

# Machine Learning Bias

---

Nick Eubank

# Supervised Machine Learning

Conceptual that because Supervised Machine Learning (SML) is built on math, it can't be biased.

SMLs just try to replicate behavior in training data.

# Supervised Machine Learning

Conceptual that because Supervised Machine Learning (SML) is built on math, it can't be biased.

SMLs just try to **replicate behavior in training data**.

If training data is biased, algorithm will be too, *even if implemented perfectly!*

# Supervised Machine Learning

We're building SML model designed to predict performance reviews using resumes.

# Supervised Machine Learning

We're building SML model designed to predict performance reviews using resumes.

Train using resumes and performance evaluations of current employees.

# Supervised Machine Learning

We're building SML model designed to predict performance reviews using resumes.

Train using resumes and performance evaluations of current employees.

- Help decide who to hire.

# Supervised Machine Learning

We're building SML model designed to predict performance reviews using resumes.

Train using resumes and performance evaluations of current employees.

- Help decide who to hire.

If supervisors tend to discriminate against women, then our SML will look for signals that an applicant is a woman, since they can use this to give women lower reviews, better matching the training data.

# Proxies

OK, but what if I don't include data on gender, race, sexuality, etc. in my model?

*Everything* in society is correlated:



# Proxies

OK, but what if I don't include data on gender, race, sexuality, etc. in my model?

*Everything* in society is correlated:

- Going to a women's college (Scripps College, Barnard College)

OK, but what if I don't include data on gender, race, sexuality, etc. in my model?

*Everything* in society is correlated:

- Going to a women's college (Scripps College, Barnard College)
- Going to a Historical Black University (Howard University)

OK, but what if I don't include data on gender, race, sexuality, etc. in my model?

*Everything* in society is correlated:

- Going to a women's college (Scripps College, Barnard College)
- Going to a Historical Black University (Howard University)
- Many activities are gender-correlated (Yoga, Football)

OK, but what if I don't include data on gender, race, sexuality, etc. in my model?

*Everything* in society is correlated:

- Going to a women's college (Scripps College, Barnard College)
- Going to a Historical Black University (Howard University)
- Many activities are gender-correlated (Yoga, Football)
- Geography is *extremely* correlated with race and income (Princeton Review)

# Proxies

OK, but what if I don't include data on gender, race, sexuality, etc. in my model?

*Everything* in society is correlated:

- Going to a women's college (Scripps College, Barnard College)
- Going to a Historical Black University (Howard University)
- Many activities are gender-correlated (Yoga, Football)
- Geography is *extremely* correlated with race and income (Princeton Review)

In COMPAS, **race wasn't in the model.**

## Target an Unbiased Outcome

## Target an Unbiased Outcome

- In hiring example, variables correlated with gender created bias because the target (performance evaluations) were biased!

## Target an Unbiased Outcome

- In hiring example, variables correlated with gender created bias because the target (performance evaluations) were biased!

Less biased targets will reduce the incentive for your algorithm to be biased.



# Target an Unbiased Outcome

- In hiring example, variables correlated with gender created bias because the target (performance evaluations) were biased!

Less biased targets will reduce the incentive for your algorithm to be biased.

Picking unbiased outcomes is not as easy as it seems...

## Target an Unbiased Outcome

COMPAS: Predicted probability of future arrest.

## Target an Unbiased Outcome

COMPAS: Predicted probability of future arrest.  
Arrests are pretty objective, right?

FIGURE 6A.

## Rates of Drug Use and Sales, by Race

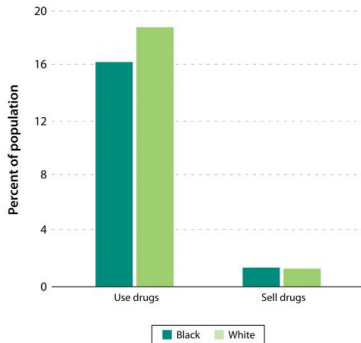
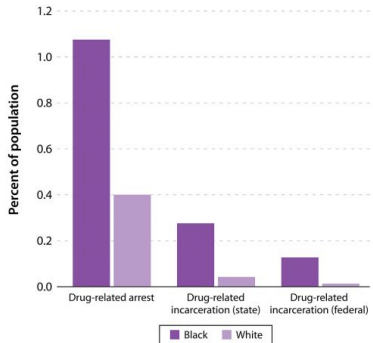


FIGURE 6B.

## Rates of Drug-Related Criminal Justice Measures, by Race



At the state level, blacks are about 6.5 times as likely as whites to be incarcerated for drug-related crimes.

Source: BLS n.d.c; Carson 2015; Census Bureau n.d.; FBI 2015; authors' calculations.

Probability of arrest  $\neq$  probability of committing a crime

## Target an Unbiased Outcome

Medical algorithm **thought** future spending on treatments was unbiased

## Target an Unbiased Outcome

Medical algorithm **thought** future spending on treatments was unbiased

The relationship between actual ailments and spending varies by race!

# Supervised Machine Learning

SML models are designed to find any patterns they can to help predict outcomes / classify records.

# Supervised Machine Learning

SML models are designed to find any patterns they can to help predict outcomes / classify records.

Because sexism, racism, xenophobia, homophobia, etc. shape outcomes in the world,



# Supervised Machine Learning

SML models are designed to find any patterns they can to help predict outcomes / classify records.

Because sexism, racism, xenophobia, homophobia, etc. shape outcomes in the world,

~> SMLs generally **perform better** when they are sexist/racist/xenophobic/homophobic!

# Biased By Design

Bias in Machine Learning isn't the result of negligence.

# Biased By Design

Bias in Machine Learning isn't the result of negligence.

So long as society has biases, Machine Learning has an affirmative incentive to be biased too!