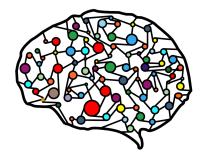# Lecture 24 – Naive Bayes

**DSC 40A, Fall 2022 @ UC San Diego**
Mahdi Soleymani, with help from **many others**

# Agenda

- ▶ Naive Bayes.

- ▶ Naive Bayes in practice — text classification.

- ▶ Practical demo.

# Naive Bayes

# Naive Bayes classifier

▶ We want to predict a class, given certain features.

▶ Using Bayes' theorem, we write

$$P(f_1, f_2 | class) = P(f_1 | class) \cdot P(f_2 | class)$$

$$P(class|features) = \frac{P(class) \cdot \overbrace{P(features|class)}}{P(features)}$$

▶ For each class, we compute the numerator using the **naive assumption of conditional independence of features given the class**.

▶ We estimate each term in the numerator based on the training data.

▶ We predict the class with the largest numerator.
  ▶ Works if we have multiple classes, too!

# na·ive

/nīˈēv/

*adjective*

- (of a person or action) showing a lack of experience, wisdom, or judgment.
  "the rather naive young man had been totally misled"
- (of a person) natural and unaffected; innocent.
  "Andy had a sweet, naive look when he smiled"

  Similar: innocent | unsophisticated | artless | ingenuous | inexperienced | ⌄

- of or denoting art produced in a straightforward style that deliberately rejects sophisticated artistic techniques and has a bold directness resembling a child's work, typically in bright colors with little or no perspective.

# Example: comic characters

*features* — male, Marvel; bad good neutral

*Label*

| ALIGN | SEX | COMPANY |
|---|---|---|
| Bad | Male | Marvel |
| Neutral | Male | Marvel |
| Good | Male | Marvel |
| Bad | Male | DC |
| Good | Female | Marvel |
| Bad | Male | DC |
| Good | Male | DC |
| Bad | Male | Marvel |
| Good | Female | Marvel |
| Bad | Female | Marvel |

My favorite character is a male Marvel character. Using Naive Bayes, would we predict that my favorite character is bad, good, or neutral?

| ALIGN | SEX | COMPANY |
|---|---|---|
| Bad | Male | Marvel |
| Neutral | Male | Marvel |
| Good | Male | Marvel |
| Bad | Male | DC |
| Good | Female | Marvel |
| Bad | Male | DC |
| Good | Male | DC |
| Bad | Male | Marvel |
| Good | Female | Marvel |
| Bad | Female | Marvel |

$$P(\text{bad} \mid m, M) \propto$$

male → $m$   Marvel → $M$

$$\alpha \cdot P(\text{bad}) \cdot P(m, M \mid \text{bad})$$

$$= P(\text{bad}) \cdot P(m \mid \text{bad}) \, P(M \mid \text{bad})$$

$$= \frac{5}{10} \cdot \frac{4}{5} \cdot \frac{3}{5} = \frac{6}{25}$$

$$P(\text{good} \mid m, M) \propto P(\text{good}) \cdot P(m \mid \text{good}) \, P(M \mid \text{good})$$

$$= \frac{4}{10} \cdot \frac{2}{4} \cdot \frac{3}{4} = \frac{3}{20} = \frac{6}{40}$$

$$P(\text{neutral} \mid m, M) \propto P(\text{neutral}) \quad P(m \mid \text{neutral}) \, P(M \mid \text{neutral})$$

$$\frac{1}{10} \cdot \frac{1}{1} \cdot \frac{1}{1} = \frac{1}{10} = \frac{6}{60}$$

predict charcter is BAD!

# Example: comic characters

| ALIGN | SEX | COMPANY |
|---|---|---|
| Bad | Male | Marvel |
| Neutral | Male | Marvel |
| Good | Male | Marvel |
| Bad | Male | DC |
| Good | Female | Marvel |
| Bad | Male | DC |
| Good | Male | DC |
| Bad | Male | Marvel |
| Good | Female | Marvel |
| Bad | Female | Marvel |

My other favorite character is a female Marvel character. Using Naive Bayes, would we predict that my favorite character is bad, good, or neutral?

| ALIGN | SEX | COMPANY |
|-------|-----|---------|
| Bad | Male | Marvel |
| Neutral | Male | Marvel |
| Good | Male | Marvel |
| Bad | Male | DC |
| Good | Female | Marvel |
| Bad | Male | DC |
| Good | Male | DC |
| Bad | Male | Marvel |
| Good | Female | Marvel |
| Bad | Female | Marvel |

$$P(\text{neutral} \mid \text{female, Marvel})$$

$$\propto P(\text{neutral}) \cdot P(\text{female} \mid \text{neutral})$$

$$P(\text{Marvel} \mid \text{neutral})$$

$$= 0$$

$$\frac{0}{1}$$

# Uh oh...

▶ There are no neutral female characters in the data set.

▶ The estimate $P(\text{female}|\text{neutral}) \approx \frac{\text{\# female neutral characters}}{\text{\# neutral characters}}$ is 0.

▶ The estimated numerator,
$P(\text{neutral}) \cdot P(\text{female, Marvel}|\text{neutral}) =$
$P(\text{neutral}) \cdot P(\text{female}|\text{neutral}) \cdot P(\text{Marvel}|\text{neutral})$,
is also 0.

▶ But just because there isn't a neutral female character in the data set, doesn't mean they don't exist!

▶ **Idea:** Adjust the numerators and denominators of our estimate so that they're never 0.

# Smoothing

*→ Laplace smoothing*

$$\frac{+\alpha}{+K\alpha}$$
↓
# Classes

▸ **Without** smoothing:

$$P(\text{female}|\text{neutral}) \approx \frac{\text{\# female neutral}}{\text{\# female neutral} + \text{\# male neutral}}$$

$\alpha = 1$

$$P(\text{male}|\text{neutral}) \approx \frac{\text{\# male neutral}}{\text{\# female neutral} + \text{\# male neutral}}$$

▸ **With** smoothing:

*Smoothing for CONDITIONAL probs.*
↓

$$P(\text{female}|\text{neutral}) \approx \frac{\text{\# female neutral} + 1}{\text{\# female neutral} + 1 + \text{\# male neutral} + 1}$$

$$P(\text{male}|\text{neutral}) \approx \frac{\text{\# male neutral} + 1}{\text{\# female neutral} + 1 + \text{\# male neutral} + 1}$$

▸ When smoothing, we add 1 to the count of every group whenever we're estimating a probability.

# Example: comic characters

Using smoothing, let's determine whether Naive Bayes would predict a female Marvel character to be bad, good, or neutral.

| ALIGN | SEX | COMPANY |
|-------|-----|---------|
| Bad | Male | Marvel |
| Neutral | Male | Marvel |
| Good | Male | Marvel |
| Bad | Male | DC |
| Good | Female | Marvel |
| Bad | Male | DC |
| Good | Male | DC |
| Bad | Male | Marvel |
| Good | Female | Marvel |
| Bad | Female | Marvel |

$$P(\text{bad} \mid f, M) \propto P(\text{bad})$$

$$P(f \mid \text{bad}) \quad P(M \mid \text{bad})$$

$$\left(\frac{5}{10}\right) \cdot \left(\frac{1+1}{1+1+4+1}\right) \cdot \left(\frac{3+1}{3+1} \atop +2+1\right)$$

$$P(\text{good} \mid f, M) \propto P(\text{good}) \cdot P(f \mid \text{good}) \cdot P(M \mid \text{good})$$

$$\left(\frac{4}{10}\right) \cdot \left(\frac{2+1}{2+1+2+1}\right) \cdot \left(\frac{3+1}{3+1+1+1}\right) = \frac{5}{13} \cdot \frac{3}{6} \cdot \frac{4}{6}$$

$$P(\text{neutral} \mid f, M) \propto P(\text{neutral}) \cdot P(f \mid \text{neutral})$$

$$\cdot P(M \mid \text{neutral})$$

$$= \left( \frac{1}{10} \right) \left( \frac{0+1}{\underbrace{0+1}_{f} + \underbrace{1+1}_{m}} \right) \cdot \left( \frac{1+1}{\underbrace{1+1}_{M} + \underbrace{0+1}_{DC}} \right)$$

$$= \left( \frac{1}{10} \right) \left( \frac{1}{3} \right) \left( \frac{2}{3} \right)$$

# Recap: Naive Bayes classifier

▶ We want to predict a class, given certain features.

▶ Using Bayes' theorem, we write

$$P(\text{class}|\text{features}) = \frac{P(\text{class}) \cdot P(\text{features}|\text{class})}{P(\text{features})}$$

▶ For each class, we compute the numerator using the **naive assumption of conditional independence of features given the class**.

▶ We estimate each term in the numerator based on the training data.

▶ We predict the class with the largest numerator.
  ▶ Works if we have multiple classes, too!

CAPE's surveys : 3 extra credit

90% complete

**Text classification**

# Text classification

- Text classification problems include:
    - Sentiment analysis (e.g. positive and negative customer reviews).

    - Determining genre (news articles, blog posts, etc.).

    - **Spam filtering.**

- **Our goal:** given the body of an email, determine whether it's **spam** or **ham** (not spam).

**Question:** How do we come up with features?

# Features

**Idea:**

▶ Choose a **dictionary** of $d$ words, e.g. "prince", "money", "free"…

▶ Represent each email with a **feature vector** $\vec{x}$:

$$\vec{x} = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ ... \\ x^{(d)} \end{bmatrix}$$

where

▶ $x^{(i)}$ = 1 if word $i$ is present in the email, and
▶ $x^{(i)}$ = 0 otherwise.

This is called the **bag-of-words** model.

# Concrete example

▸ Dictionary: "prince", "money", "free", and "xxx".

▸ Dataset of 5 emails (red are spam, green are ham):
  ▸ "I am the prince of UCSD and I demand money."
  ▸ "Tapioca Express: redeem your free Thai Iced Tea!"
  ▸ "DSC 40A: free points if you fill out CAPEs!"
  ▸ "Click here to make a tax-free donation to the IRS."
  ▸ "Free COVID-19 tests at Price Center."

| | prince | money | free | x x x | class |
|---|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 | spam |
| 2 | 0 | 0 | 1 | 0 | ham |
| 3 | 0 | 0 | 1 | 0 | ham |
| 4 | 0 | 0 | 1 | 0 | spam |
| 5 | 0 | 0 | 1 | 0 | ham |

# Naive Bayes for spam classification

$$P(\text{class} \mid \text{features}) = \frac{P(\text{class}) \cdot P(\text{features} \mid \text{class})}{P(\text{features})}$$

▶ To classify an email, we'll use Bayes' theorem to calculate the probability of it belonging to each class:
  ▶ $P(\text{spam} \mid \text{features})$.
  ▶ $P(\text{ham} \mid \text{features})$.

▶ We'll predict the class with a larger probability.

# Naive Bayes for spam classification

$$P(\text{class} \mid \text{features}) = \frac{P(\text{class}) \cdot P(\text{features} \mid \text{class})}{P(\text{features})}$$

▶ Note that the formulas for $P(\text{spam} \mid \text{features})$ and $P(\text{ham} \mid \text{features})$ have the same denominator, $P(\text{features})$.

▶ Thus, we can find the larger probability just by comparing numerators:
  ▶ $P(\text{spam}) \cdot P(\text{features} \mid \text{spam})$.
  ▶ $P(\text{ham}) \cdot P(\text{features} \mid \text{ham})$.

# Naive Bayes for spam classification

*(handwritten)* $A$ : spam
$\bar{A}$ : ham

## Discussion Question

We need to determine four quantities:

1. $P(\text{features} \mid \text{spam})$.
2. $P(\text{features} \mid \text{ham})$.
3. $P(\text{spam})$.
4. $P(\text{ham})$.

*(handwritten)*
$$P(\text{spam} \mid \text{features}) + P(\text{ham} \mid \text{features}) = 1$$

Which of these probabilities should add to 1?

*(handwritten)* $P(A) + P(\bar{A}) = 1$

A) 1, 2

B) 3, 4

C) Both A and B

D) Neither A nor B

*(handwritten)*
$P(A \mid \text{ham})$
$= 0.1 + 0.2 = 0.3$

$5685 \quad 0753$

**To answer, go to** `menti.com` **and enter** ~~7055 7461~~

# Estimating probabilities with training data

▶ To estimate $P(\text{spam})$, we compute

$$P(\text{spam}) \approx \frac{\text{\# spam emails in training set}}{\text{\# emails in training set}}$$

▶ To estimate $P(\text{ham})$, we compute

$$P(\text{ham}) \approx \frac{\text{\# ham emails in training set}}{\text{\# emails in training set}}$$

▶ What about $P(\text{features} \mid \text{spam})$ and $P(\text{features} \mid \text{ham})$?

# Assumption of conditional independence

▶ Note that $P(\text{features} \mid \text{spam})$ looks like

$$P(x^{(1)} = 0, x^{(2)} = 1, ..., x^{(d)} = 0 \mid \text{spam})$$

word is not , word2 yes , ⋯ , word no

▶ Recall: the key assumption that the Naive Bayes classifier makes is that **the features are conditionally independent given the class**.

▶ This means we can estimate $P(\text{features} \mid \text{spam})$ as

$$P(x^{(1)} = 0, x^{(2)} = 1, ..., x^{(d)} = 0 \mid \text{spam})$$
$$= P(x^{(1)} = 0 \mid \text{spam}) \cdot P(x^{(2)} = 1 \mid \text{spam}) \cdot ... \cdot P(x^{(d)} = 0 \mid \text{spam})$$

# Concrete example

▶ Dictionary: "prince", "money", "free", and "xxx".

▶ Dataset of 5 emails (red are spam, green are ham):
  ▶ "I am the prince of UCSD and I demand money."
  ▶ "Tapioca Express: redeem your free Thai Iced Tea!"
  ▶ "DSC 40A: free points if you fill out CAPEs!"
  ▶ "Click here to make a tax-free donation to the IRS."
  ▶ "Free COVID-19 tests at Prince Center."

# Concrete example

▶ New email to classify: "Download a free copy of the Prince of Persia."

| | prince | money | free | x x x | class |
|---|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 | spam |
| 2 | 0 | 0 | 1 | 0 | ham |
| 3 | 0 | 0 | 1 | 0 | ham |
| 4 | 0 | 0 | 1 | 0 | spam |
| 5 | 1 | 0 | 1 | 0 | ham |

$$P(\text{spam} \mid \text{features}) \propto P(\text{spam}) \cdot P\left(\begin{array}{c} x_1 = 1 \\ \text{yes} \\ \text{prince} \end{array} \mid \text{spam}\right)$$

$$P\left(\begin{array}{c} \text{no} \\ \text{money} \end{array} \mid \text{spam}\right) P\left(\begin{array}{c} \text{yes} \\ \text{free} \end{array} \mid \text{spam}\right) P\left(\begin{array}{c} \text{no} \\ x x x \end{array} \mid \text{spam}\right)$$

$$= \left(\frac{2}{5}\right)\left(\frac{1}{2}\right)\left(\frac{1}{2}\right)\left(\frac{1}{2}\right)\left(\frac{2}{2}\right) = \frac{1}{20}$$

| | prince | money | free | xxx | class |
|---|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 | spam |
| 2 | 0 | 0 | 1 | 0 | ham |
| 3 | 0 | 0 | 1 | 0 | ham |
| 4 | 0 | 0 | 1 | 0 | spam |
| 5 | 1 | 0 | 1 | 0 | ham |

$$P(\text{ham} \mid \text{features}) \propto P(\text{ham}) \cdot P\left(\begin{array}{c} x_1 = 1 \\ \text{yes} \\ \text{prince} \end{array} \mid \text{ham}\right)$$

$$P\left(\begin{array}{c} \text{no} \\ \text{money} \end{array} \mid \text{ham}\right) P\left(\begin{array}{c} \text{yes} \\ \text{free} \end{array} \mid \text{ham}\right) P\left(\begin{array}{c} \text{no} \\ xxx \end{array} \mid \text{ham}\right)$$

$$= \left(\frac{3}{5}\right)\left(\frac{1}{3}\right)\left(\frac{3}{3}\right)\left(\frac{3}{3}\right)\left(\frac{3}{3}\right) = \frac{1}{5}$$

predict ham !

## Uh oh…

▶ What happens if we try to classify the email "xxx what's your price, prince"?

| | prince | money | free | xxx | class |
|---|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 | spam |
| 2 | 0 | 0 | 1 | 0 | ham |
| 3 | 0 | 0 | 1 | 0 | ham |
| 4 | 0 | 0 | 1 | 0 | spam |
| 5 | 1 | 0 | 1 | 0 | ham |

$$P(\text{spam} \mid \text{features}) \propto P(\text{spam}) \cdot P\left(\begin{array}{c}\text{yes}\\\text{prince}\end{array} \middle| \text{spam}\right)$$

$$P\left(\begin{array}{c}\text{no}\\\text{money}\end{array} \middle| \text{spam}\right) P\left(\begin{array}{c}\text{no}\\\text{free}\end{array} \middle| \text{spam}\right) P\left(\begin{array}{c}\text{yes}\\\text{xxx}\end{array} \middle| \text{spam}\right)$$

$$P\left(\text{ham} \mid \text{features}\right) \propto \quad \cdot \; - \; - \; - \quad P\left(\begin{array}{c}\text{yes}\\\text{xxx}\end{array} \middle| \text{ham}\right)$$

# Smoothing

▸ **Without** smoothing:

$$P(x^{(i)} = 1 \mid \text{spam}) \approx \frac{\text{\# spam containing word } i}{\text{\# spam containing word } i + \text{\# spam not containing word } i}$$

▸ **With** smoothing:

$$P(x^{(i)} = 1 \mid \text{spam}) \approx \frac{(\text{\# spam containing word } i) + 1}{(\text{\# spam containing word } i) + 1 + (\text{\# spam not containing word } i) + 1}$$

▸ When smoothing, we add 1 to the count of every group whenever we're estimating a conditional probability.

  ▸ **Don't** smooth the estimates of unconditional probabilities (e.g. $P(\text{spam})$).

# Concrete example with smoothing

▶ What happens if we try to classify the email "xxx what's your price money"?

| | prince, money | money, prince | free | x x x | class |
|---|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 | spam |
| 2 | 0 | 0 | 1 | 0 | ham |
| 3 | 0 | 0 | 1 | 0 | ham |
| 4 | 0 | 0 | 1 | 0 | spam |
| 5 | 1 | 0 | 1 | 0 | ham |

$$P(\text{spam} \mid \text{features}) \propto P(\text{spam}) \cdot P\left(\underset{\text{prince}}{\text{yes}} \mid \text{spam}\right)$$

$$P\left(\underset{\text{money}}{\text{no}} \mid \text{spam}\right) P\left(\underset{\text{free}}{\text{no}} \mid \text{spam}\right) P\left(\underset{\text{xxx}}{\text{yes}} \mid \text{spam}\right)$$

$$= \left(\frac{2}{5}\right)\left(\frac{1+1}{1+1+1+1}\right)\left(\frac{1+1}{1+1+1+1}\right)\left(\frac{1+1}{1+1+1+1}\right)\left(\frac{0+1}{0+1+2+1}\right)$$

$$\frac{1}{80} \approx 0.0125$$

$$P(\text{ham} \mid \text{features}) \propto P(\text{ham}) \, P\left(\begin{smallmatrix}\text{yes}\\\text{prince}\end{smallmatrix} \mid \text{ham}\right)$$

$$P\left(\begin{smallmatrix}\text{no}\\\text{money}\end{smallmatrix} \mid \text{ham}\right) \, P\left(\begin{smallmatrix}\text{no}\\\text{free}\end{smallmatrix} \mid \text{ham}\right)$$

**Practical demo** $\quad P\left(\begin{smallmatrix}\text{yes}\\\text{xxx}\end{smallmatrix} \mid \text{ham}\right)$

$$= \left(\frac{3}{5}\right) \left(\frac{1+1}{1+1+2+1}\right) \left(\frac{3+1}{0+1+3+1}\right) \left(\frac{0+1}{0+1+3+1}\right)$$

$$\left(\frac{0+1}{0+1+3+1}\right) = \left(\frac{3}{5}\right) \left(\frac{2}{5}\right) \left(\frac{4}{5}\right) \left(\frac{1}{5}\right) \left(\frac{1}{5}\right)$$

$$\approx 0.0077$$

Follow along with the demo by clicking the **code** link on the course website next to Lecture 24.

**Summary**

## Summary

▶ The Naive Bayes classifier works by estimating the numerator of $P(\text{class}|\text{features})$ for all possible classes.

▶ It uses Bayes' theorem:

$$P(\text{class}|\text{features}) = \frac{P(\text{class}) \cdot P(\text{features}|\text{class})}{P(\text{features})}$$

▶ It also uses a simplifying assumption, that features are conditionally independent given a class:

$P(\text{features}|\text{class}) = P(\text{feature}_1|\text{class}) \cdot P(\text{feature}_2|\text{class}) \cdot \ldots$

▶ The Naive Bayes classifier can be used for text classification, using the bag-of-words model.