# CS 260: Foundations of Data Science

Prof. Sara Mathieson

Fall 2023

HAVERFORD
COLLEGE

# Admin

- **Midterm 1**
  - due at the *beginning* of class on Thursday
- **Lab today:** continue midterm review
  - Attendance optional but recommended!
- **Zoom appointment slots** tomorrow (Wed)
  - Email me *today* to arrange
- **Lab 3 grades** should be up tomorrow
- **Lab 5** released tomorrow
  - due Wednesday after fall break

# Midterm 1 Notes

- Timed exam: **3 hour limit**. DO NOT open the exam until you are ready to take it for 3 hours!

- You may use a one page (front and back) "study sheet", handwritten, created by you

- Outside of your "study sheet" and calculator, **no other notes or resources**

- As per the Honor Code, all work must be your own

# Feedback forms (thank you!)

General workload/difficult

1 x

2 x

3 xxxxxxxxxxxx

4 xxxxxxxxx

5

# Feedback forms (thank you!)

- Different office hours times, zoom?
  - Will try to shift a bit later on Monday and add a Friday zoom when I can
- Diversity of opinions about lab deadlines
  - Will generally try to keep earlier in the week
- More prep for the labs in class
  - For Naïve Bayes we'll discuss implementation
- Half and half more slides vs. more board
- Lab instructions could be clearer

# Outline for October 10

- Recap Bayesian models

- Naïve Bayes algorithm

- Thurs: algorithmic bias

# Outline for October 10

- Recap Bayesian models


- Naïve Bayes algorithm


- Thurs: algorithmic bias

# Informal Quiz (discuss with a partner)

1. How would you say $P(A, B)$ in words?

2. Based on class on Tuesday, what is Bayes rule?

$$P(A, B) =$$

4. If I want to predict the label $(y)$ of an example based on its features $(\vec{x})$, which of the following expressions would I want to compute? (circle the best one)

   (a) $p(\vec{x}, y)$
   (b) $p(\vec{x} \mid y)$
   (c) $p(y \mid \vec{x})$

# Informal Quiz (discuss with a partner)

1. How would you say $P(A, B)$ in words?     Probability of A and B

2. Based on class on Tuesday, what is Bayes rule?

$$P(A, B) = \quad \text{P(A) P(B|A)} \qquad \text{or} \qquad \text{P(B) P(A|B)}$$

4. If I want to predict the label $(y)$ of an example based on its features $(\vec{x})$, which of the following expressions would I want to compute? (circle the best one)

    (a)  $p(\vec{x}, y)$
    (b)  $p(\vec{x} \mid y)$
    (c)  $p(y \mid \vec{x})$

# Components of a Bayesian Model

- Identify the evidence, prior, posterior, and likelihood in the equation below

$$p(y = k|\boldsymbol{x}) = \frac{p(y = k)p(\boldsymbol{x}|y = k)}{p(\boldsymbol{x})}$$

# Components of a Bayesian Model

- Identify the evidence, prior, posterior, and likelihood in the equation below

$$p(y = k|\boldsymbol{x}) = \frac{p(y = k)p(\boldsymbol{x}|y = k)}{p(\boldsymbol{x})}$$

- **Evidence**: this is the data (features) we actually observe, which we think will help us predict the outcome we're interested in

# Components of a Bayesian Model

- Identify the evidence, prior, posterior, and likelihood in the equation below

$$p(y = k|\boldsymbol{x}) = \frac{p(y = k)p(\boldsymbol{x}|y = k)}{p(\boldsymbol{x})}$$

- **Prior**: without seeing any evidence (data), what is our prior believe about each outcome (intuition: what is the outcome in the population as a whole?)

# Components of a Bayesian Model

- Identify the evidence, prior, posterior, and likelihood in the equation below

$$p(y = k|\boldsymbol{x}) = \frac{p(y = k)p(\boldsymbol{x}|y = k)}{p(\boldsymbol{x})}$$

- **Posterior**: this is the quantity we are actually interested in. *Given* the evidence, what is the probability of the outcome?

# Components of a Bayesian Model

- Identify the evidence, prior, posterior, and likelihood in the equation below

$$p(y = k|\boldsymbol{x}) = \frac{p(y = k)p(\boldsymbol{x}|y = k)}{p(\boldsymbol{x})}$$

- **Likelihood**: given an outcome, what is the probability of observing this set of features?

# Outline for October 10

- Recap Bayesian models

- Naïve Bayes algorithm

- Thurs: algorithmic bias

# Real-world example of Naïve Bayes

"A Comparison of Event Models for Naive Bayes Text Classification" (5649 citations!)

http://www.cs.cmu.edu/~knigam/papers/multinomial-aaaiws98.pdf

Goal: text classification (classify documents into topics based on the words as features)

# Naive Bayes

Single example    $\vec{x} = [x_1, x_2 \ldots x_p]^T$    (word counts)

multi-class response    $y \in \{1, 2, \ldots, K\}$

$\boxed{\text{code: } 0, 1, \ldots, K-1}$

· goal : multi-class classification (ie. $\hat{y}$)

Bayesian model

$$\underbrace{p(y=k \mid \vec{x})}_{\text{posterior}} = \frac{\overbrace{p(y=k)}^{\text{prior}} \; \overbrace{p(\vec{x} \mid y=k)}}{\underbrace{p(\vec{x})}_{\substack{\text{ignore?} \\ \text{same for all } k}}}$$

prediction

$$\boxed{\hat{y} = \underset{k=1, \ldots, K}{\text{argmax}} \; p(y=k \mid \vec{x})}$$

$$\star \; p(\vec{x} \mid y=k) = p(x_1, x_2, x_3 \cdots x_p \mid y=k)$$

$$\underbrace{\phantom{x_1}}_{A} \quad \underbrace{\phantom{x_2 x_3 \cdots x_p}}_{B}$$

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

$$= p(\underbrace{x_2, x_3 \cdots x_p}_{B} \mid y=k) \, p(\underbrace{x_1}_{A} \mid \underbrace{x_2 \cdots x_p, y=k}_{B})$$

$$\underbrace{\phantom{x_2 x_3}}_{C} \quad \underbrace{\phantom{\cdots}}_{D}$$

$$= p(x_3, \cdots x_p \mid y=k) \, p(\boxed{x_2} \mid \cancel{x_3 \cdots} x_p, y=k) \, p(x_1 \mid x_2 \cdots x_p, y=k)$$

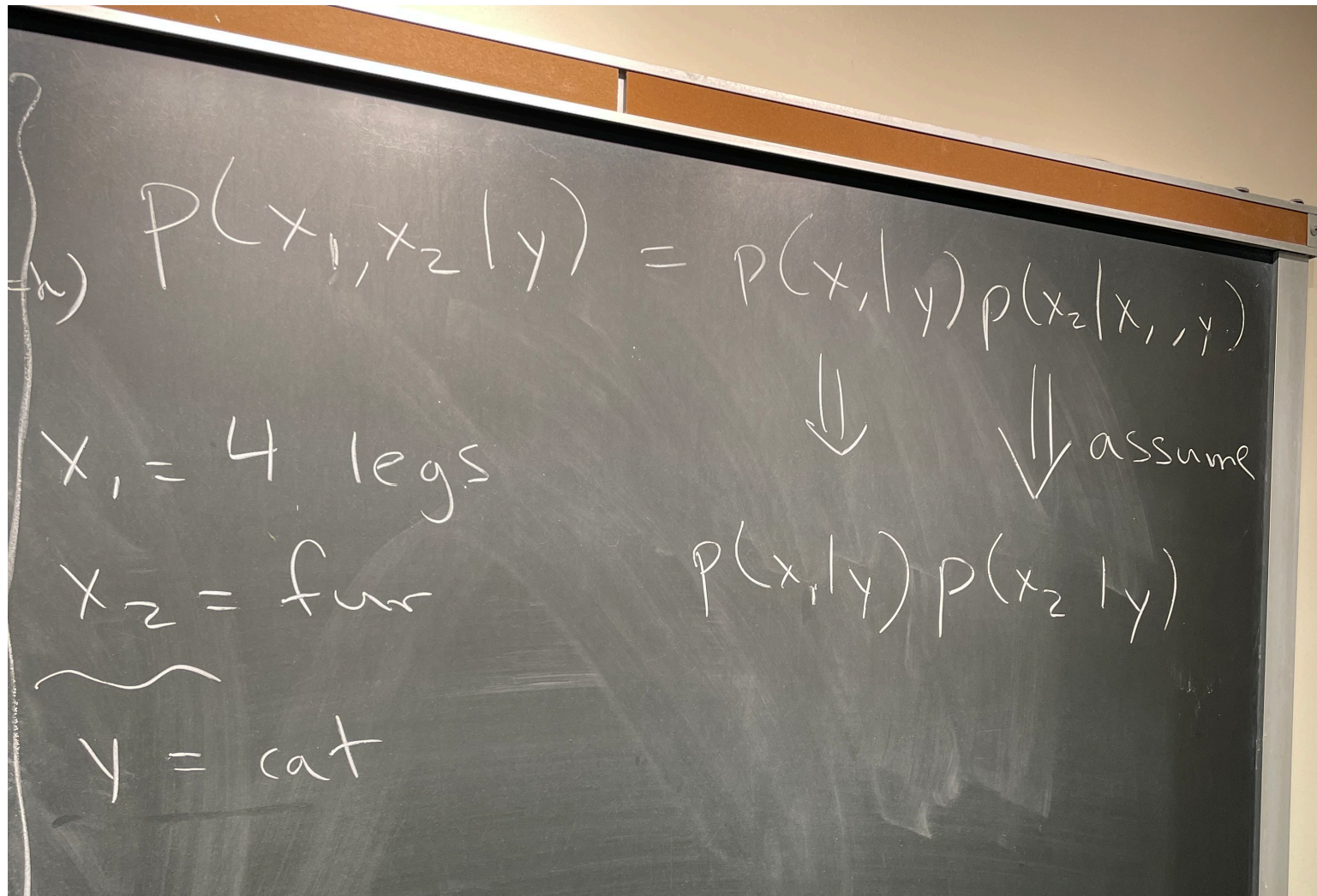$$\underbrace{\phantom{x_1 \cdots x_p}}_{\text{assumption} \rightarrow \; p(x_1 \mid y=k)}$$

conditional independence assumption
(Naive Bayes)

"feature $j$ is independent from all other features <u>given</u> <u>label</u> $k$"

# Conditional independence example

$$P(x_1, x_2 \mid y) = P(x_1 \mid y) \, p(x_2 \mid x_1, y)$$

$$\Downarrow \qquad \Downarrow \text{ assume}$$

$$p(x_1 \mid y) \, p(x_2 \mid y)$$

$x_1 = 4 \text{ legs}$

$x_2 = \text{fur}$

$y = \text{cat}$

$P(B) \, P(A|B)$

$\rightarrow p(\vec{x} \, | \, y=k) = p(x_p | y=k) \, p(x_{p-1} | y=k) \cdots p(x_2 | y=k) \, p(x_1 | y=k)$

product (like $\sum$ for sum)

$$= \prod_{j=1}^{p} p(x_j | y=k)$$

Naive Bayes Model

"proportional to"

$p(y=k | \vec{x}) \, \propto \, p(y=k) \prod_{j=1}^{p} p(x_j | y=k)$

$p(\vec{x})$

pred $\Rightarrow$ $\hat{y} = \text{argmax}$

$p(x_1, \ldots$

$x_1 = 4 \, \text{le}$

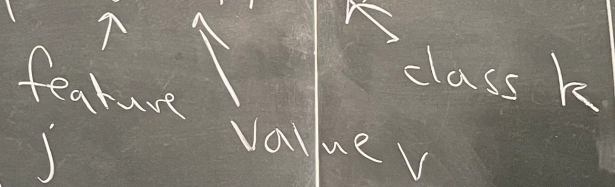$x_2 = \text{fu}$

$\sim$

$y = \text{cat}$

What are $p(y=k)$ & $p(x_j | y=k)$?

estimate based on training data!

$\Theta_k$ = estimate for $p(y=k)$

$\Theta_{k,j,v}$ = estimate for $p(x_j = v | y = k)$

feature $j$ &uarr; value $v$ &uarr; class $k$

ex

$$\boxed{p\left(x_{outside} = sun \mid tennis = yes\right)}$$

let

let $N_k = \#$ examples with label $k$

$$\Theta_k = \frac{N_k}{n}$$

$$\Theta_1 = \frac{875}{1000}$$

$$\Theta_2 = \frac{125}{1000}$$

add to 1
(probability distribution)

$1 =$ healthy
$2 =$ disease

$N_1 = 875$

$N_2 = 125$

$n = 1000$

$$\Theta_k = \frac{N_k + 1}{n + K}$$ ← LaPlace count

implementation

$$\sum \Theta_k = \sum \frac{N_k + 1}{n + K} = \frac{1}{n+K}(n+K)$$

let $N_{kj \cdot v} = \#$ of examples with feature $j =$ value $v$ & class label $k$

$$\theta_{k,j,v} = \frac{N_{k,j,v} + 1}{N_k + |f_j|}$$

$\nearrow$ # of features
values for
feature $j$

$N(\text{tennis} = \text{yes}) = 7$

$N_{\text{yes, outlook, sun}} = 4$

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

Naive Bayes Model

proportional to"

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

$p(\bar{x})$

pred $\Rightarrow$  $\hat{y} = \text{argmax}$

# Handout 11

Say we have two tests for a specific disease. Each test (features $f_1$, $f_2$) can come back either positive "pos" or negative "neg", and the true underlying condition of the patient is represented by $y$ ($y = 1$ is "healthy" and $y = 2$ is "disease"). We observe this training data where $n = 7$ and $p = 2$:

| $\boldsymbol{x}$ | $f_1$ | $f_2$ | $y$ |
|---|---|---|---|
| $\boldsymbol{x}_1$ | pos | neg | 1 |
| $\boldsymbol{x}_2$ | pos | pos | 2 |
| $\boldsymbol{x}_3$ | pos | neg | 2 |
| $\boldsymbol{x}_4$ | neg | neg | 1 |
| $\boldsymbol{x}_5$ | pos | neg | 2 |
| $\boldsymbol{x}_6$ | neg | neg | 1 |
| $\boldsymbol{x}_7$ | neg | pos | 2 |

1. To estimate the probability $p(y = k)$, for $k = 1, 2, \cdots, K$, we will use the formula:

$$\theta_k = \frac{N_k + 1}{n + K}$$

where $N_k$ is the count ("Number") of data points where $y = k$. Compute $\theta_1$ and $\theta_2$. What would $\theta_1$ and $\theta_2$ be if we in fact had *no* training data?

| $\bar{X}$ | $f_1$ $f_2$ | $Y$ |
|---|---|---|
| $\bar{X}_1$ | pos neg | 1 |
| $\bar{X}_2$ | pos pos | 2 |
| | pos neg | 2 |
| $\bar{X}_3$ | neg neg | 1 |
| $\bar{X}_4$ | | 2 |
| $\bar{X}_5$ | pos neg | 2 |
| $\bar{X}_6$ | neg neg | 1 |
| $\bar{X}_7$ | neg pos | 2 |

$$\Theta_1 = \frac{3+1}{7+2}$$

$$\frac{4}{9}$$

$$\Theta_2 = \frac{5}{9}$$

| $Y=1$ | pos | neg |
|---|---|---|
| $f_1$ | $\frac{1+1}{3+2}$ | |
| $f_2$ | | |

| $Y=2$ | pos | neg |
|---|---|---|
| $f_1$ | | |
| $f_2$ | | |

# Handout 11

2. To estimate the probabilities $p(x_j = v | y = k)$ for all features $j$, values $v$, and class label $k$, we will use the formula:

$$\theta_{k,j,v} = \frac{N_{k,j,v} + 1}{N_k + |f_j|}$$

where $N_{k,j,v}$ is the count of data points where $y = k$ and $x_j = v$, and $|f_j|$ is the number of possible values that $f_j$ (feature $j$) can take on. Fill in the following tables with these $\theta$ values.

| $y = 1$ | pos | neg |
|---------|-----|-----|
| $f_1$ |  |  |
| $f_2$ |  |  |

| $y = 2$ | pos | neg |
|---------|-----|-----|
| $f_1$ |  |  |
| $f_2$ |  |  |