Tufts COMP 135: Introduction to Machine Learning https://www.cs.tufts.edu/comp/135/2019s/

Fairness, Ethics, and Machine Learning





Many ideas/slides attributable to: Alexandra Chouldechova Moritz Hardt Prof. Mike Hughes

Fairness: Unit Objectives

- How to think systematically about end-to-end ML
 - Where does data come from?
 - What features am I measuring? What protected information can leak in unintentionally?
 - Who will be impacted?
- How to define and measure notions fairness
 - Use concepts: accuracy, TPR, FPR, PPV, NPV
 - What is achievable? What is impossible?

Example Concerns about Fairness

Unfair image search

Who's a CEO? Google image results can shift gender biases

UNIVERSITY OF WASHINGTON



🖨 PRINT 🛛 🖾 E-MAIL



Unfair Word Embeddings

$\overrightarrow{\mathrm{man}} - \overrightarrow{\mathrm{woman}} \approx \overrightarrow{\mathrm{king}} - \overrightarrow{\mathrm{queen}}$ $\overrightarrow{\mathrm{man}} - \overrightarrow{\mathrm{woman}} \approx \overrightarrow{\mathrm{computer programmer}} - \overrightarrow{\mathrm{homemaker}}$

Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings

Tolga Bolukbasi¹, Kai-Wei Chang², James Zou², Venkatesh Saligrama^{1,2}, Adam Kalai² ¹Boston University, 8 Saint Mary's Street, Boston, MA ²Microsoft Research New England, 1 Memorial Drive, Cambridge, MA tolgab@bu.edu, kw@kwchang.net, jamesyzou@gmail.com, srv@bu.edu, adam.kalai@microsoft.com

Unfair Hiring?

; TheUpshot

ROBO RECRUITING

Can an Algorithm Hire Better Than a Human?



Claire Cain Miller @clairecm JUNE 25, 2015

"[H]iring could become faster and less expensive, and [...] lead recruiters to more highly skilled people who are better matches for their companies. Another potential result: a more diverse workplace. The software relies on data to surface candidates from a wide variety of places and match their skills to the job requirements, free of human biases."

Miller (2015)

Job Ad Classifier: Is this fair?

Google's online advertising system, for instance, showed an ad for highincome jobs to men much more often than it showed the ad to women, <u>a</u> <u>new study</u> by Carnegie Mellon University researchers found.

Unfair Recidivism Prediction

Two Petty Theft Arrests



Focus: Binary Classifier

- Let's say we have two groups, A and B
 - Could be any protected group (race / gender / age)
- We're trying to build a binary classifier that will predict individuals as HIGH or LOW risk
 - Likelihood of recidivism
 - Ability to pay back a loan

Group Discussion

- When should protected information (gender, race, age, etc) be provided as input to a predictor?
 - Can you build a "race-blind" classifier?
- How could we measure if the predictions are fair?
 - Is it enough to ensure accuracy parity?
 - ACC(group A) = ACC(group B)

Notation for Binary Classifier

		classifier calls "negative" "positive" C=0 C=1	
true outcome	Y=0	TN	FP
	Y=1	FN	TP

Example of Accuracy Parity

		Group A	Group B
true outcomes 1 = would fail to appear in court	Y	0011	0011
classifier prediction 1 = too risky for bail	С	0 0 0 0	1 1 1 1

Is this fair?

Case Study: The COMPAS future crime prediction algorithm



COMPAS classifier



other features (e.g. demographics, questionnaire answers, family history)

HIGH RISK of future crime *hold in jail before trial*

LOW RISK of future crime *release before trial*

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

Monkey Cage

A computer program used for bail and sentencing decisions was labeled biased against blacks. It's actually not that clear.

By Avi Feller, Emma Pierson, Sam Corbett-Davies and Sharad Goel October 17, 2016

The COMPAS tool assigns defendants scores from 1 to 10 that indicate how likely they are to reoffend based on more than 100 <u>factors</u>, including age, sex and criminal history. Notably, race is not used. These scores profoundly affect defendants' lives: defendants who are defined as medium or high risk, with scores of 5-10, are more likely to be detained while awaiting trial than are low-risk defendants, with scores of 1-4.

We obtained the risk scores assigned to more than 7,000 people arrested in Broward County, Florida, in 2013 and 2014 and checked to see how many were charged with new crimes over the next two years, the <u>same benchmark used</u> by the creators of the algorithm.

The score proved remarkably unreliable in forecasting violent crime: Only 20 percent of the people predicted to commit violent crimes actually went on to do so.

When a full range of crimes were taken into account — including misdemeanors such as driving with an expired license — the algorithm was somewhat more accurate than a coin flip. Of those deemed likely to re-offend, 61 percent were arrested for any subsequent crimes within two years. Northpointe's core product is a set of scores derived from <u>137 questions</u> that are either answered by defendants or pulled from criminal records. Race is not one of the questions.

Exclude the current case for these questions.

- How many times has this person been arrested before as an adult or juvenile (criminal arrests only)?
- How many prior juvenile felony offense arrests?
 0 1 2 3 4 5+

Family Criminality

- 32. If you lived with both parents and they later separated, how old were you at the time? ✓ Less than 5 5 to 10 11 to 14 15 or older Does Not Apply
- 33. Was your father (or father figure who principally raised you) ever arrested, that you know of?
 ☑ No □ Yes
- 34. Was your mother (or mother figure who principally raised you) ever arrested, that you know of?
 ☑ No □ Yes

Social Environment

- 69. Is it easy to get drugs in your neighborhood?
 ☑ No □ Yes
- 70. Are there gangs in your neighborhood?
 No Yes

Anger

- 121. "Some people see me as a violent person."
 □ Strongly Disagree □ Disagree ☑ Not Sure 0
- 122. "I get into trouble because I do things without to ☐ Strongly Disagree ☐ Disagree ☑ Not Sure

Full Document: <u>https://www.documentcloud.org/documents/2702103-Sample-Risk-Assessment-COMPAS-CORE.html</u>ike Hughes - Tufts COMP 135 - Spring 2019

ProPublica says: "Groups have different **False Pos. Rates**"

Prediction Fails Differently for Black Defendants				
WHITE AFRICAN AMERIC				
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%		
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%		

Compas Team Says: "Groups have same **predictive value**"



 Within each risk category, the proportion of defendants who reoffend is approximately the same regardless of race; this is Northpointe's definition of fairness.

False Positive Rate =

н

FP + TN

FP

• When true outcome is o, how often does classifier say "1".

		classifier calls	
		"negative" "positive" C=0 C=1	
true outcome	Y=0	TN	FP
	Y=1	FN	TP

True Positive Rate = $\frac{TP}{TP + FN}$

• When true outcome is 1, how often does classifier say "1".

1

		classifier calls	
		"negative" "positive" C=0 C=1	
true outcome	Y=0	TN	FP
	Y=1	FN	TP

Positive Predictive Value = $\frac{TP}{TP + FP}$

When classifier says "1", how often is true label 1.

		classifi	er calls
		"negative" C=0	"positive" C=1
true outcome	Y=0	TN	FP
	Y=1	FN	TP

Negative Predictive Value = $\frac{TN}{TN + FN}$

When classifier says "0", how often is true label 0.

		classifi	er calls
		"negative" C=0	"positive" C=1
true outcome	Y=0	TN	FP
	Y=1	FN	TP

ProPublica says: "Groups have different **False Pos. Rates**"

Prediction Fails Differently for Black Defendants				
WHITE AFRICAN AMERICA				
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%		
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%		

		classifier calls	
		"negative" C=0	"positive" C=1
true outcome	Y=0	TN	FP
	Y=1	FN	TP



Compas Team Says: "Groups have same **predictive value**"



<u>Worksheet</u>

True Positive Rate (TPR)	TP TP + FN	subject who is positive will be called positive
False Positive Rate (FPR)	FP FP + TN	subject who is negative will be called positive
Positive Predictive Value (PPV)	TP TP + FP	subject called positive will actually be positive

		classifie	er calls
		"negative" C=0	"positive" C=1
true	Y=0	TN	FP
outcome	Y=1	FN	TP

Equation of the Day

$FPR = \frac{p}{1-p} \frac{1-PPV}{PPV} TPR$

where prevalence p = Pr(Y = 1)

If two groups have different p values, can we simultaneously have TPR parity AND FPR parity AND PPV parity AND NPV parity?

Bias in Criminal Risk Scores Is Mathematically Inevitable, Researchers Say

The scholars set out to address this question: Since blacks are re-arrested more often than whites, is it possible to create a formula that is equally predictive for all races without disparities in who suffers the harm of incorrect predictions?

Working separately and using different methodologies, four groups of scholars all reached the same conclusion. It's not.

<u>https://www.propublica.org/article/bias-in-criminal-risk-scores-is-mathematically-inevitable-researchers-say</u>

Unless classifier is perfect, **must** chose one: Disparate Treatment (PPV or NPV not equal) or Disparate Impact (FPR or TPR not equal)

<u>Try demo of making decisions from</u> <u>risk scores:</u>

goo.gl/P8rmA3

Fairness: Unit Objectives

- How to think systematically about end-to-end ML
 - Where does data come from?
 - What features am I measuring? What protected information can leak in unintentionally?
 - Who will be impacted?
- How to define and measure notions fairness
 - Use concepts: accuracy, TPR, FPR, PPV, NPV
 - What is achievable? What is impossible?