# CS 260: Foundations of Data Science

Prof. Sara Mathieson

Fall 2023

HAVERFORD
COLLEGE

# Admin

- **Midterm 1** due TODAY

- **Lab 5** due Wednesday after fall break
  - Naïve Bayes

- **Lab 3** grades posted

# Outline for Oct 12

- Intro to Algorithmic Bias

- Disparate Impact

- Handout 11/12, clinical example

- Naïve Bayes implementation

- Handout 12, tennis example

# Outline for Oct 12

- Intro to Algorithmic Bias


- Disparate Impact


- Handout 11/12, clinical example


- Naïve Bayes implementation


- Handout 12, tennis example

# What does it mean to claim that algorithms are biased (or racist or political...)?

```
3   model = initialization(...)
4   n_epochs = ...
5   train_data = ...
6   for i in n_epochs:
7       train_data = shuffle(train_data)
8       X, y = split(train_data)
9       predictions = predict(X, model)
        error = calculate_error(y, predictions)
        model = update_model(model, error)
```

*Pseudocode from [A Gentle Introduction to Mini-Batch Gradient Descent and How to Configure Batch Size](#)*
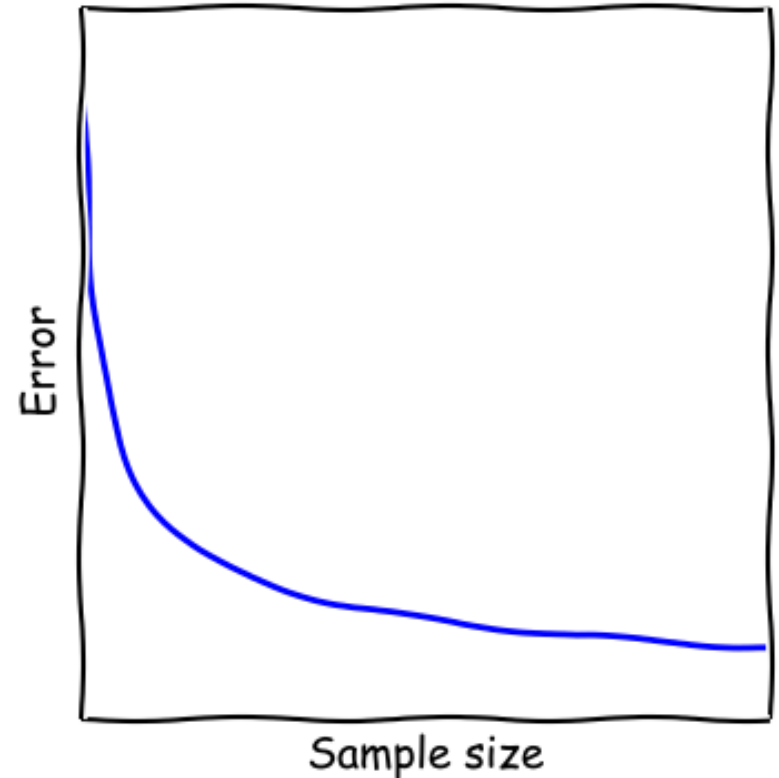
# Are algorithms fair by default?

"After all, as the former CPD [Chicago Police Department] computer experts point out, the algorithms in themselves are neutral. 'This program had absolutely nothing to do with race... but multi-variable equations,' argues Goldstein. Meanwhile, the potential benefits of predictive policing are profound."

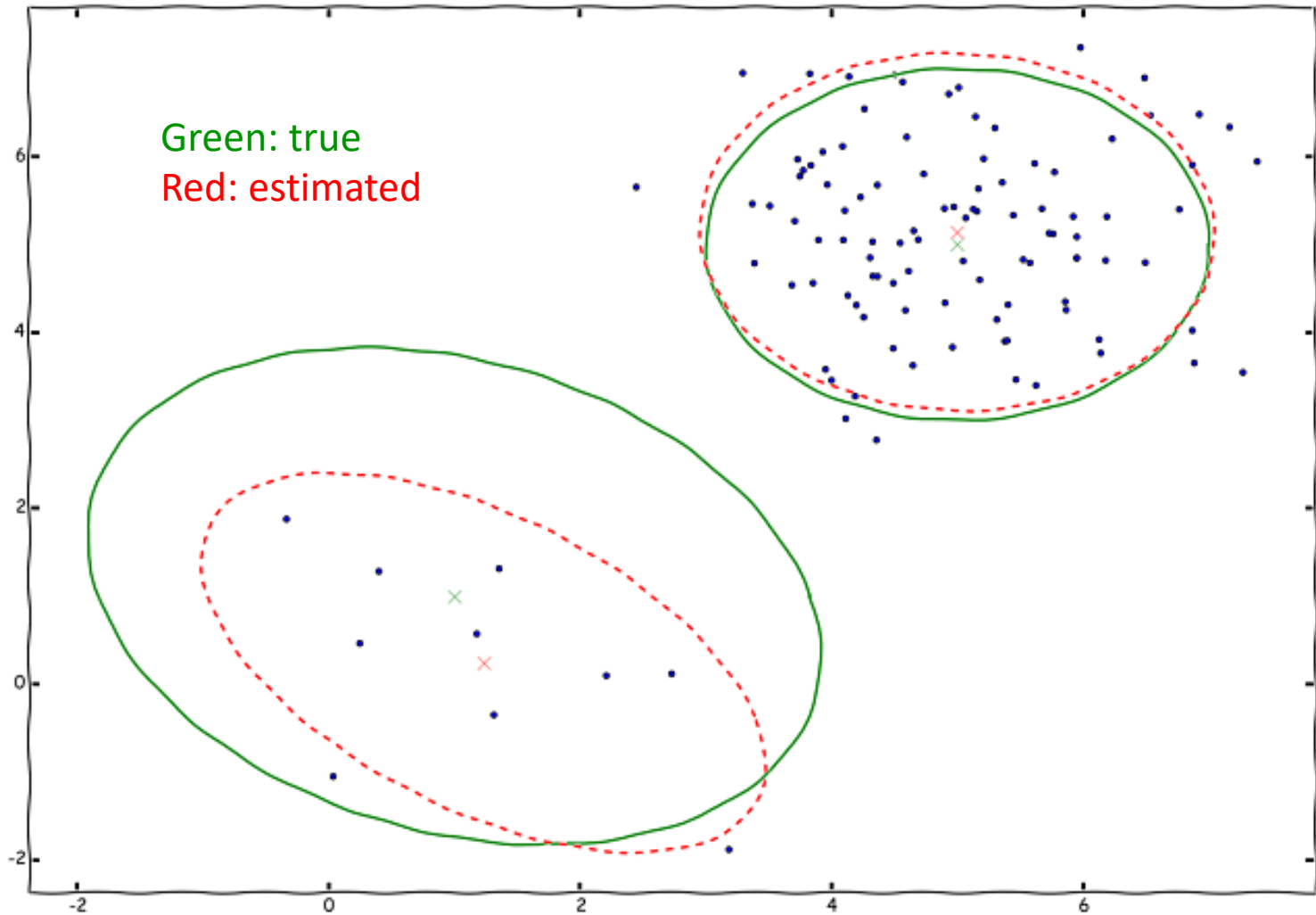-Gilian Tett

# Sample size disparity

- More data from majority will make results more accurate for that group

- Less accurate for the minority



"The error of a classifier often decreases as the inverse square root of the sample size. Four times as many samples means halving the error rate."
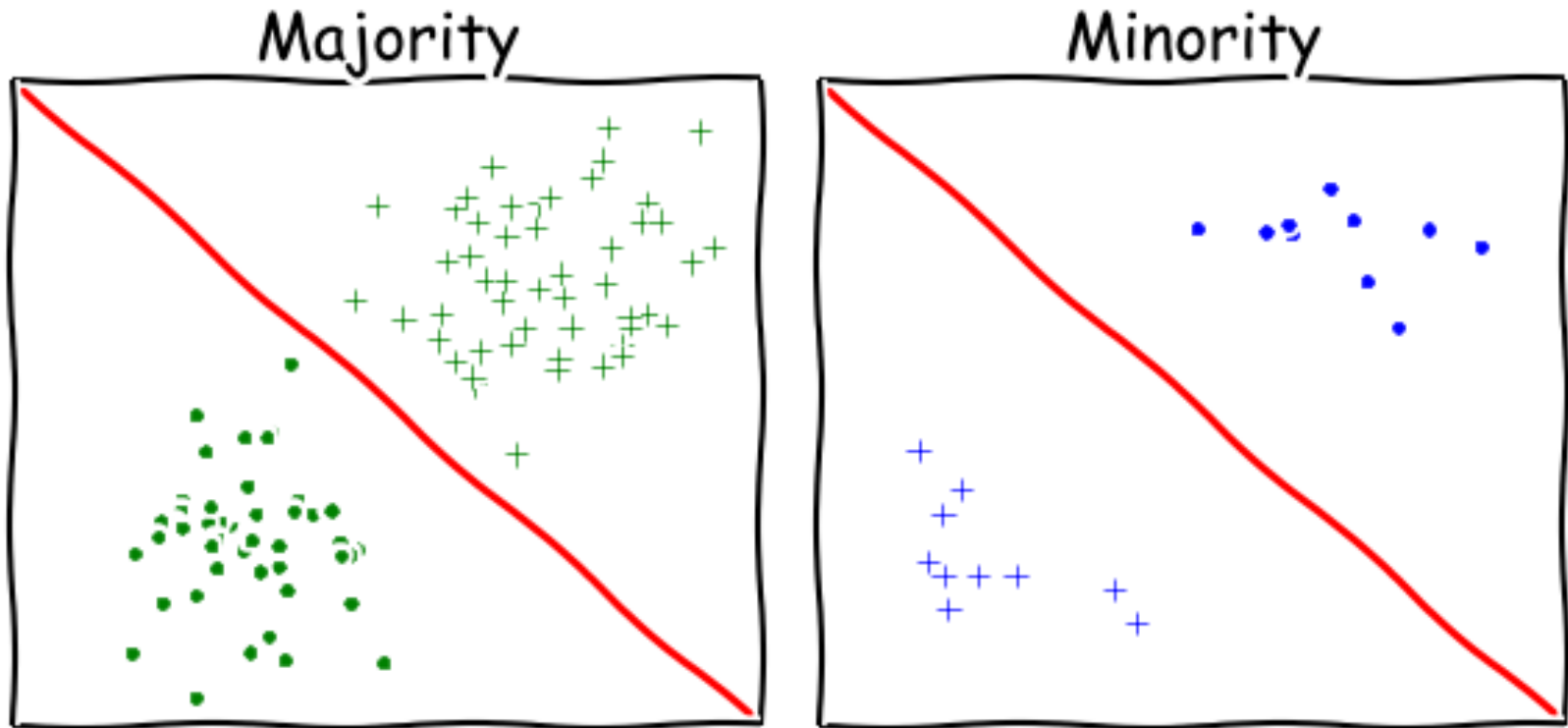Image: Moritz Hardt

# Sample size disparity



Green: true
Red: estimated

"Modeling a heterogeneous population as a gaussian mixture and learning its parameters using the EM algorithm. As expected, the estimates for the smaller group are significantly worse than for the larger. Dashed red ellipsoids describe the estimated covariance matrices. Solid green defines the correct covariance matrices. The green and red crosses indicate correct and estimated means, respectively." Image: Moritz Hardt
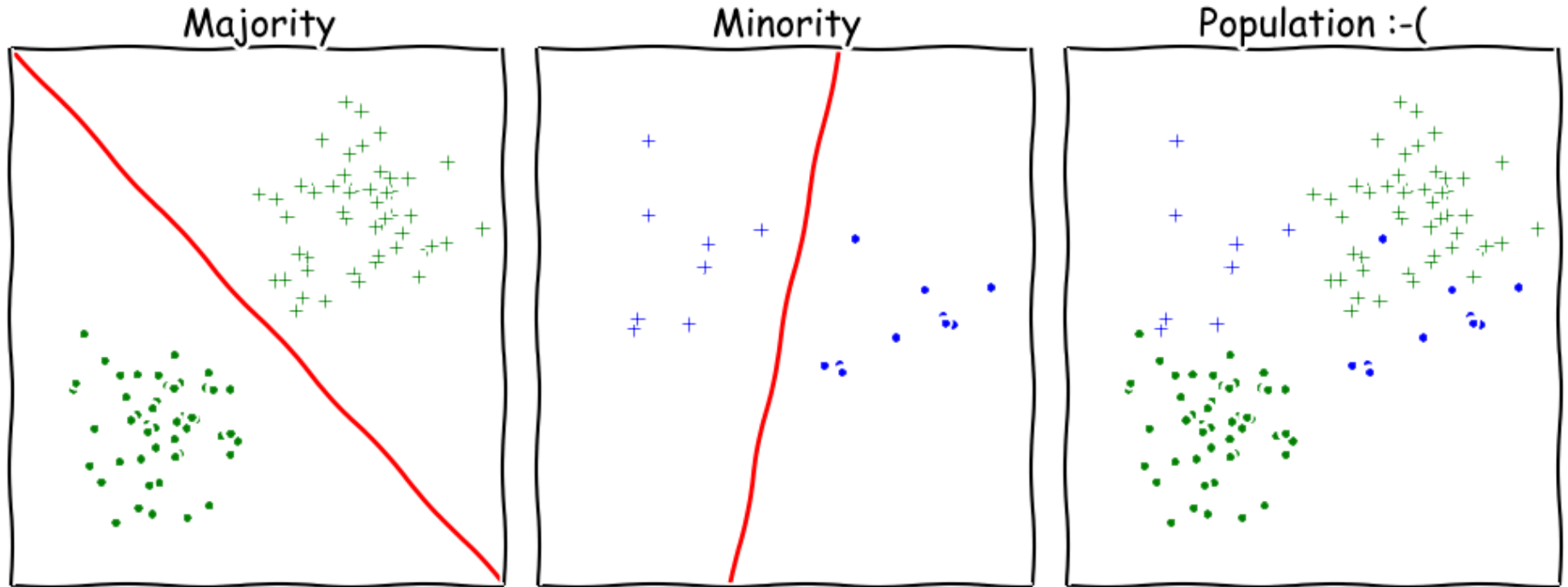
# Cultural Differences



"Positively labeled examples are on opposite sides of the classifier for the two groups." Image: Moritz Hardt

Goal: determine if a user profile (on Facebook, Twitter, etc) is genuine
- positive: real profile
- negative: fake profile

Feature: length of name

# Undesired Complexity



"Even if two groups of the population admit simple classifiers, the whole population may not."
Image: Moritz Hardt

# "How big data is unfair" (takeaways)



- ML is not fair by default, even though it relies on "neutral" multi-variable equations

- If training data reflects social biases, algorithm will likely incorporate them

- "Protected" attributes (race, gender, religion, sexual orientation, etc) often redundantly encoded

# Example: machine translation



Turkish - detected ▾

o bir aşçı
o bir mühendis
o bir doktor
o bir hemşire
o bir temizlikçi
o bir polis
o bir asker
o bir öğretmen

English ▾

# Example: machine translation



| Turkish - detected | English |
|---|---|
| o bir aşçı | she is a cook |
| o bir mühendis | he is an engineer |
| o bir doktor | he is a doctor |
| o bir hemşire | she is a nurse |
| o bir temizlikçi | he is a cleaner |
| o bir polis | He-she is a police |
| o bir asker | he is a soldier |
| o bir öğretmen | She's a teacher |

Slide: Ameet Soni

# Challenges

Algorithms do not exist in a bubble

- Inherit the prejudices of their designers
- Reflect cultural biases
- Difficult to identify - can entrench/enhance issues
- Deny historically disadvantaged groups full participation

# Outline for Oct 12

- Intro to Algorithmic Bias

- Disparate Impact

- Handout 11/12, clinical example

- Naïve Bayes implementation

- Handout 12, tennis example

# 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)

# 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

$X$: protected attribute

features $\begin{cases} \boxed{\begin{array}{l} 0 \quad \text{minority group} \\ 1 \quad \text{majority group} \end{array}} \\ Y: \text{other} \\ \\ C: \text{binary outcome} \in \{0, 1\} \end{cases}$

$\underset{\text{not hired}}{\swarrow} \quad \underset{\text{hired}}{\swarrow}$

## Disparate Impact

(legal definition)

$\underline{\text{if}}$

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

$\Rightarrow$ disparate impact

$\underline{ex}$ 40% of women hired

30% of men hired

$0.4 \not\leq 0.8(0.3) \Rightarrow \boxed{no}$

$\underset{0.24}{\underbrace{\qquad}}$ ?

$0.4 < 0.8(0.6) \Rightarrow \boxed{yes}$

## Naive Bayes

Idea $\Rightarrow$ if we can predict $X$ from $Y$, could be disparate impact.

$\rightsquigarrow$ predictor $f: Y \rightarrow X$

Balanced error rate    BER

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

$\bigstar$ ① train classifier $f(y) \Rightarrow X$

② if BER is low could be disparate impact

Want high!

indicates confusion

# Outline for Oct 12

- Intro to Algorithmic Bias

- Disparate Impact

- Handout 11/12, clinical example

  Handout 11 (#2), Handout 12 (#1)

- Naïve Bayes implementation

- Handout 12, tennis example

$f_1$ $f_2$ | $y$

| $f_1$ | $f_2$ | $y$ |
|---|---|---|
| p | n | 1 |
| p | p | 2 |
| p | n | 2 |
| n | n | 1 |
| p | n | 2 |
| n | n | 2 |
| n | p | 2 |

$X$

likelihood

$\big\{$

$P(\vec{x}\,|\,y=1)$

feature values

| $y=1$ | p | n |
|---|---|---|
| $f_1$ | $\frac{1+1}{3+2}$ | $\frac{2+1}{3+2}$ |
| $f_2$ | $\frac{0+1}{3+2}$ | $\frac{3+1}{3+2}$ |

$=\frac{M}{5}$

$\frac{}{5}$

Prior

values

| $y=2$ | p | n |
|---|---|---|
| $f_1$ | $\frac{4}{6}$ | $\frac{2}{6}$ |
| $f_2$ | $\frac{3}{6}$ | $\frac{3}{6}$ |

$\theta_1 = \frac{3+1}{7+2} = \frac{4}{9}$

$\theta_2 = \frac{4+1}{7+2} = \frac{5}{9}$

$\big\}$ add to 1

Handout 12

$$\vec{x} = [\underset{f_1}{neg}, \underset{f_2}{pos}]$$

① $p(y=1 \mid \vec{x}) \approx p(y=1)\, p(f_1=neg \mid y=1)\, p(f_2=pos \mid y=1)$

(healthy)

$$\approx \frac{4}{9} \cdot \frac{3}{5} \cdot \frac{1}{5} = \frac{4}{75}$$

$p(y=2 \mid \vec{x}) \approx \frac{5}{9} \cdot \frac{1}{3} \cdot \frac{1}{2} = \frac{5}{54}$

(disease)

$$\left( \frac{1}{3}, \frac{1}{5}, \underset{\bigcirc}{\frac{1}{2}} \right)$$

$\begin{matrix} 0 & 1 & ② \end{matrix}$ ← argmax

$p(y=k)$

max? $\boxed{\dfrac{5}{54}} \Rightarrow \boxed{\hat{y}=2}$

↑ argmax

$$\left[ \frac{4}{75}, \frac{5}{54} \right] \Rightarrow \left[ 37\%, 63\% \right]$$

$\frac{4}{75} + \frac{5}{54}$

☆

# Outline for Oct 12

- Intro to Algorithmic Bias

- Disparate Impact

- Handout 11/12, clinical example

- Naïve Bayes implementation

- Handout 12, tennis example

$$\prod_{j=1}^{P} P(x_j = v \mid y = k) \approx P(y = k \mid \bar{x})$$

$$P(A, B) = P(B) P(A \mid B)$$

$$\cancel{P(\bar{x})}$$

$$\log\left(\frac{a}{b}\right) = \log a - \log b$$
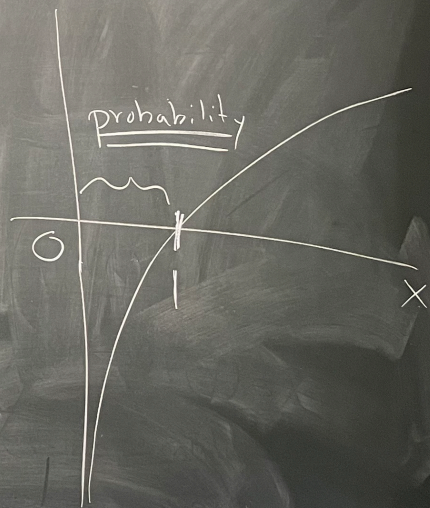
<u>issue is underflow!</u>

$$\frac{1}{1000} \cdot \frac{1}{10} \cdot \frac{1}{10} \cdot \frac{1}{10} \cdots \approx 0$$

$$\log(\Theta_k) = \log\left(\frac{N_k + 1}{n + K}\right) = \log(N_k + 1) - \log(n + K)$$



probability

# Outline for Oct 12

- Intro to Algorithmic Bias

- Disparate Impact

- Handout 11/12, clinical example

- Naïve Bayes implementation

- Handout 12, tennis example

# Data Structure idea
## (tennis example)

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis $(y)$ |
|---|---|---|---|---|---|
| $x_1$ | Sunny | Hot | High | Weak | No |
| $x_2$ | Sunny | Hot | High | Strong | No |
| $x_3$ | Overcast | Hot | High | Weak | Yes |
| $x_4$ | Rain | Mild | High | Weak | Yes |
| $x_5$ | Rain | Cool | Normal | Weak | Yes |
| $x_6$ | Rain | Cool | Normal | Strong | No |
| $x_7$ | Overcast | Cool | Normal | Strong | Yes |
| $x_8$ | Sunny | Mild | High | Weak | No |
| $x_9$ | Sunny | Cool | Normal | Weak | Yes |
| $x_{10}$ | Rain | Mild | Normal | Weak | Yes |
| $x_{11}$ | Sunny | Mild | Normal | Strong | Yes |
| $x_{12}$ | Overcast | Mild | High | Strong | Yes |
| $x_{13}$ | Overcast | Hot | Normal | Weak | Yes |
| $x_{14}$ | Rain | Mild | High | Strong | No |

# Data Structure idea

## (tennis example)

**Condition on y=No**

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis ($y$) |
|---|---|---|---|---|---|
| $x_1$ | Sunny | Hot | High | Weak | No |
| $x_2$ | Sunny | Hot | High | Strong | No |
| $x_3$ | Overcast | Hot | High | Weak | Yes |
| $x_4$ | Rain | Mild | High | Weak | Yes |
| $x_5$ | Rain | Cool | Normal | Weak | Yes |
| $x_6$ | Rain | Cool | Normal | Strong | No |
| $x_7$ | Overcast | Cool | Normal | Strong | Yes |
| $x_8$ | Sunny | Mild | High | Weak | No |
| $x_9$ | Sunny | Cool | Normal | Weak | Yes |
| $x_{10}$ | Rain | Mild | Normal | Weak | Yes |
| $x_{11}$ | Sunny | Mild | Normal | Strong | Yes |
| $x_{12}$ | Overcast | Mild | High | Strong | Yes |
| $x_{13}$ | Overcast | Hot | Normal | Weak | Yes |
| $x_{14}$ | Rain | Mild | High | Strong | No |

# Data Structure idea
## (tennis example)

**y=No (0)**

| outlook | **Sunny:** | **Overcast:** | **Rain:** |
|---|---|---|---|

| temperature | **Cool:** | **Mild:** | **Hot:** |
|---|---|---|---|

| humidity | **Normal:** | **High:** |
|---|---|---|

| wind | **Weak:** | **Strong:** |
|---|---|---|

**y=Yes (1)**

| outlook | **Sunny:** | **Overcast:** | **Rain:** |
|---|---|---|---|

| temperature | **Cool:** | **Mild:** | **Hot:** |
|---|---|---|---|

| humidity | **Normal:** | **High:** |
|---|---|---|

| wind | **Weak:** | **Strong:** |
|---|---|---|

# Discussion: admissions at Haverford

- Haverford has suddenly started receiving 10x more applications than usual

- You are tasked with creating an algorithm to determine whether or not an applicant should be admitted

- Questions:
  - How would you encode features?
  - How would you use past admission data to train?
  - What loss function are you trying to optimize?