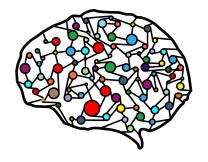
#### **Lecture 24 – More Naive Bayes**



DSC 40A, Fall 2022 @ UC San Diego
Dr. Truong Son Hy, with help from many others

#### **Announcements**

- Look at the readings linked on the course website!
- ► We will have the Thanksgiving break, so there is no class on Friday this week.
- ▶ The final is coming, so there will be a review session.

## **Agenda**

- ► Naive Bayes.
- ► Text classification.
- ► Practical demo.

# **Naive Bayes**

### **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!



/nī'ēv/

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"

unsophisticated innocent artless Similar:

· 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.

ingenuous

inexperienced

## **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 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		
D/I II I 14				

$$P(\text{bad}|\text{male, Marvel}) \propto P(\text{bad}) \cdot P(\text{male, Marvel}|\text{bad})$$
  
 $P(\text{male, Marvel}|\text{bad}) = P(\text{male}|\text{bad}) \cdot P(\text{Marvel}|\text{bad})$ 

 $P(\text{bad}) = \frac{5}{10}$ 

$$P(\text{male}|\text{bad}) = \frac{3}{5}$$
$$P(\text{Marvel}|\text{bad}) = \frac{2}{5}$$

$$P(\text{bad}|\text{male, Marvel}) \propto \frac{5 \cdot 3 \cdot 2}{10 \cdot 5 \cdot 5} = \frac{3}{25}$$

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

 $P(\text{good}|\text{male, Marvel}) \propto P(\text{good}) \cdot P(\text{male, Marvel}|\text{good})$ 

 $P(\text{male, Marvel}|\text{good}) = P(\text{male}|\text{good}) \cdot P(\text{Marvel}|\text{good})$ 

$$P(good) = \frac{4}{10}$$

$$P(\text{male}|\text{good}) = \frac{2}{4}$$

$$P(Marvel|good) = \frac{3}{4}$$

$$P(\text{good}|\text{male, Marvel}) \propto \frac{4 \cdot 2 \cdot 3}{10 \cdot 4 \cdot 4} = \frac{3}{20}$$

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

$$P(\text{neutral}|\text{male, Marvel}) \propto P(\text{neutral}) \cdot P(\text{male, Marvel}|\text{neutral})$$

$$P(\text{male}, \text{Marvel}|\text{neutral}) = P(\text{male}|\text{neutral}) \cdot P(\text{Marvel}|\text{neutral})$$

$$P(\text{neutral}) = \frac{1}{10}$$

$$P(\text{male}|\text{neutral}) = \frac{1}{1} = 1$$

$$P(Marvel|neutral) = \frac{1}{1} = 1$$

$$P(\text{neutral}|\text{male, Marvel}) \propto \frac{1}{10}$$

### **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 **male** Marvel character. Using Naive Bayes, would we predict that my favorite character is bad, good, or neutral? Good!

## **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. What is the probability that this character is 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
		•

 $P(\text{neutral}|\text{female, Marvel}) \propto P(\text{neutral}) \cdot P(\text{female, Marvel}|\text{neutral})$ 

 $P(\text{female}, \text{Marvel}|\text{neutral}) = P(\text{female}|\text{neutral}) \cdot P(\text{Marvel}|\text{neutral})$ 

$$P(\text{neutral}) = \frac{1}{10}$$

$$P(\text{female}|\text{neutral}) = \frac{0}{1} = 0$$

$$P(Marvel|neutral) = \frac{1}{1} = 1$$

$$P(\text{neutral}|\text{female, Marvel}) \propto 0$$

#### 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(neutral) · P(female, Marvel|neutral) = P(neutral) · P(female|neutral) · P(Marvel|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**

Without smoothing:

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

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

► With smoothing:

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

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

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

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

 $P(\text{neutral}|\text{female, Marvel}) \propto P(\text{neutral}) \cdot P(\text{female, Marvel}|\text{neutral})$ 

 $P(female, Marvel|neutral) = P(female|neutral) \cdot P(Marvel|neutral)$ 

$$P(\text{neutral}) = \frac{1}{10}$$

$$P(\text{female}|\text{neutral}) = \frac{1}{2}$$

$$P(Marvel|neutral) = \frac{2}{3}$$

$$P(\text{neutral}|\text{female, Marvel}) \propto \frac{1}{30}$$

#### **Summary: Naive Bayes classifier**

- In classification, our goal is to predict a discrete category, called a class, given some features.
- We want to predict a class, given certain 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!

#### **Summary: Naive Bayes classifier**

- ► The Naive Bayes classifier works by estimating the numerator of *P*(class|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:

$$P(\text{feature}_1|\text{class}) \cdot P(\text{feature}_2|\text{class}) \cdot \dots$$

### **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).

#### Shutterfly

11/3/21

Thank us later—snag an EXTRA 20% OFF your holiday card an... Plus, claim your 4 freebies (today only)! > | View web version floorer cards and gifts now to avoid delays UP TO 50% OFF...

#### Alumni Alliances

11/2/21

Univ. of Cal. Berkeley Alumni Club Invites Suraj from Halicioğl...
Have you claimed your members-only access? Hi Suraj, You're
Invited to Join Alumni Alliances. an invitation-only alumni club....

#### IRS.gov

11/1/21

Re: You are Eligible For a Tax Return on Nov 1, 06:01:52 pm Third Round of Economic Impact Payments Status Available.

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**  $\vec{x}$ :

$$\vec{X} = \begin{bmatrix} X^{(1)} \\ X^{(2)} \\ \dots \\ X^{(d)} \end{bmatrix}$$

#### 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.

- 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."

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

- To classify an email, we'll use Bayes' theorem to calculate the probability of it belonging to each class:
  - P(spam | features).
  - P(ham | features).
- We'll predict the class with a larger probability.

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

- Note that the formulas for P(spam | features) and P(ham | features) have the same denominator, P(features).
- Thus, we can find the larger probability just by comparing numerators:
  - $\triangleright$   $P(\text{spam}) \cdot P(\text{features} \mid \text{spam}).$
  - $\triangleright$   $P(\text{ham}) \cdot P(\text{features} \mid \text{ham}).$

#### **Discussion Question**

We need to determine four quantities:

- 1. P(features | spam).
- 2. P(features | ham).
- 3. *P*(spam).
- 4. *P*(ham).

Which of these probabilities should add to 1?

- A) 1, 2
- B) 3, 4
- C) Both A and B
- D) Neither A nor B

#### **Discussion Question**

We need to determine four quantities:

- 1. P(features | spam).
- 2. P(features | ham).
- 3. *P*(spam).
- 4. P(ham).

Which of these probabilities should add to 1?

- A) 1, 2
- B) 3, 4
- C) Both A and B
- D) Neither A nor B

**Answer:** B) P(spam) + P(ham) = 1.

### Estimating probabilities with training data

► To estimate *P*(spam), we compute

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

► To estimate P(ham), we compute

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

▶ What about P(features | spam) and P(features | ham)?

### **Assumption of conditional independence**

▶ Note that *P*(features | spam) looks like

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

- 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(features | spam) as

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

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."

	prince	money	free	XXX	Label
Sentence 1	1	1	0	0	spam
Sentence 2	0	0	1	0	ham
Sentence 3	0	0	1	0	ham
Sentence 4	0	0	1	0	spam
Sentence 5	1	0	1	0	ham

	prince	money	free	XXX	Label
Sentence 1	1	1	0	0	spam
Sentence 2	0	0	1	0	ham
Sentence 3	0	0	1	0	ham
Sentence 4	0	0	1	0	spam
Sentence 5	1	0	1	0	ham

$$x^{(1)}$$
 = prince,  $x^{(2)}$  = money,  $x^{(3)}$  = free,  $x^{(4)}$  = xxx

#### **Prior:**

$$P(\text{spam}) = \frac{2}{5}$$
$$P(\text{ham}) = \frac{3}{5}$$

	prince	money	free	XXX	Label
Sentence 1	1	1	0	0	spam
Sentence 2	0	0	1	0	ham
Sentence 3	0	0	1	0	ham
Sentence 4	0	0	1	0	spam
Sentence 5	1	0	1	0	ham

$$x^{(1)}$$
 = prince,  $x^{(2)}$  = money,  $x^{(3)}$  = free,  $x^{(4)}$  = xxx

$$P(x^{(1)} = 0 | \text{spam}) = \frac{1}{2}, \quad P(x^{(1)} = 1 | \text{spam}) = \frac{1}{2},$$

	prince	money	free	XXX	Label
Sentence 1	1	1	0	0	spam
Sentence 2	0	0	1	0	ham
Sentence 3	0	0	1	0	ham
Sentence 4	0	0	1	0	spam
Sentence 5	1	0	1	0	ham

$$x^{(1)}$$
 = prince,  $x^{(2)}$  = money,  $x^{(3)}$  = free,  $x^{(4)}$  = xxx

$$P(x^{(1)} = 0|\text{spam}) = \frac{1}{2}, \quad P(x^{(1)} = 1|\text{spam}) = \frac{1}{2},$$

$$P(x^{(2)} = 0 | \text{spam}) = \frac{1}{2}, \quad P(x^{(2)} = 1 | \text{spam}) = \frac{1}{2},$$

	prince	money	free	XXX	Label
Sentence 1	1	1	0	0	spam
Sentence 2	0	0	1	0	ham
Sentence 3	0	0	1	0	ham
Sentence 4	0	0	1	0	spam
Sentence 5	1	0	1	0	ham

$$x^{(1)}$$
 = prince,  $x^{(2)}$  = money,  $x^{(3)}$  = free,  $x^{(4)}$  = xxx

$$P(x^{(1)} = 0 | \text{spam}) = \frac{1}{2}, \quad P(x^{(1)} = 1 | \text{spam}) = \frac{1}{2},$$

$$P(x^{(2)} = 0 | \text{spam}) = \frac{1}{2}, \quad P(x^{(2)} = 1 | \text{spam}) = \frac{1}{2},$$

$$P(x^{(3)} = 0 | \text{spam}) = \frac{1}{2}, \quad P(x^{(3)} = 1 | \text{spam}) = \frac{1}{2},$$

	prince	money	free	XXX	Label
Sentence 1	1	1	0	0	spam
Sentence 2	0	0	1	0	ham
Sentence 3	0	0	1	0	ham
Sentence 4	0	0	1	0	spam
Sentence 5	1	0	1	0	ham
(1)	(2)	1.	٠	//\	

$$x^{(1)}$$
 = prince,  $x^{(2)}$  = money,  $x^{(3)}$  = free,  $x^{(4)}$  = xxx

ponal probability on spam:  

$$P(x^{(1)} = 0 | \text{spam}) = \frac{1}{2}, \quad P(x^{(1)} = 1 | \text{spam}) = \frac{1}{2},$$

$$P(x^{(2)} = 0 | \text{spam}) = \frac{1}{2}, \quad P(x^{(2)} = 1 | \text{spam}) = \frac{1}{2},$$

$$P(x^{(3)} = 0 | \text{spam}) = \frac{1}{2}, \quad P(x^{(3)} = 1 | \text{spam}) = \frac{1}{2},$$

$$P(x^{(4)} = 0 | \text{spam}) = 1, P(x^{(4)} = 1 | \text{spam}) = 0.$$

	prince	money	free	XXX	Label
Sentence 1	1	1	0	0	spam
Sentence 2	0	0	1	0	ham
Sentence 3	0	0	1	0	ham
Sentence 4	0	0	1	0	spam
Sentence 5	1	0	1	0	ham

$$x^{(1)}$$
 = prince,  $x^{(2)}$  = money,  $x^{(3)}$  = free,  $x^{(4)}$  = xxx

$$P(x^{(1)} = 0 | \text{ham}) = \frac{2}{3}, \quad P(x^{(1)} = 1 | \text{ham}) = \frac{1}{3},$$
 $P(x^{(2)} = 0 | \text{ham}) = 1, \quad P(x^{(2)} = 1 | \text{ham}) = 0,$ 
 $P(x^{(3)} = 0 | \text{ham}) = 0, \quad P(x^{(3)} = 1 | \text{ham}) = 1,$ 
 $P(x^{(4)} = 0 | \text{ham}) = 1, \quad P(x^{(4)} = 1 | \text{ham}) = 0.$ 

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

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

prince	money	free	XXX
1	0	1	0

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

prince	money	free	XXX
1	0	1	0

Probability of **spam**:

= 
$$P(x^{(1)} = 1 | \text{spam}) P(x^{(2)} = 0 | \text{spam}) P(x^{(3)} = 1 | \text{spam}) P(x^{(4)} = 0 | \text{spam})$$
  
=  $\frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot 1 = \frac{1}{8}$ 

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

prince	money	free	XXX
1	0	1	0

Probability of spam:

= 
$$P(x^{(1)} = 1 | \text{spam}) P(x^{(2)} = 0 | \text{spam}) P(x^{(3)} = 1 | \text{spam}) P(x^{(4)} = 0 | \text{spam})$$
  
=  $\frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot 1 = \frac{1}{8}$ 

$$P(\text{spam}|\text{features}) \propto P(\text{features}|\text{spam}) \cdot P(\text{spam}) = \frac{1}{8} \cdot \frac{2}{5} = \frac{1}{20}$$

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

prince	money	free	XXX
1	0	1	0

#### Probability of ham:

= 
$$P(x^{(1)} = 1 | \text{ham}) P(x^{(2)} = 0 | \text{ham}) P(x^{(3)} = 1 | \text{ham}) P(x^{(4)} = 0 | \text{ham})$$
  
=  $\frac{1}{3} \cdot 1 \cdot 1 \cdot 1 = \frac{1}{3}$ 

$$P(\text{ham}|\text{features}) \propto P(\text{features}|\text{ham}) \cdot P(\text{ham}) = \frac{1}{3} \cdot \frac{3}{5} = \frac{1}{5}$$

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

prince	money	free	XXX
1	0	1	0

**Because** 

$$P(\text{ham}|\text{features}) = \frac{1}{5} > P(\text{spam}|\text{features}) = \frac{1}{20},$$

this sentence is classified as ham.

#### Uh oh...

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

#### Uh oh...

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

prince	money	free	XXX
1	0	0	1

There is a keyword "xxx" and the sentence is likely **spam**. But:

$$P(x^{(4)} = 1|\text{spam}) = 0$$

Thus:

$$P(\text{features}|\text{spam}) = 0$$

Then, it will be classified as **ham** with absolute certainty.

# **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*(spam)).

	prince	money	free	XXX	Label
Sentence 1	1	1	0	0	spam
Sentence 2	0	0	1	0	ham
Sentence 3	0	0	1	0	ham
Sentence 4	0	0	1	0	spam
Sentence 5	1	0	1	0	ham

$$x^{(1)}$$
 = prince,  $x^{(2)}$  = money,  $x^{(3)}$  = free,  $x^{(4)}$  = xxx

#### **Conditional probability on spam:**

$$P(x^{(1)} = 0 | \text{spam}) = \frac{1}{2}, \quad P(x^{(1)} = 1 | \text{spam}) = \frac{1}{2},$$
 $P(x^{(2)} = 0 | \text{spam}) = \frac{1}{2}, \quad P(x^{(2)} = 1 | \text{spam}) = \frac{1}{2},$ 
 $P(x^{(3)} = 0 | \text{spam}) = \frac{1}{2}, \quad P(x^{(3)} = 1 | \text{spam}) = \frac{1}{2},$ 
 $P(x^{(4)} = 0 | \text{spam}) = \frac{2}{3}, \quad P(x^{(4)} = 1 | \text{spam}) = \frac{1}{3}.$ 

	prince	money	free	XXX	Label
Sentence 1	1	1	0	0	spam
Sentence 2	0	0	1	0	ham
Sentence 3	0	0	1	0	ham
Sentence 4	0	0	1	0	spam
Sentence 5	1	0	1	0	ham

$$x^{(1)}$$
 = prince,  $x^{(2)}$  = money,  $x^{(3)}$  = free,  $x^{(4)}$  = xxx

#### **Conditional probability on ham:**

$$P(x^{(1)} = 0 | \text{ham}) = \frac{3}{5}, \quad P(x^{(1)} = 1 | \text{ham}) = \frac{2}{5},$$
 $P(x^{(2)} = 0 | \text{ham}) = \frac{2}{3}, \quad P(x^{(2)} = 1 | \text{ham}) = \frac{1}{3},$ 
 $P(x^{(3)} = 0 | \text{ham}) = \frac{1}{3}, \quad P(x^{(3)} = 1 | \text{ham}) = \frac{2}{3},$ 
 $P(x^{(4)} = 0 | \text{ham}) = \frac{2}{3}, \quad P(x^{(4)} = 1 | \text{ham}) = \frac{1}{3}.$ 

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

prince	money	free	XXX
1	0	0	1

Probability of **spam**:

= 
$$P(x^{(1)} = 1 | \text{spam}) P(x^{(2)} = 0 | \text{spam}) P(x^{(3)} = 0 | \text{spam}) P(x^{(4)} = 1 | \text{spam})$$
  
=  $\frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{3} = \frac{1}{24}$ 

$$P(\text{spam}|\text{features}) \propto P(\text{features}|\text{spam}) \cdot P(\text{spam}) = \frac{1}{24} \cdot \frac{2}{5} = \frac{1}{60} \approx 0.0166$$

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

prince	money	free	XXX
1	0	0	1

#### Probability of ham:

= 
$$P(x^{(1)} = 1 | \text{ham}) P(x^{(2)} = 0 | \text{ham}) P(x^{(3)} = 0 | \text{ham}) P(x^{(4)} = 1 | \text{ham})$$
  
=  $\frac{2}{5} \cdot \frac{2}{3} \cdot \frac{1}{3} \cdot \frac{1}{3} = \frac{4}{135}$ 

$$P(\text{ham}|\text{features}) \propto P(\text{features}|\text{ham}) \cdot P(\text{ham}) = \frac{4}{135} \cdot \frac{3}{5} \approx 0.0177$$

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

We have:

 $P(\text{spam}|\text{features}) \approx 0.0166$ 

 $P(\text{ham}|\text{features}) \approx 0.0177$ 

Probability of spam: 48.3% Probability of ham: 51.7%

This is a confusing case for Naive Bayes classifier. We need more data!

# **Practical demo**

## More realistic example

**My source code in Java** (it is easier to do in Python):

https://github.com/HyTruongSon/Spambase-filtering

#### Data:

https://archive.ics.uci.edu/ml/datasets/Spambase

**Classifiers:** Linear/RBF Support Vector Machine, Logistic Regression and Multilayer Perceptron.