# CS 360: Machine Learning

## Prof. Sara Mathieson

## Fall 2020

HAVERFORD
COLLEGE

# Admin

- Office hours **today 4:30-6pm**

- **Lab 8** posted today, due Friday Nov 20
  - Keep working on getting logged in to lab machines
  - Can start Part 1 (data pre-processing)

- After Thanksgiving – **two options for capstone**
  - Midterm 2
  - Final project (posted soon)

# Outline for November 10

- Ways to repair biased data algorithmically

- Validation best practices

- Introduction to neural networks

# Outline for November 10

- Ways to repair biased data algorithmically

- Validation best practices

- Introduction to neural networks

# How can we tell if an algorithm is biased?

D: dataset with attributes X, Y

    \* X is protected

    \* Y is unprotected (other features)

"Certifying and Removing Disparate Impact" Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, Suresh Venkatasubramanian

# How can we tell if an algorithm is biased?

D: dataset with attributes X, Y

      * X is protected

      * Y is unprotected (other features)

Goal: determine outcome C (hired, admitted, etc)

"Certifying and Removing Disparate Impact" Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, Suresh Venkatasubramanian

# How can we tell if an algorithm is biased?

D: dataset with attributes X, Y

* X is protected
* Y is unprotected (other features)

Goal: determine outcome C (hired, admitted, etc)

Direct discrimination: C = f(X)

* Female instrumentalist not hired for orchestra
* Some ethnic groups not allowed to eat at a restaurant

# How can we tell if an algorithm is biased?

D: dataset with attributes X, Y

   * X is protected
   * Y is unprotected (other features)

Goal: determine outcome C (hired, admitted, etc)

Indirect discrimination: C = f(Y)

   * but strong correlation between X and Y

   * Ex: housing loans
   * Ex: programming experience

"Certifying and Removing Disparate Impact" Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, Suresh Venkatasubramanian

# Disparate Impact

X = protected attribute ⟵
Y = other attributes

$$\begin{cases} X = 0 & \text{minority group} \\ X = 1 & \text{majority group} \end{cases}$$

C = binary outcome ⟵ $\underbrace{hired,\ admitted}_{C=1}$, or $\underbrace{not}_{C=0}$

legal definition

$$P(C=1 \mid X=0) \leq 0.8\ P(C=1 \mid X=1)$$

example    40% women hired
           60% men hired    } years (?)

$$0.40 \leq 0.8\ (0.6)$$
$$\underbrace{\qquad}_{0.48}$$

yes

is disparate impact

# Disparate Impact

$C$ = outcome

**Idea: if we can predict X from Y, there could be disparate impact**

predictor → "try"  $f: Y \rightarrow X$  ← could be more than one

Balanced Error Rate (BER), $\varepsilon$ majority

$$\varepsilon = BER = \frac{P[f(y) = 0 \mid X = 1] + P[f(y) = 1 \mid X = 0]}{2}$$

error          error    minority

threshold

want: high!  $\}$  max = 0.5

|  f  | pred 0 | 1 |
|-----|--------|-----|
| true 0 | 0.5 | (0.5) |
| 1 | (0.5) | 0.5 |

$\dfrac{0.5 + 0.5}{2} = 0.5$

# Disparate Impact

$C$ = outcome

Idea: if we can predict X from Y, there could be disparate impact

Predictor "try" $\rightarrow$ $\boxed{f: Y \rightarrow X}$ $\leftarrow$ could be more than one

Balanced Error Rate (BER), $\varepsilon$ majority

relationship between X & Y

$\boxed{\varepsilon} = BER = \dfrac{P[\,f(Y) = 0 \mid X = 1\,] + P[\,f(Y) = 1 \mid X = 0\,]}{2}$

error          error          minority

actually happened

| outcome | X = 0 | X = 1 |
|---------|-------|-------|
| C = 0   | a     | b     |
| C = 1   | c     | d     |

$\beta = \dfrac{c}{a+c}$

$\varepsilon' = \dfrac{1}{2} - \dfrac{\beta}{8}$

if $\boxed{\varepsilon > \varepsilon'}$ $\Rightarrow$

no disparate impact

# Example of repair



**Figure 1: Consider the fake probability density functions shown here where the blue curve shows the distribution of SAT scores ($Y$) for $X = $ `female`, with $\mu = 550, \sigma = 100$, while the red curve shows the distribution of SAT scores for $X = $ `male`, with $\mu = 400, \sigma = 50$. The resulting fully repaired data is the distribution in black, with $\mu = 475, \sigma = 75$. Male students who originally had scores in the 95th percentile, i.e., had scores of 500, are given scores of 625 in the 95th percentile of the new distribution in $\bar{Y}$, while women with scores of 625 in $\bar{Y}$ originally had scores of 750.**

"Certifying and Removing Disparate Impact" Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, Suresh Venkatasubramanian

# Example of repair



Figure 1: Consider the fake probability density functions shown here where the blue curve shows the distribution of SAT scores ($Y$) for $X =$ female, with $\mu = 550, \sigma = 100$, while the red curve shows the distribution of SAT scores for $X =$ male, with $\mu = 400, \sigma = 50$. The resulting fully repaired data is the distribution in black, with $\mu = 475, \sigma = 75$. Male students who originally had scores in the 95th percentile, i.e., had scores of 500, are given scores of 625 in the 95th percentile of the new distribution in $\bar{Y}$, while women with scores of 625 in $\bar{Y}$ originally had scores of 750.

"Certifying and Removing Disparate Impact" Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, Suresh Venkatasubramanian

# Discussion Questions

1) What are our responsibilities as engineers to ensure that our algorithms are fair?

2) How would you handle a situation where you felt you didn't have enough data (or the right data) necessary to build your algorithm?

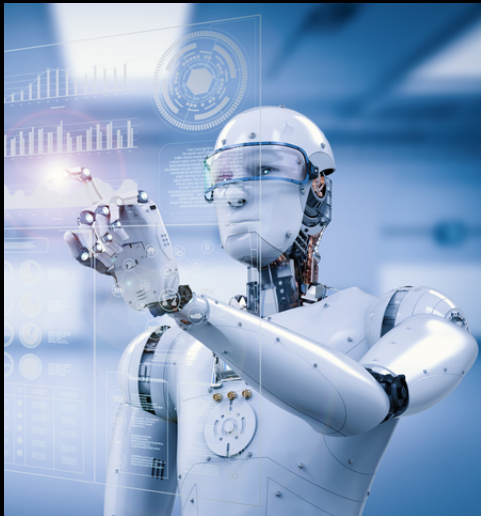3) How would you try to detect if your algorithm was making biased decisions during deployment?

# Outline for November 10

- Ways to repair biased data algorithmically

- Validation best practices

- Introduction to neural networks

# Evaluation in Practice



All Data

Training Data

Test Data

# Evaluation in Practice

# Evaluation in Practice



All Data

Training Data

learn

Test Data

evaluate

Repeat until happy

Modified from Jessica Wu

# Evaluation in Practice



NO!  Using test data as part of the model selection process

Modified from Jessica Wu

# Better: use a *validation* dataset



Modified from Jessica Wu

# Outline for November 10

- Ways to repair biased data algorithmically

- Validation best practices

- Introduction to neural networks

# MACHINE LEARNING

What society thinks I do



What my boss thinks I do



What other computer scientists think I do



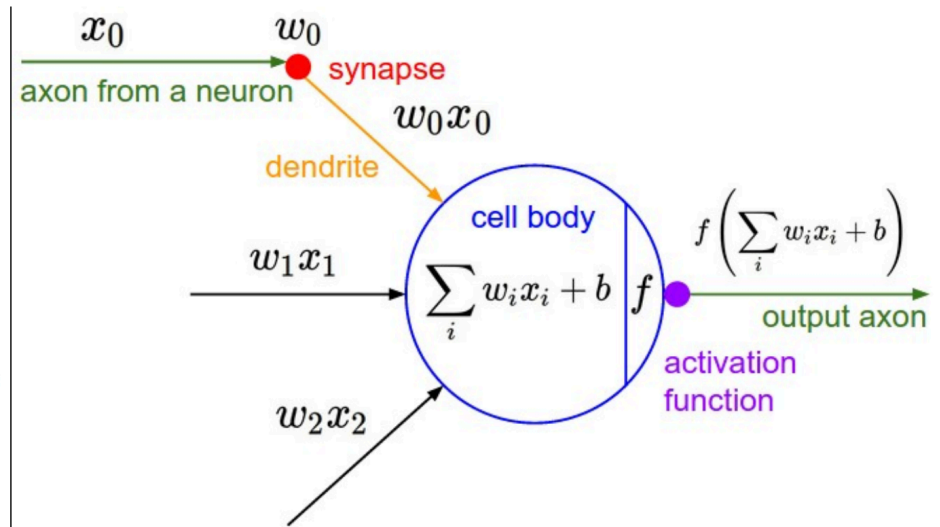What mathematicians think I do



What I think I do

```
>>> from sklearn import svm
>>> import tensorflow as tf
```

What I really do

# MACHINE LEARNING

What society thinks I do

What other computer scientists think I do

Takeaway: we should understand the methods we are using!
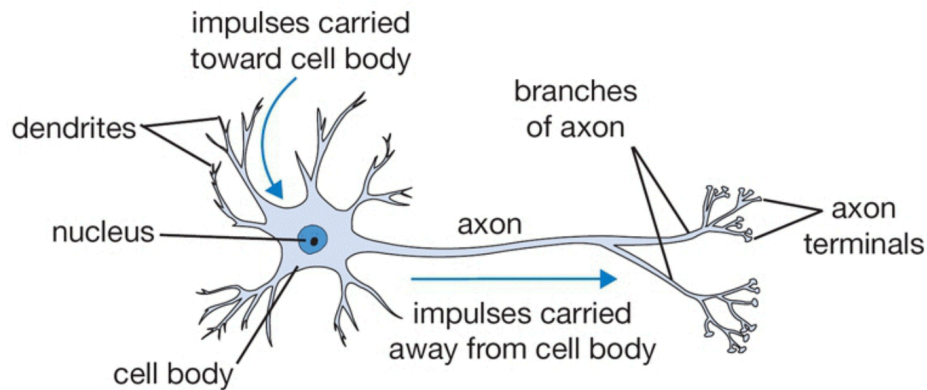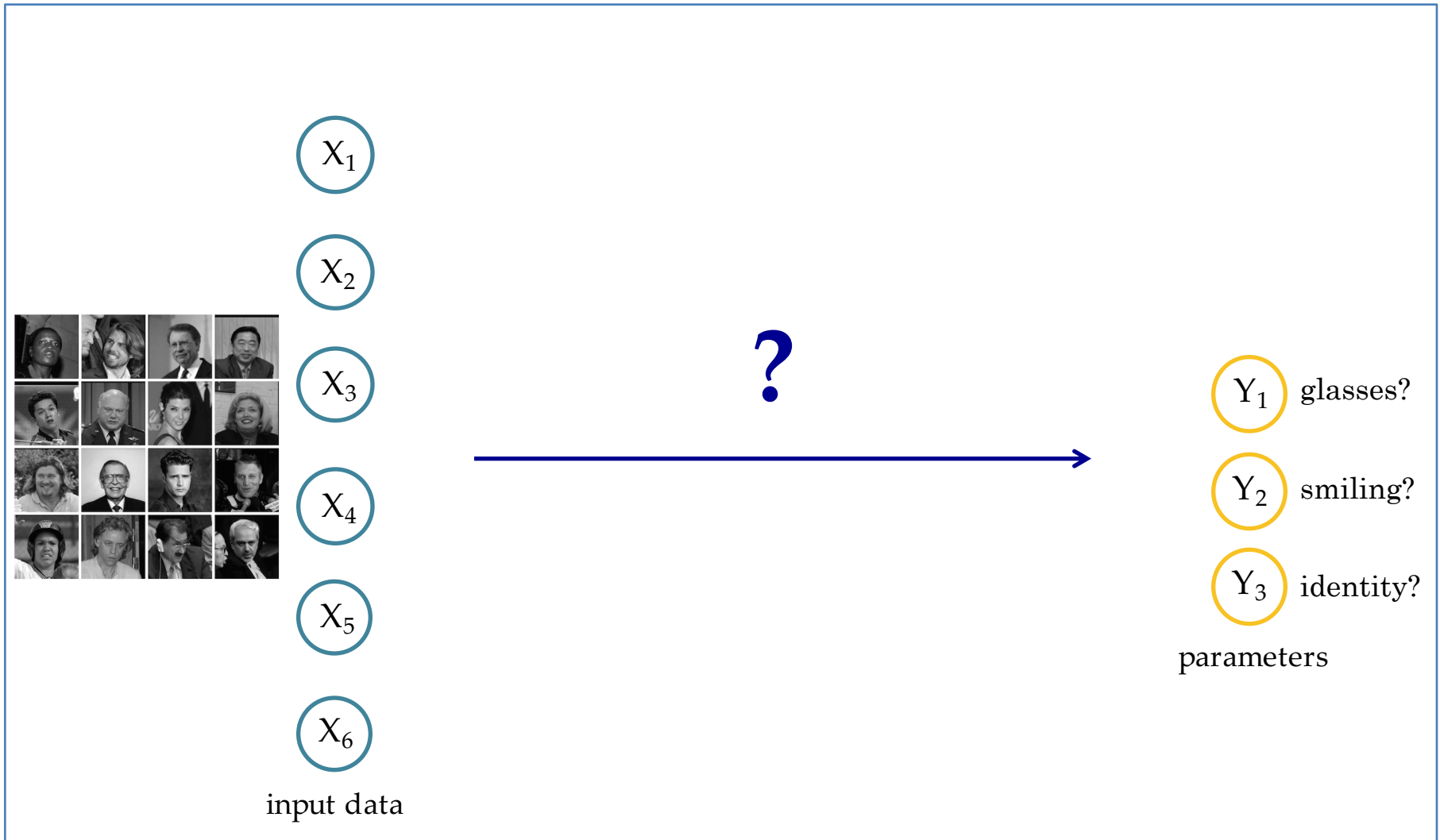
What I think I do

What mathematicians think I do

```
>>> from sklearn import svm
>>> import tensorflow as tf
```

What I really do

# Biological Inspiration

# Goal: learn from complicated inputs



$X_1$

$X_2$

$X_3$

$X_4$

$X_5$

$X_6$

input data

**?**

$Y_1$ glasses?

$Y_2$ smiling?

$Y_3$ identity?

parameters

Image: Labeled Faces in the Wild (UMass)

# Idea: transform data into lower dimension

# Multi-layer networks = "deep learning"



input data

hidden layer 1

hidden layer 2

$Y_1$ glasses?

$Y_2$ smiling?

$Y_3$ identity?

parameters

Image: Labeled Faces in the Wild (UMass)
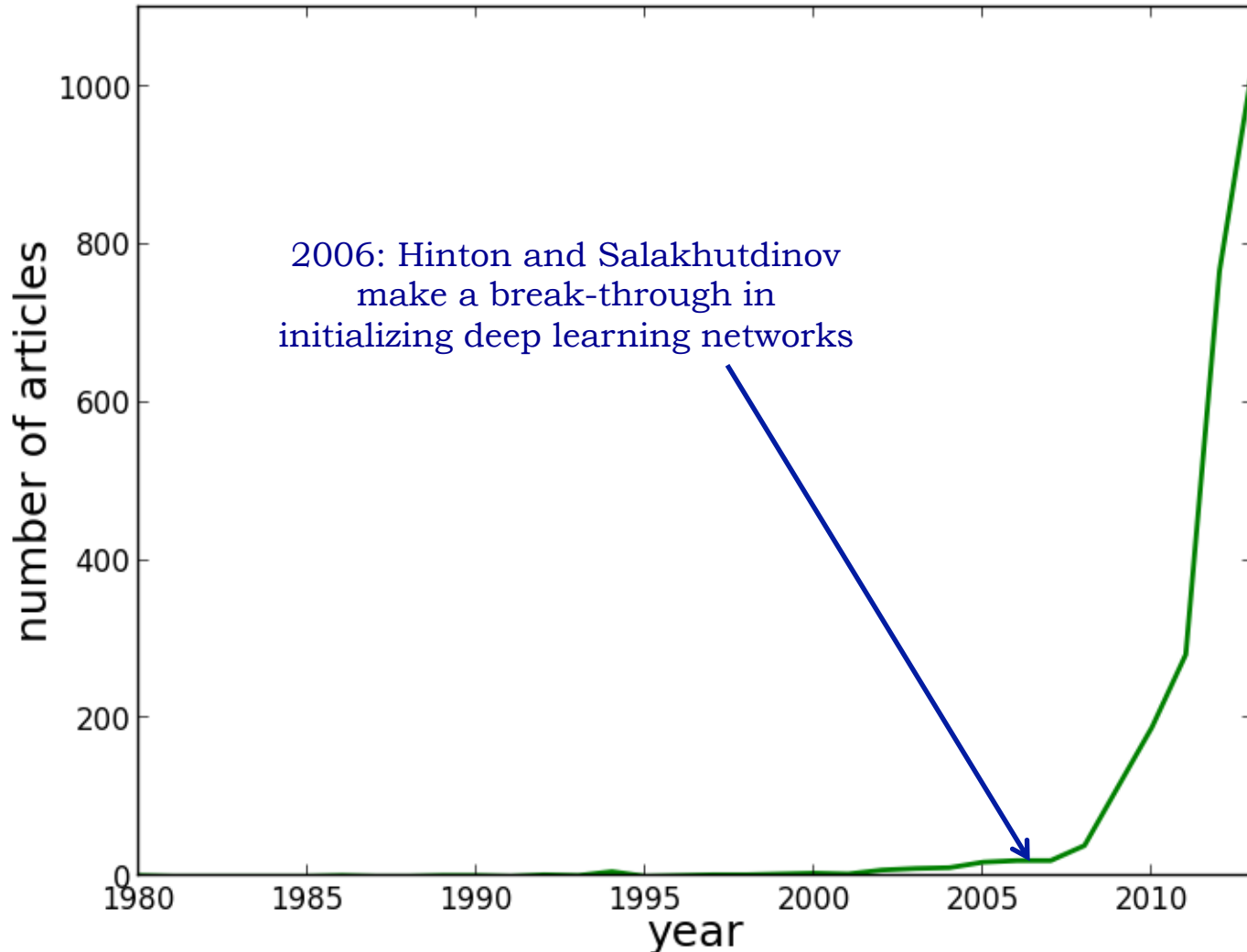
# History of Neural Networks

- Perceptron can be interpreted as a simple neural network

- Misconceptions about the weaknesses of perceptrons contributed to declining funding for NN research

- Difficulty of training multi-layer NNs contributed to second setback

- Mid 2000's: breakthroughs in NN training contribute to rise of "deep learning"

# Number of papers that mention "deep learning" over time

# Big picture for today

- Neural networks can approximate any function!

# Big picture for today

- Neural networks can approximate any function!

- For our purposes in ML, we want to use them to approximate a function from our inputs to our outputs

# Big picture for today

- Neural networks can approximate any function!

- For our purposes in ML, we want to use them to approximate a function from our inputs to our outputs

- We will train our network by asking it to minimize the loss between its output and the true output

# Big picture for today

- Neural networks can approximate any function!

- For our purposes in ML, we want to use them to approximate a function from our inputs to our outputs

- We will train our network by asking it to minimize the loss between its output and the true output

- We will use SGD-like approaches to minimize loss