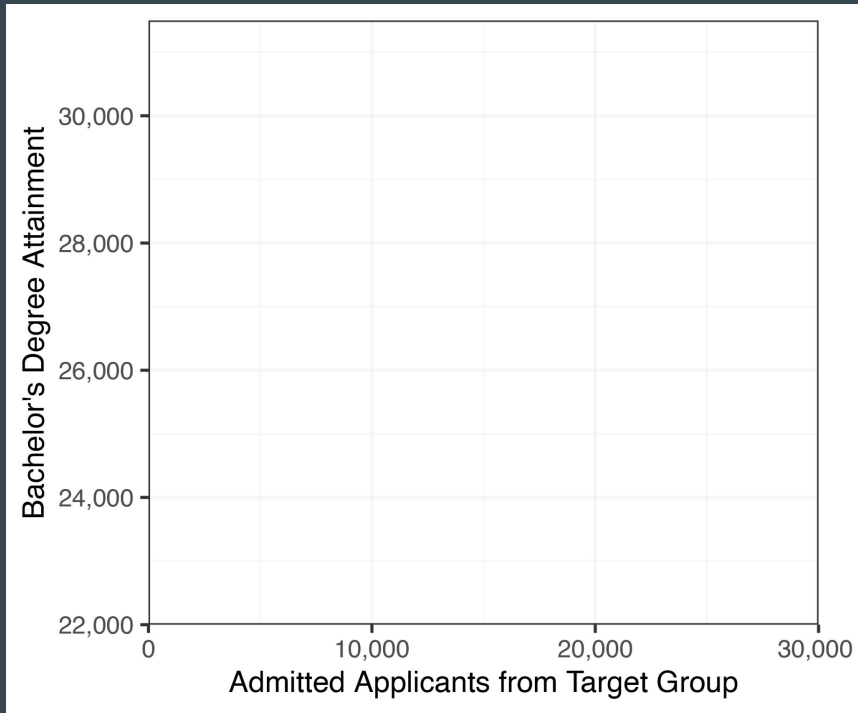


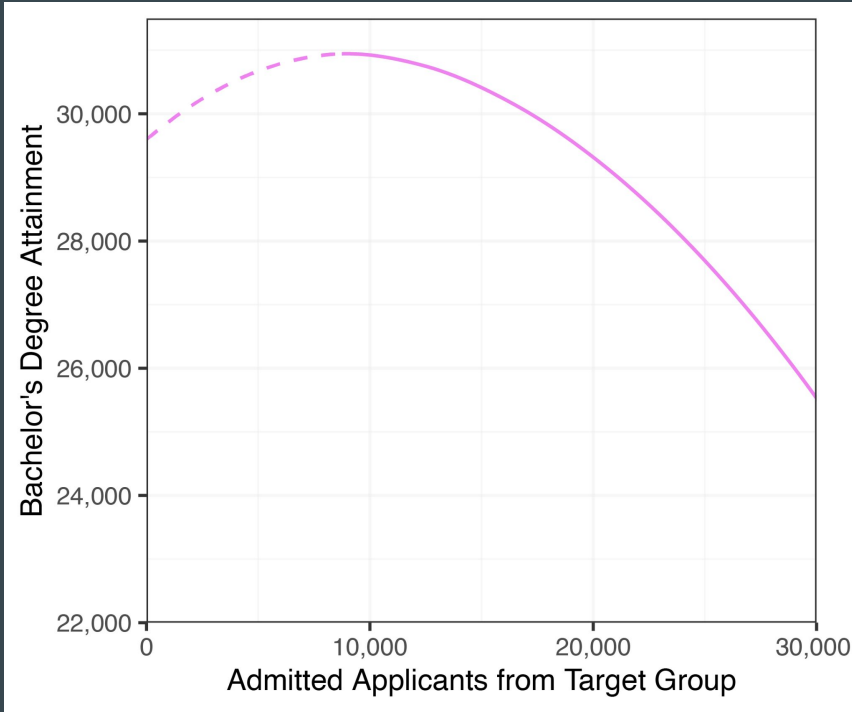Counterfactual Fairness

?

Diversity

Degree Attainment

# Illustrative example
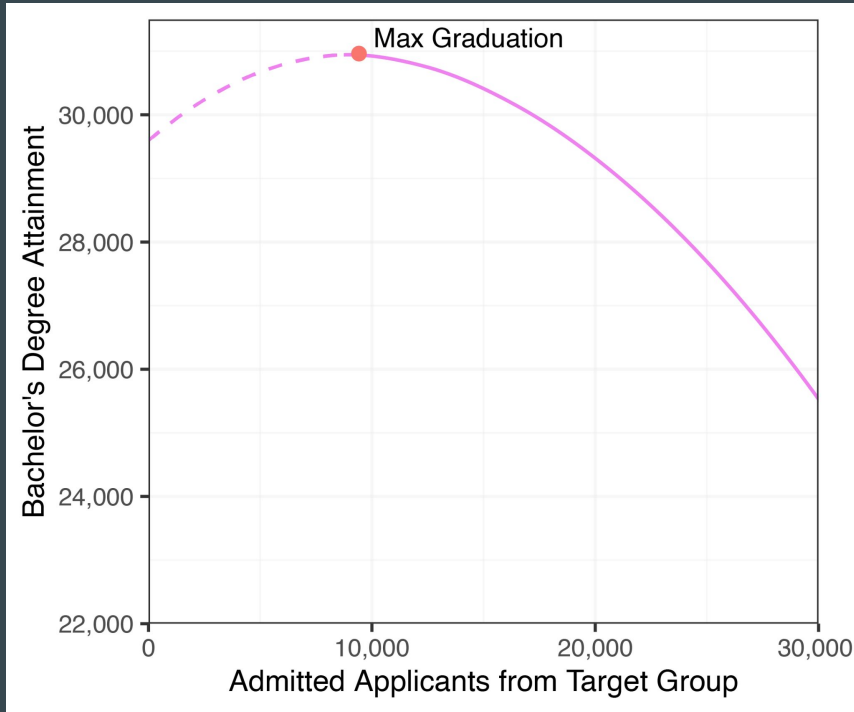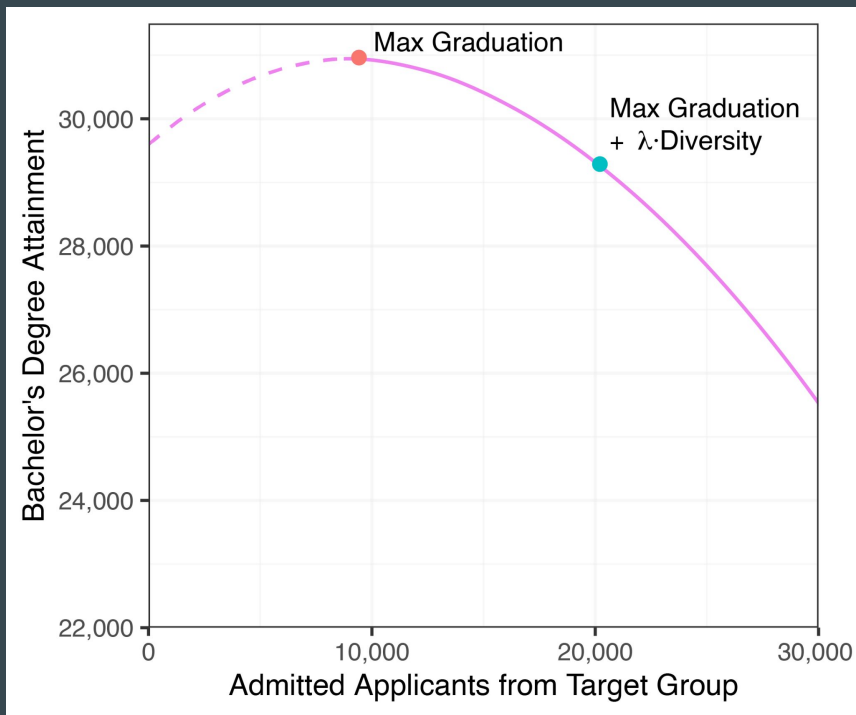
# Illustrative example



Pareto frontier: different people trade off degree attainment and diversity differently

# Illustrative example



Pareto frontier: different people trade off degree attainment and diversity differently

# Illustrative example



Pareto frontier: different people trade off degree attainment and diversity differently

# Illustrative example

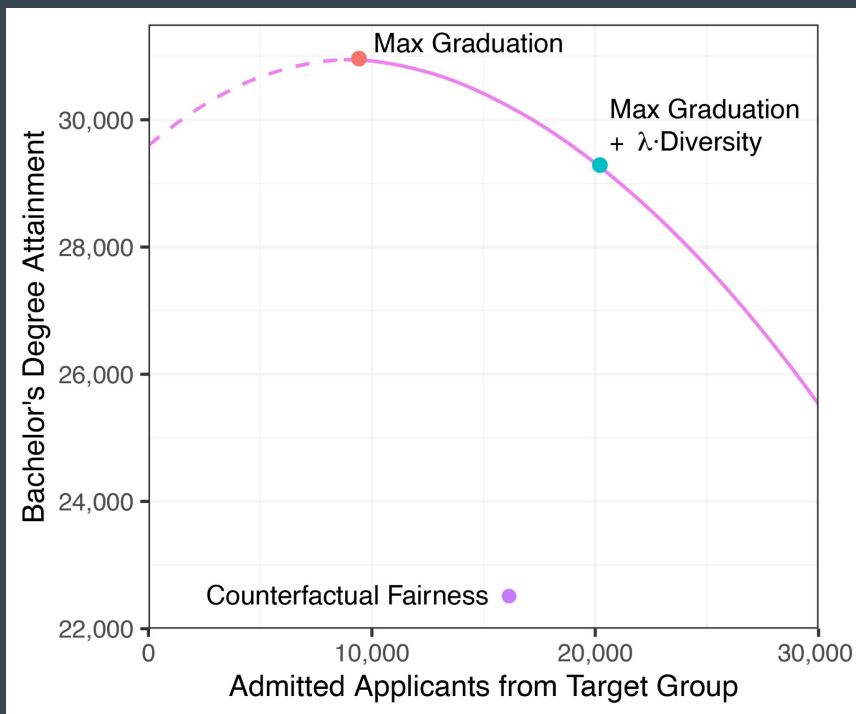# Illustrative example

<u>Theoretical result</u>: Under mild assumptions, counterfactual fairness requires decisions to ignore race and all downstream covariates



(different opportunities)

Race

Test Score

Decision

**Theoretical result**: Under mild assumptions, counterfactual fairness requires decisions to ignore race and all downstream covariates

# Theoretical result: Under mild assumptions, counterfactual fairness requires decisions to ignore race and all downstream covariates



(different opportunities)
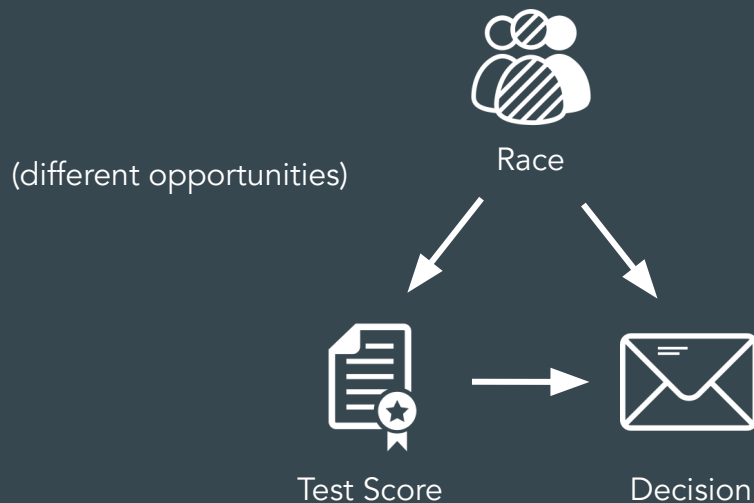
Race

Test Score

Decision
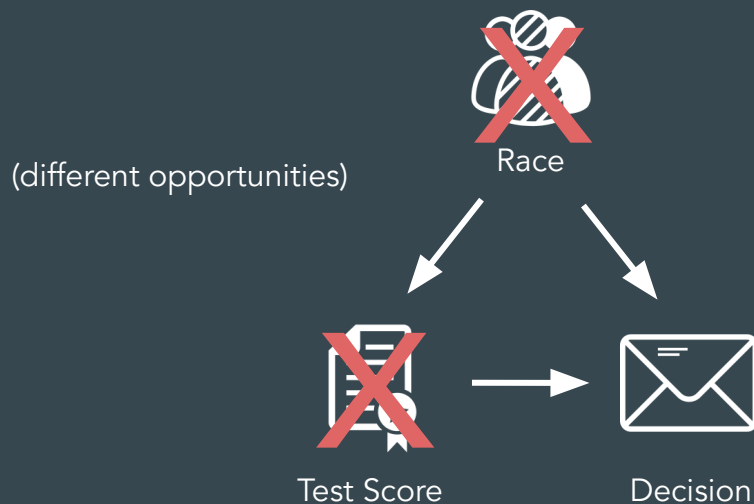
$$D(\ \ \ ,\ \ \ ) =$$

Test Score    Race

Randomized Lottery

Theoretical result: Under mild assumptions, counterfactual fairness requires decisions to ignore race and all downstream covariates



(different opportunities)

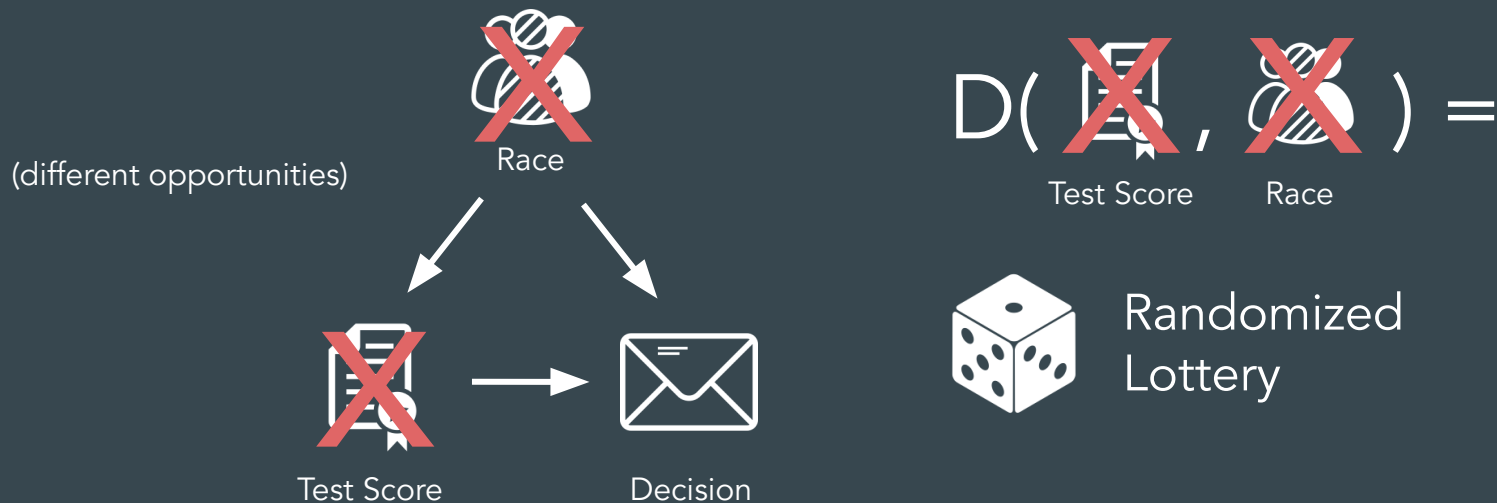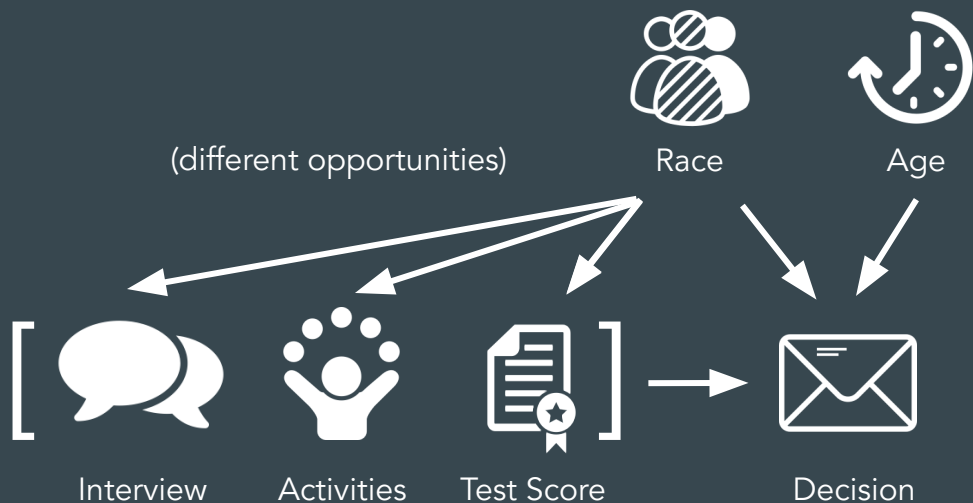Race          Age

[ Interview   Activities   Test Score ] → Decision

**Theoretical result**: Under mild assumptions, counterfactual fairness requires decisions to ignore race and all downstream covariates
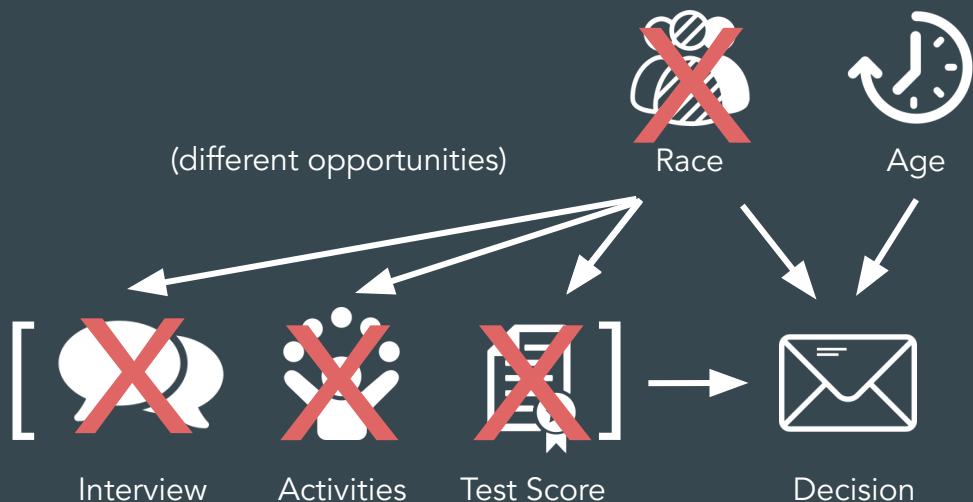
(different opportunities)

Race

Age

Interview    Activities    Test Score    Decision

Decisions based exclusively on age

# Proof sketch

D(T = Low, Race = Majority)

D(T = Med., Race = Majority)

D(T = High, Race = Majority)

D(T = Low, Race = Minority)

D(T = Med., Race = Minority)

D(T = High, Race = Minority)

# Proof sketch

D(T = Low, Race = Majority)

D(T = Low, Race = Minority)

D(T = Med., Race = Majority)

D(T = Med., Race = Minority)

D(T = High, Race = Majority)

D(T = High, Race = Minority)

# Proof sketch



D(T = Low, Race = Majority)

D(T = Med., Race = Majority)

D(T = High, Race = Majority)

D(T = Low, Race = Minority)

D(T = Med., Race = Minority)

D(T = High, Race = Minority)

# Proof sketch



D(T = Low, Race = Majority)

D(T = Low, Race = Minority)

D(T = Med., Race = Majority)

D(T = Med., Race = Minority)

D(T = High, Race = Majority)

D(T = High, Race = Minority)
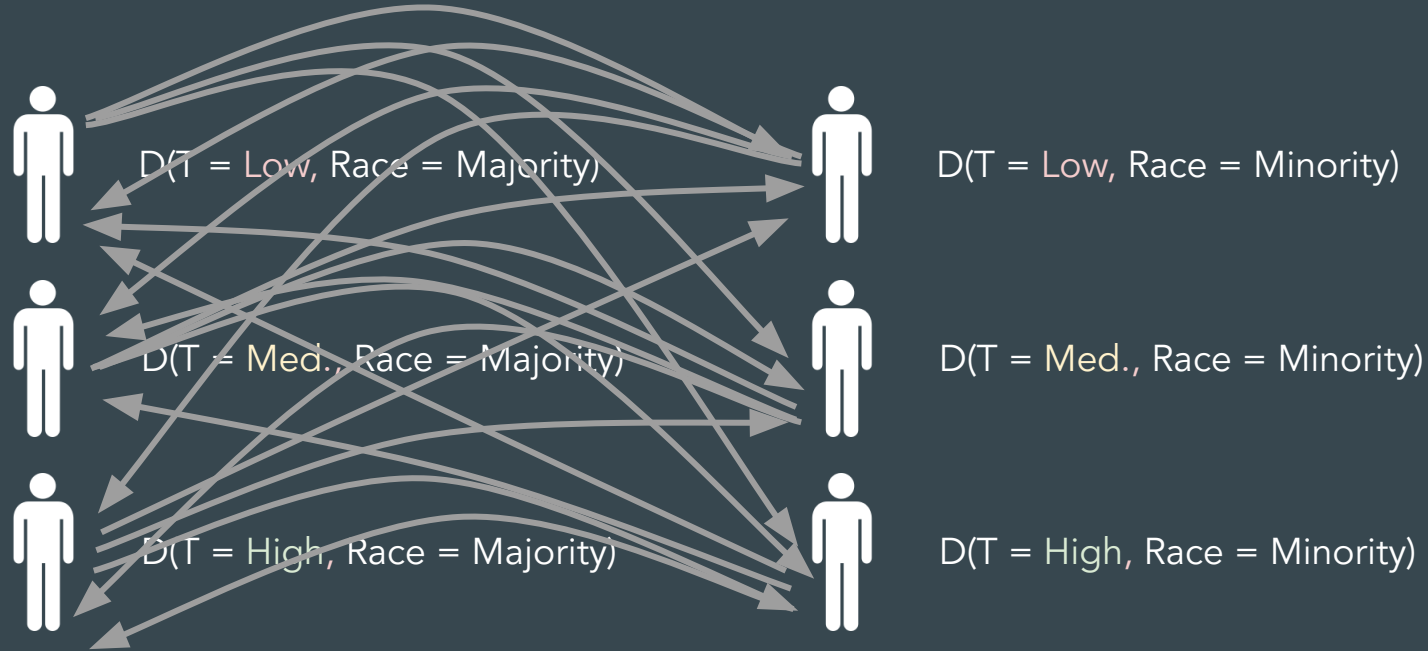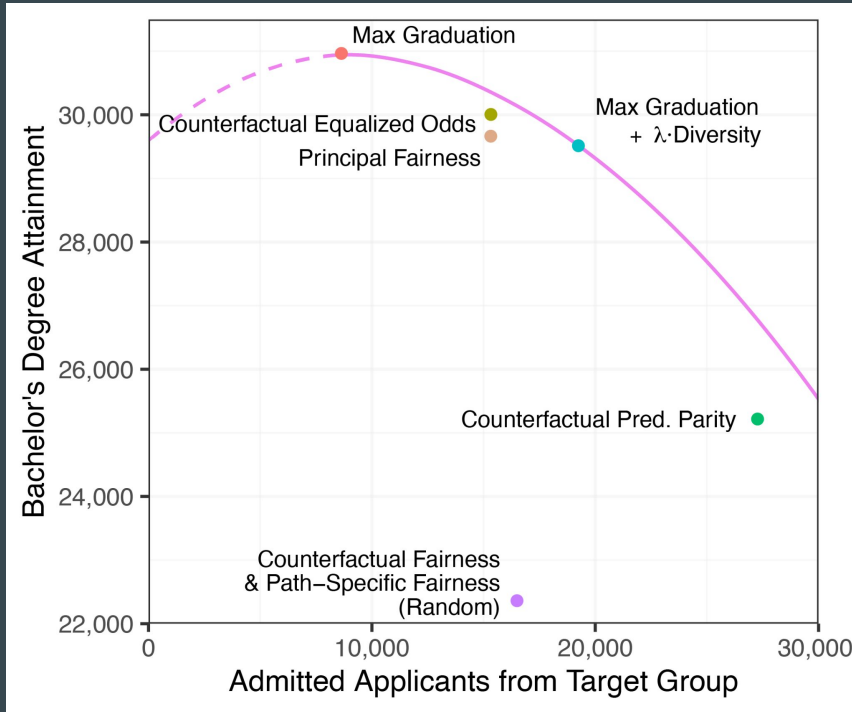
# Causal fairness taxonomy [see paper]

Family 1: Limit direct and indirect effects of race on decision

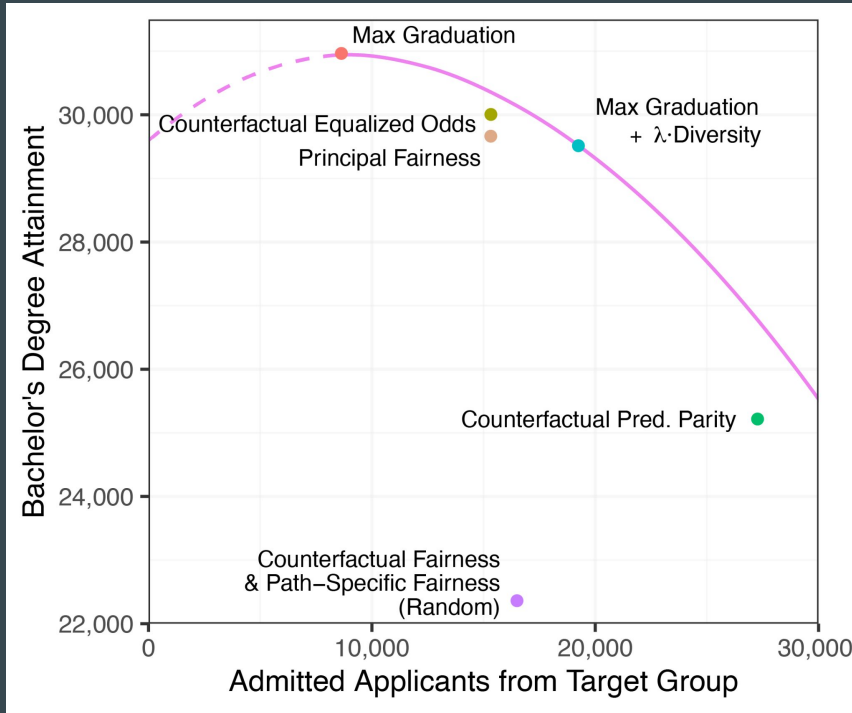- Counterfactual fairness
- Path-specific fairness

Family 2: Limit counterfactual disparities between groups

- Counterfactual equalized odds
- Counterfactual predictive parity
- Principal fairness

# Key theoretical result #2

# Key theoretical result #2



Causal Fairness
(Family 1 and 2)

Decreased
Degree
Attainment

Decreased
Class Diversity

# Key theoretical result #2

In *almost every* case (in a measure theoretic sense) it is <u>impossible</u> to satisfy prominent causal fairness definitions and be Pareto optimal



Causal Fairness

(Family 1 and 2)

Decreased Degree Attainment

Decreased Class Diversity

# Summary

- Causal fairness definitions lead to Pareto inefficient decisions, perversely harming the groups they were designed to protect

- Directly optimizing for desired outcomes (e.g. degree attainment, diversity) may be preferable

# Thank You!



Full Paper

H. Nilforoshan*, J. Gaebler*, R. Shroff, & S. Goel. "Causal Conceptions of Fairness and their Consequences." *International Conference on Machine Learning* (ICML 2022).



Technical Blog Post

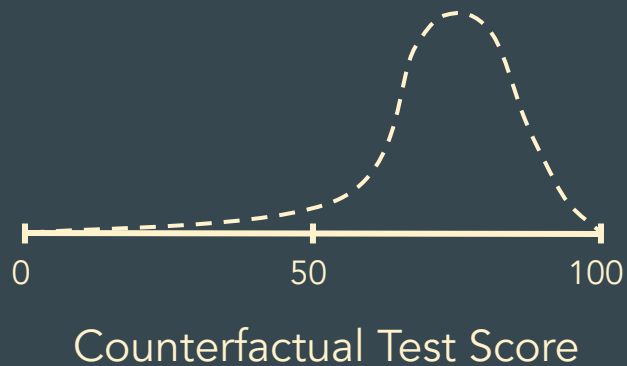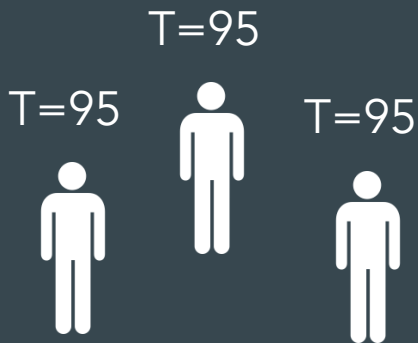jgaeb.com/2022/07/18/prevalence.html

[jgaeb.com                                    jgaeb@stanford.edu]
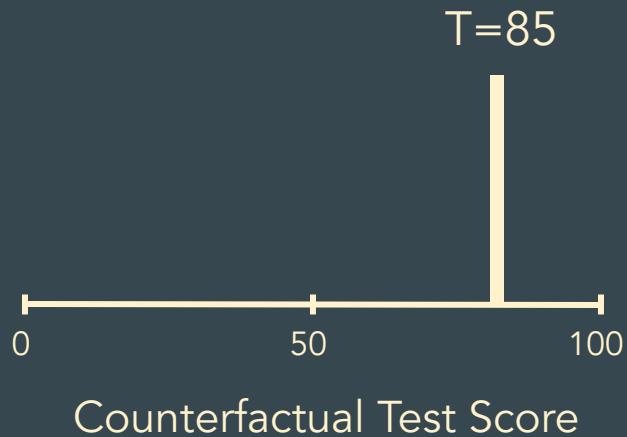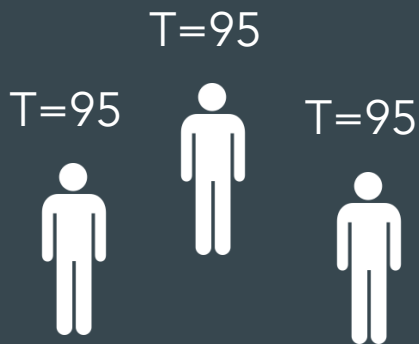
[hamedn.com                              hamedn@cs.stanford.edu]

# Assumptions

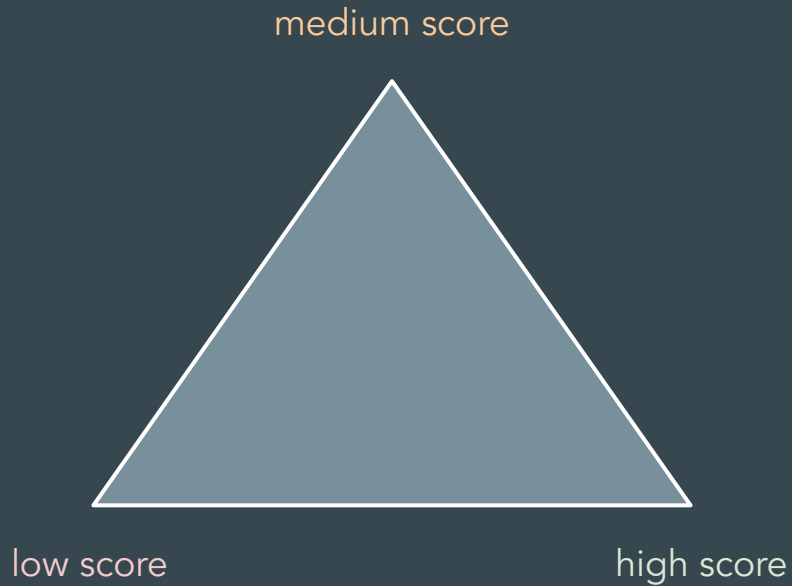There is variance in the counterfactual distribution of covariates



Counterfactual Test Score

# Assumptions

There is variance in the counterfactual distribution of covariates

# What do we mean by "almost every"?

medium score

low score                    high score

# What do we mean by "almost every"?

# What do we mean by "almost every"?



medium score
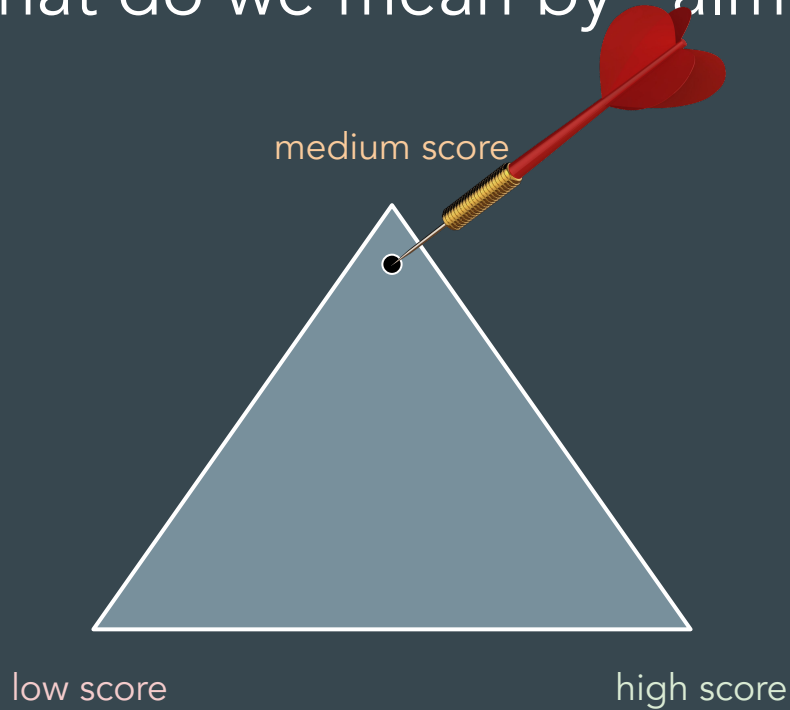
low score          high score

$P(\text{📄} = \text{low}) \qquad = 0.05$

$P(\text{📄} = \text{medium}) \quad = 0.05$

$P(\text{📄} = \text{high}) \qquad = 0.90$
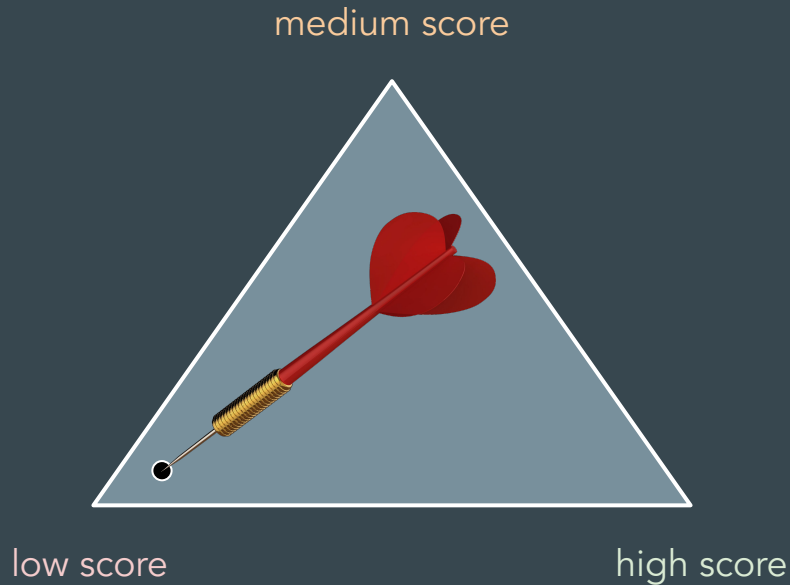
# What do we mean by "almost every"?



medium score

low score          high score

$$P(\text{📄} = \text{low}) \quad\quad = 0.05$$

$$P(\text{📄} = \text{medium}) \quad = 0.90$$

$$P(\text{📄} = \text{high}) \quad\quad = 0.05$$
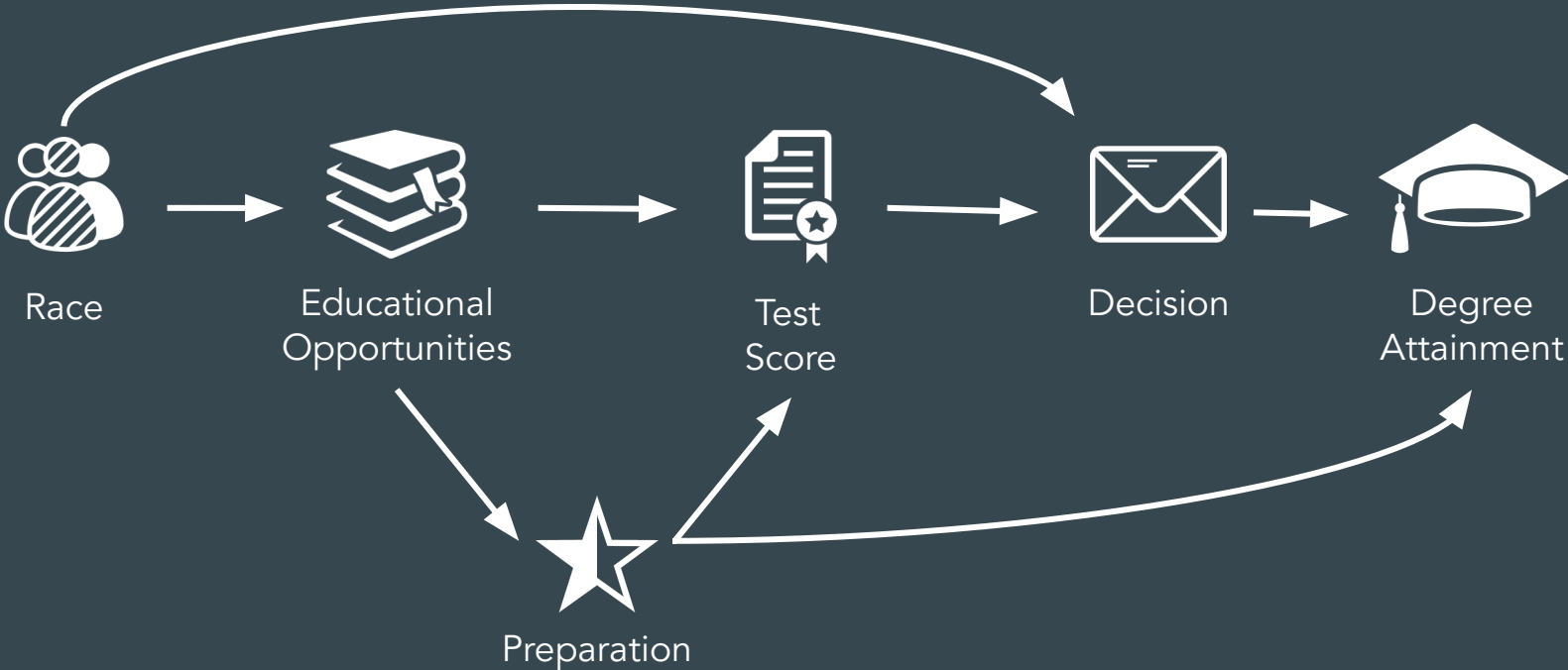
# What do we mean by "almost every"?



medium score

low score

high score

$P(\text{📄} = \text{low}) = 0.90$

$P(\text{📄} = \text{medium}) = 0.05$

$P(\text{📄} = \text{high}) = 0.05$

$P(\text{👥 Pareto Inefficient} \mid \text{🎯 Randomly Chosen Distribution}) = 1.0$

# Simulation variables



Race

Educational
Opportunities

Test
Score

Decision

Degree
Attainment

Preparation

# Key idea

# Key idea

# Key idea