



Department of Computer Science
UNIVERSITY OF COLORADO BOULDER



Machine Learning: Yoshinari Fujinuma

University of Colorado Boulder
LECTURE 28

Slides adapted from Chenhao Tan, Jordan Boyd-Graber, Chris Ketelsen

Logistics

- No class on Friday April 30th
- Final project write-up is deadline is May 1st
 - Submission link on Canvas will be open today
- An optional makeup quiz on May 2nd 1:30 PM - 2:30 PM
 - All topics covered in this class

Learning Objectives

- Ethics, transparency, fairness in machine learning models

Outline

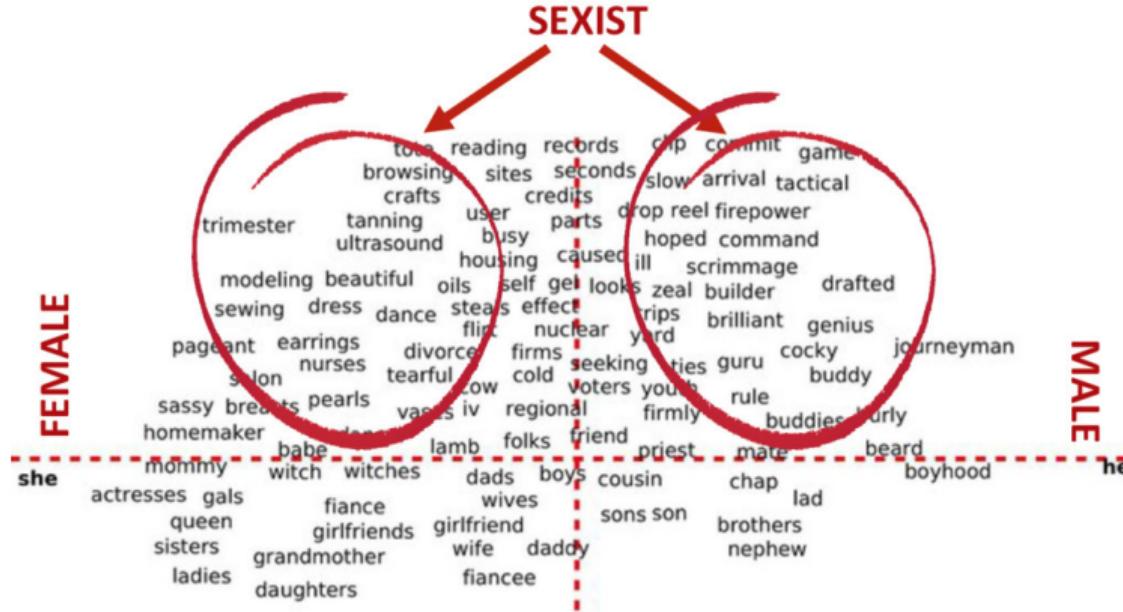
Fairness

Ethics

Transparency

Our data reflect our world ...

- Word representations learned from massive amounts of data
- Reflect prejudices and messiness of our world
- But learned representations used for many tasks
 - Detecting “bad” behavior online
 - Matching resumes to jobs
 - Recommendations
- Let’s first see an example of gender bias in word embeddings



The embedding captures gender stereotypes *and* sexism.

DEFINITIONAL

(related [Schmidt '15])

SEXIST

Easier to debias an embedding
than to debias a human

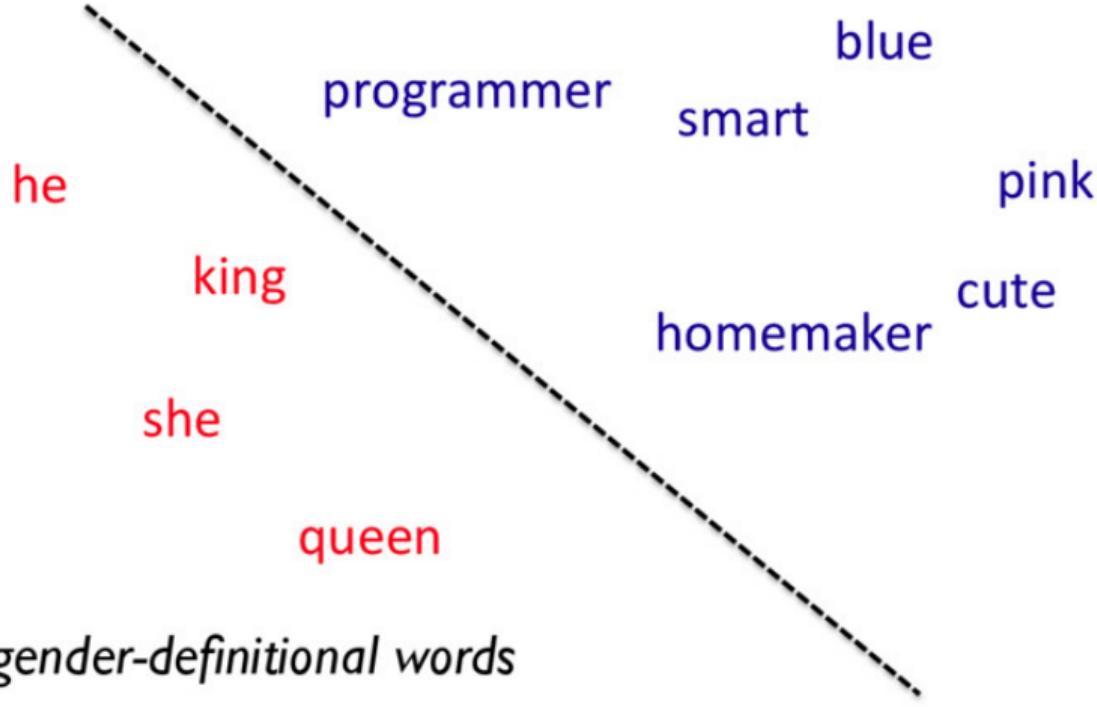
**DEFINITIONAL**

(related [Schmidt '15])

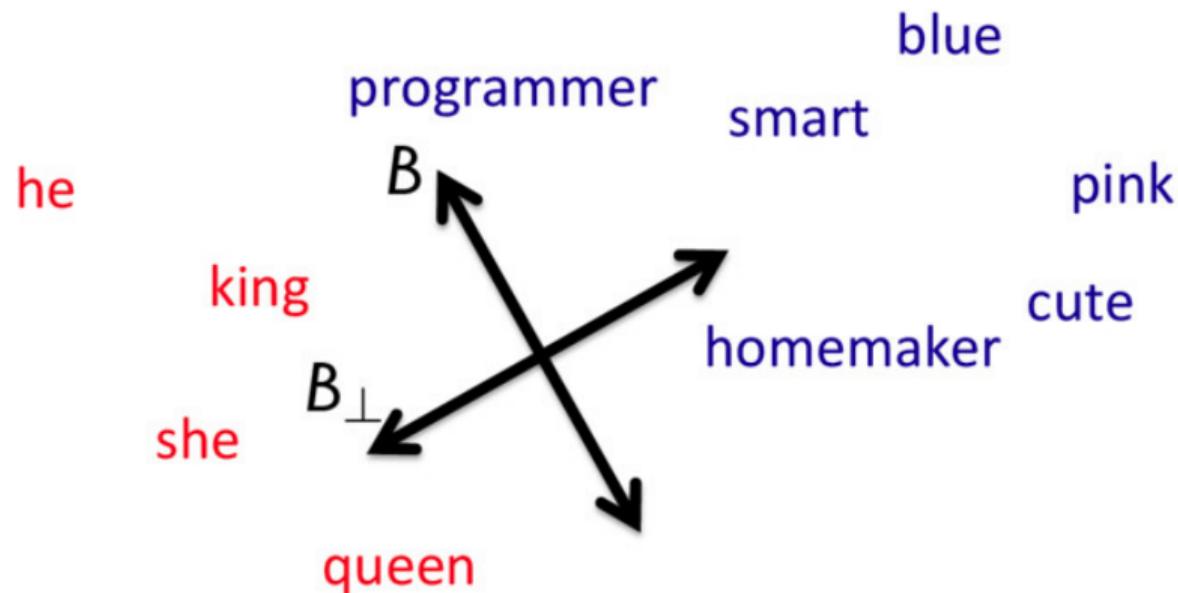
Bias encoded in some dimensions



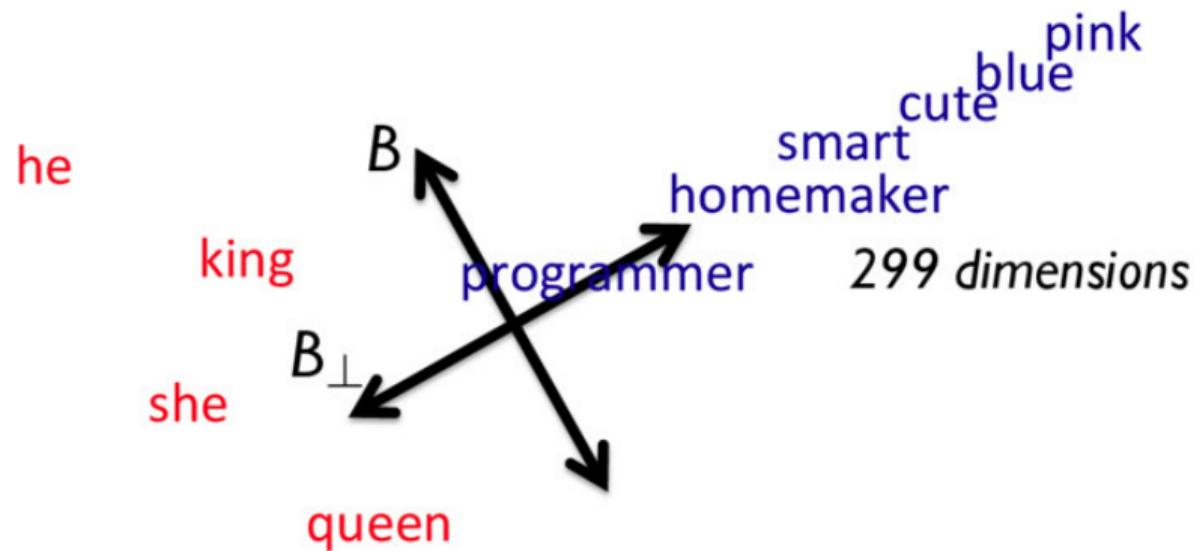
Debiasing



Debiasing



Debiasing



Outline

Fairness

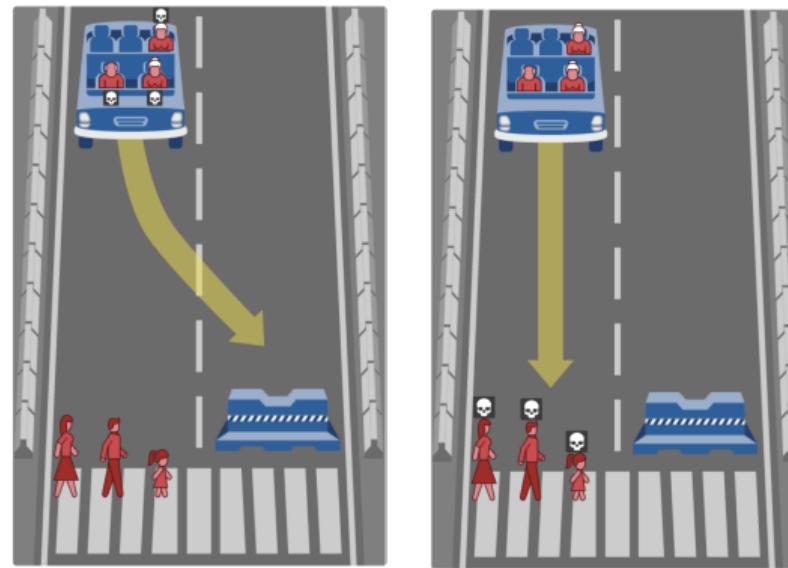
Ethics

Transparency

Case Study: What Should the Self-Driving Cars Do?

- Interactive study by MIT Moral Machine
- Main idea: In the event of an inevitable crash leading to likely loss of life, what should the car do?
- Example: Car crash will either result in death of
 - driver and several passengers
 - several pedestrians
- Debate: How do we choose?

Case Study: What Should the Self-Driving Cars Do?



<https://www.moralmachine.net/>

Case Study: Recommendation System

Filter bubble [Pariser, 2011]

- Main idea: Personalized search determines what information you see and what information you don't see
- Example: Google, Facebook, Netflix
- Debate: Is this a good thing or bad thing?

Case Study: Recommendation System

Example: Blue feed vs. Red feed

<https://graphics.wsj.com/blue-feed-red-feed/>

LIBERAL

SHOWING POSTS ABOUT:
"PRESIDENT TRUMP"

CONSERVATIVE

Allen West about an hour ago
Gotta admit, this is pretty good... Well played, President Trump.
Watch: Trump Gives a Special 'Present' to Obama... And boom goes the dynamite...
ALLENWEST.AMERICANNEWSHUB.COM

1.8K 78 240

GP Gateway Pundit 2 hours ago
VA has Alec Baldwin hanging instead of President Trump!
OUTRAGE: Veterans Affairs Offic... A photo posted by a Veterans Affairs physician sho...
THEGATEWAYPUNDIT.COM

69 301 430

CT Conservative Tribune 3 hours ago
In 8 years, the best Obama could do was give the whole Middle East over to ISIS.
In just 11 months, Donald Trump utterly crushed them.
Imagine how many lives Obama lost over those years and how many Trump just saved.

THERE'S A WAY TO REMOVE

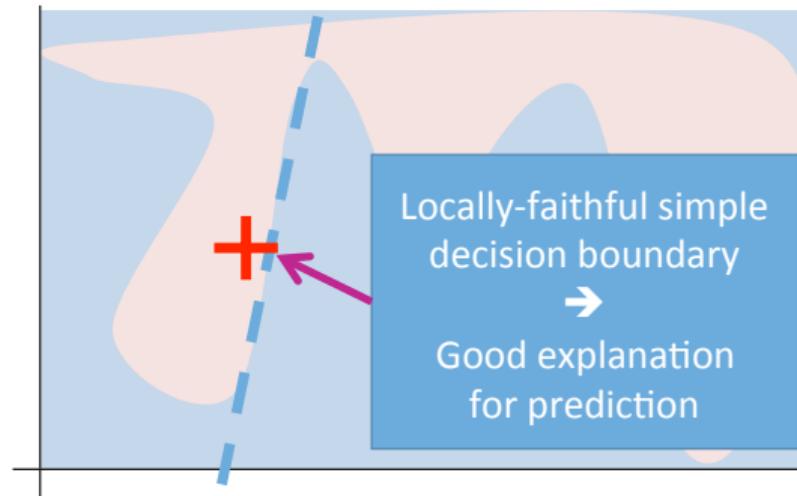
Outline

Fairness

Ethics

Transparency

LIME



Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why Should I Trust You?" Explaining the Predictions of Any Classifier. KDD 2016.
LIME: Local Interpretable Model-Agnostic Explanations

What's an Explanation

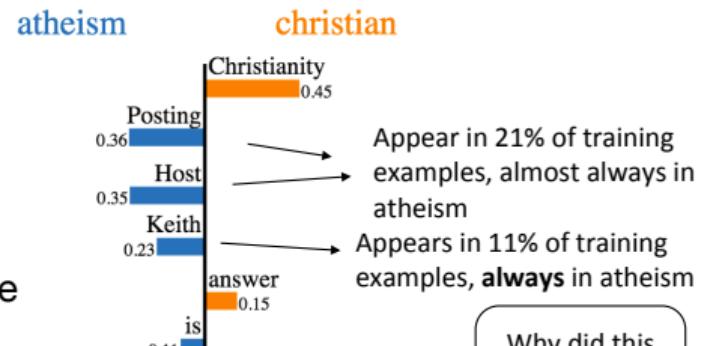
Example: Binary classification of an email

From: Keith Richards
 Subject: Christianity is the answer
 NTTP-Posting-Host: x.x.com

I think Christianity is the one true religion.
 If you'd like to know more, send me a note



Prediction probabilities



What's an Explanation

Example: Image classification



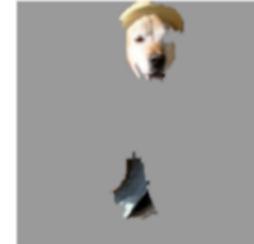
$P(\text{guitar}) = 0.32$



$P(\text{guitar}) = 0.24$



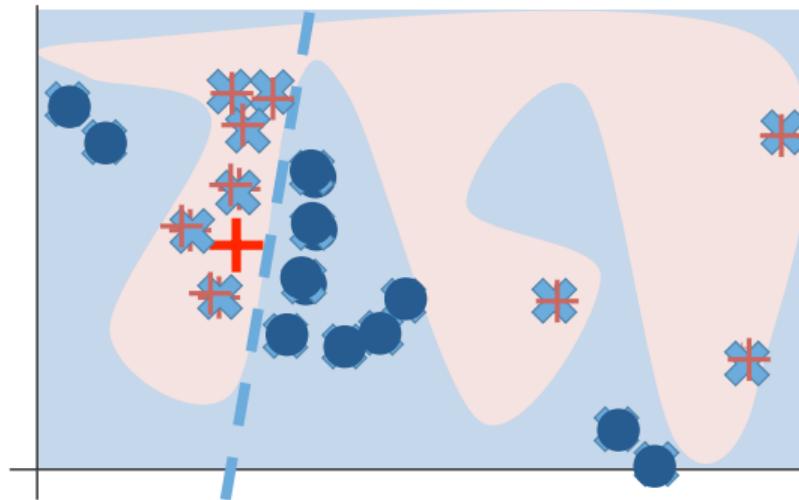
$P(\text{dog}) = 0.21$



What makes good Explanation?

- Interpretable: Humans can Understand
- Faithful: Describes Model
- Model Agnostic: Generalize to Many Models

Method



- Complicated model predicts “near” example
- Simple model explains **local variation**
- **Explains what complicated model focused on**

Is this a good Classifier?



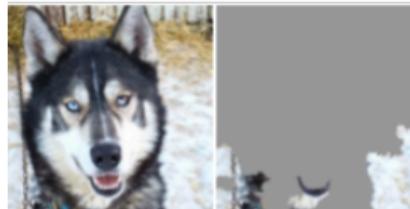
Predicted: **wolf**
True: **wolf**



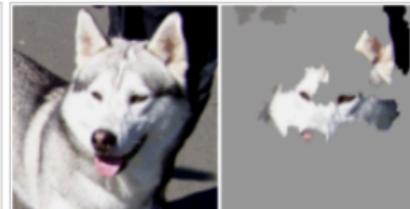
Predicted: **husky**
True: **husky**



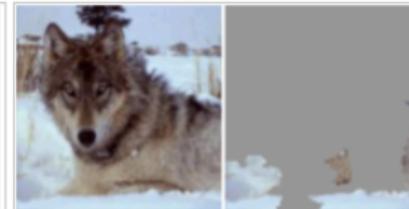
Predicted: **wolf**
True: **wolf**



Predicted: **wolf**
True: **husky**

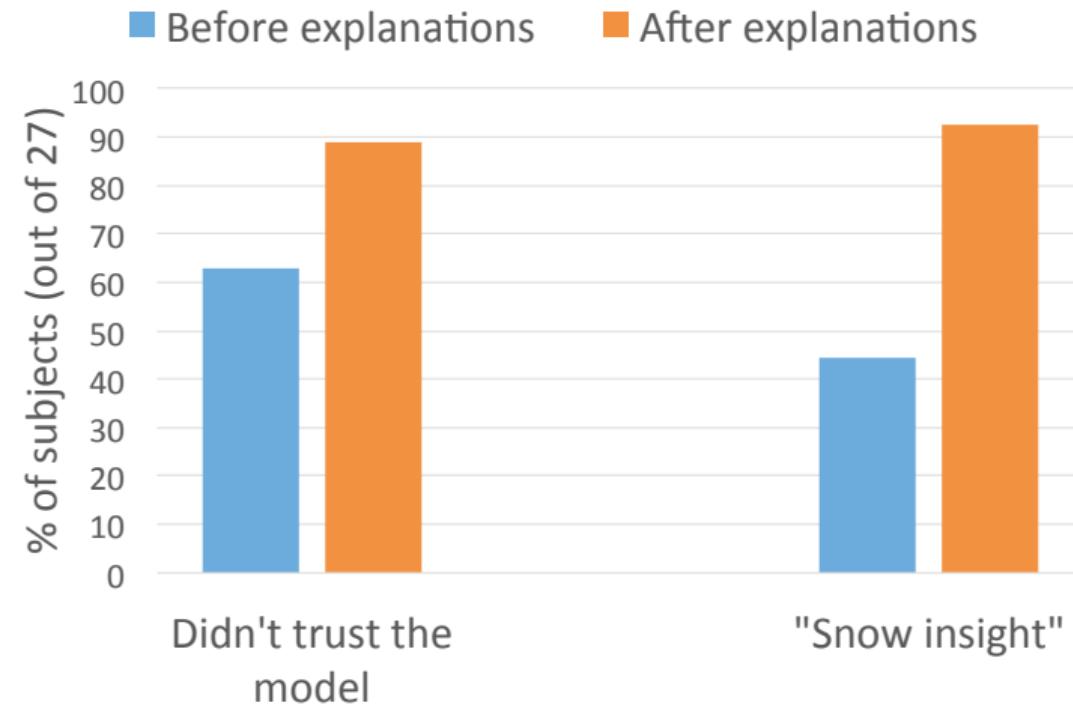


Predicted: **husky**
True: **husky**



Predicted: **wolf**
True: **wolf**

Is this a good Classifier?



Wrapping up

- Be aware of the data that you use. What data is OK to use?
- Be aware of the decision made by machine learning models.
- Be aware of the of complex models. Is it OK to use black-box models?

Finally

Thanks for participating in this class!
Special thanks to Saumya Sinha and Vignesh Karthikeyan