

# Machine Learning

## Lecture 4. Supervised learning Naïve Bayes

**Alireza Rezvanian**

Fall 2022

Amirkabir University of Technology (Tehran Polytechnic)

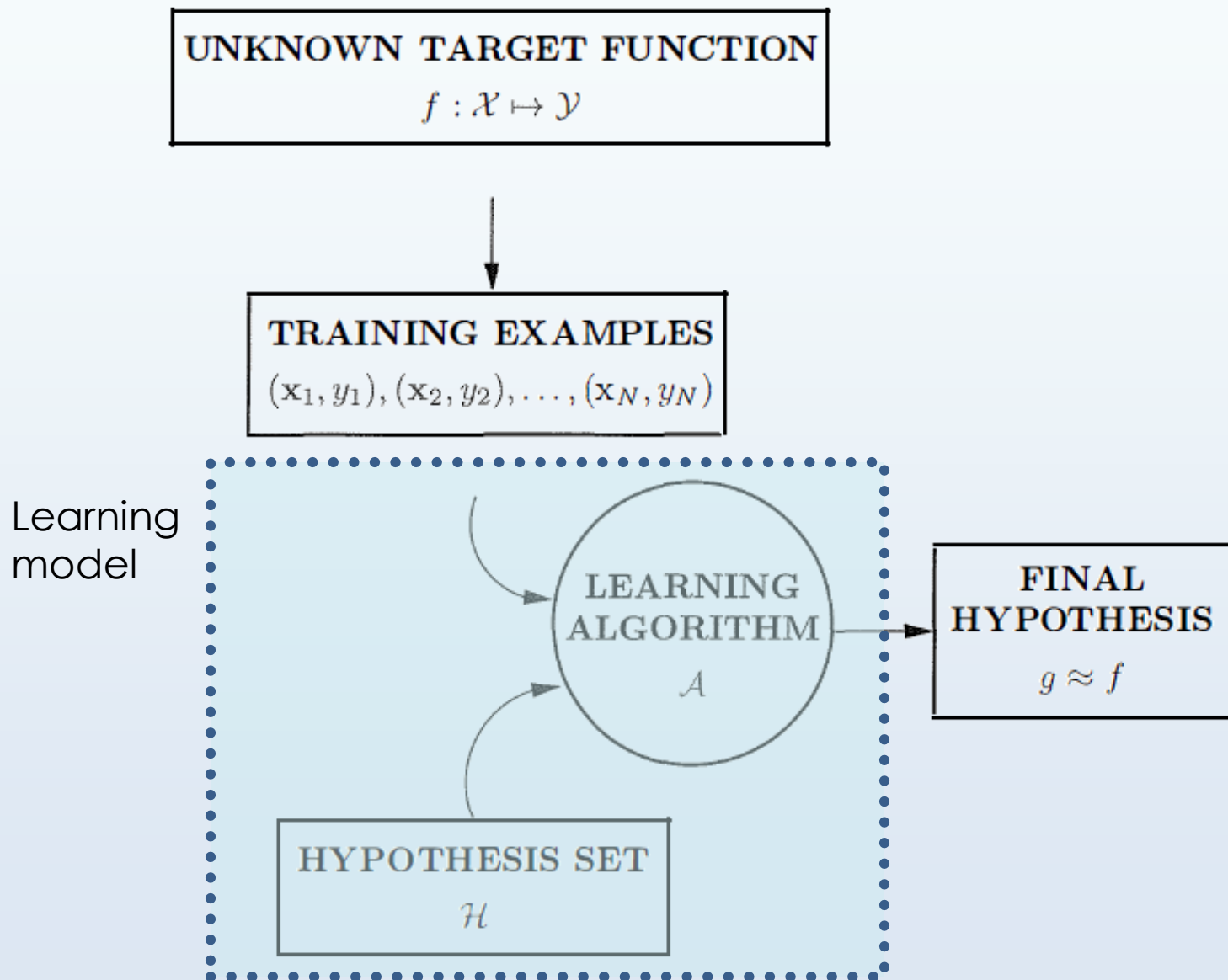
Last update: Oct. 22, 2021



# Outline

- ➡ Probability and Inference
- ➡ Bayes' Rule
- ➡ MAP & MLE
- ➡ Naive Bayes Classifier
- ➡ m-estimate
- ➡ Text classification

# Supervised Learning



# Statistical Estimation View

## ► Probabilities to rescue:

- $x$  and  $y$  are *random variables*
- $D = (x_1, y_1), (x_2, y_2), \dots, (x_N, y_N) \sim P(X, Y)$

## ► IID: **I**ndependent **I**dentically **D**istributed

- Both training & testing data sampled IID from  $P(X, Y)$
- Learn on training set
- Have some hope of *generalizing* to test set

# Interpreting Probabilities

- ➡ What does  $P(A)$  mean?
- ➡ Frequentist View
  - $\lim_{N \rightarrow \infty} \#(A \text{ is true})/N$
  - limiting frequency of a repeating non-deterministic event
- ➡ Bayesian View
  - $P(A)$  is your “belief” about  $A$
- ➡ Market Design View
  - $P(A)$  tells you how much you would bet

# Concepts

## ► Likelihood: $P(D|h)$

- How much does a certain hypothesis explain the data?

## ► Prior: $P(h)$

- What do you believe before seeing any data?

## ► Posterior: $P(h|D)$

- What do we believe after seeing the data?

# Classification

► **Learn:**  $h: \mathbf{X} \mapsto Y$

- $\mathbf{X}$  – features
- $Y$  – target classes

► Suppose you know  $P(Y | \mathbf{X})$  exactly, how should you classify?

- Bayes classifier:

► **Why?**

# Probability and Inference

➤ Result of tossing a coin is  $\in \{\text{Heads}, \text{Tails}\}$

➤ Random var  $X \in \{1, 0\}$

$$\text{Bernoulli: } P\{X=1\} = p_o^X (1 - p_o)^{(1-X)}$$

➤ Sample:  $\mathbf{X} = \{x^t\}_{t=1}^N$

$$\text{Estimation: } p_o = \# \{\text{Heads}\} / \# \{\text{Tosses}\} = \sum_t x^t / N$$

➤ Prediction of next toss:

Heads if  $p_o > 1/2$ , Tails otherwise



# Classification

- **Credit scoring:** Inputs are income and savings.  
Output is low-risk vs. high-risk
- **Input:**  $\mathbf{x} = [x_1, x_2]^T$ , Output:  $C \in \{0, 1\}$
- **Prediction:**

$$\text{choose } \begin{cases} C = 1 \text{ if } P(C = 1 | x_1, x_2) > 0.5 \\ C = 0 \text{ otherwise} \end{cases}$$

or

$$\text{choose } \begin{cases} C = 1 \text{ if } P(C = 1 | x_1, x_2) > P(C = 0 | x_1, x_2) \\ C = 0 \text{ otherwise} \end{cases}$$

# Bayes Theorem

The diagram shows the Bayes Theorem equation with labels and arrows indicating the components:

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

Labels and arrows:

- posterior* points to  $P(h|D)$
- likelihood* points to  $P(D|h)$
- prior* points to  $P(h)$
- evidence* points to  $P(D)$

**P(h)**: Prior probability of hypothesis h (**Prior**)

**P(D)**: Prior probability of training data D (**Evidence**)

**P(h | D)**: Probability of h given D (**Posterior**)

**P(D | h)**: Probability of D given h (**Likelihood**)

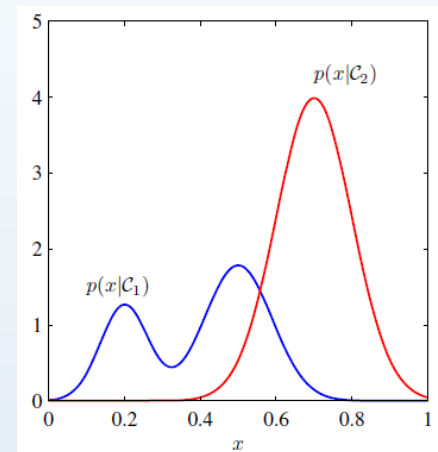
# Generative vs. Discriminative

- Using Bayes rule, optimal classifier

$$h^*(x) = \underset{c}{\operatorname{argmax}} \{ \log p(x|y=c) + \log p(y=c) \}$$

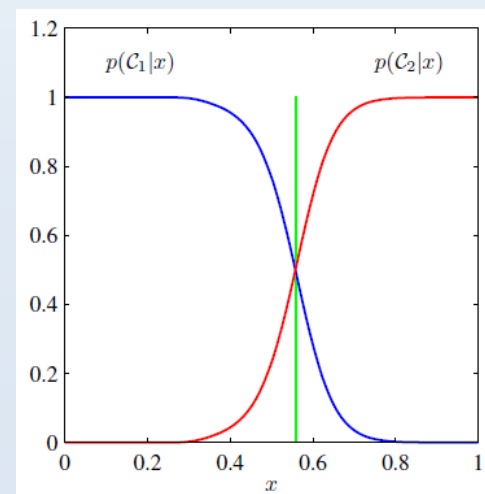
- Generative Approach

- Estimate  $\mathbf{p(x | y)}$  and  $\mathbf{p(y)}$
- Use Bayes Rule to predict  $\mathbf{y}$



- Discriminative Approach

- Estimate  $\mathbf{p(y | x)}$  directly OR
- Learn “discriminant” function  $\mathbf{h(x)}$



# Choosing hypothesis

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

- Generally we want the most probable hypothesis given the training data

- Maximum A Posteriori (MAP) hypothesis  $\mathbf{h}_{MAP}$

$$\begin{aligned} h_{MAP} &\equiv \arg \max_{h \in H} P(h | D) \\ &= \arg \max_{h \in H} \frac{P(D | h)P(h)}{\cancel{P(D)}} \quad \leftarrow \text{Evidence is a Constant independent of } \mathbf{h} \\ &= \arg \max_{h \in H} P(D | h)P(h) \end{aligned}$$

- If assume  $P(h_i) = P(h_j)$  then can further simplify, and choose

- Maximum likelihood (ML) hypothesis

$$\begin{aligned} h_{ML} &\equiv \arg \max_{h \in H} P(h | D) \\ &= \arg \max_{h \in H} \frac{P(D | h)\cancel{P(h)}}{\cancel{P(D)}} \quad \leftarrow \begin{array}{l} \text{No prior hypothesis} \\ \text{OR} \\ \text{equally probable} \\ \text{a priori } (P(h_i) = P(h_j)) \end{array} \\ &= \arg \max_{h \in H} P(D | h) \end{aligned}$$

## Example: a medical diagnosis problem

- ➡ There are two alternative hypotheses:
  - $h(1)$  that the patient has a particular form of cancer. (+)
  - $h(2)$  that the patient does no (-)
- ➡ Prior knowledge indicates that over the entire population only **0.008** have this disease.
  - $P(\text{cancer}) = 0.008$        $P(\sim\text{cancer}) = 0.992$
- ➡ The test returns a correct **positive result** in only **98%** of the actual patient
  - $P(+ \mid \text{cancer}) = 0.98$        $P(- \mid \text{cancer}) = 0.02$
- ➡ The test returns a correct **negative result** in only **97%** of the not patient
  - $P(- \mid \sim\text{cancer}) = 0.97$        $P(+ \mid \sim\text{cancer}) = 0.03$

# Medical diagnosis problem

- ➡ Suppose we now observe a new patient for whom the lab test returns a **positive result**. Should we diagnose the patient as having cancer or not?  $P(\text{cancer} | +) = ?$

- The maximum a posteriori (MAP) hypothesis

$$P(+ | \text{cancer}) \times P(\text{cancer}) = (0.98) \times (0.008) = 0.0078$$

$$P(+ | \sim\text{cancer}) \times P(\sim\text{cancer}) = (0.03) \times (0.992) = 0.0298$$

- ➡ Thus,  $h_{MAP} = \sim\text{cancer}$

Note: The exact posterior probabilities by normalizing

$$P(\text{cancer} | +) = \frac{0.0078}{0.0078 + 0.0298} = 0.21$$

# Brute-force MAP Learning

1) For each hypothesis  $h$  in  $H$ , calculate the posterior probability

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

2) Output the hypothesis  $\mathbf{h}_{MAP}$  with the highest posterior probability

$$h_{MAP} = \operatorname{argmax}_{h \in H} P(h|D)$$

# Naive Bayes Classifier

- Target function:  $f: X \rightarrow V$
- tuple of attribute values  $\langle \mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n \rangle$
- **Goal:** predict the target value, or classification, for this new instance.
- The Bayesian approach is to assign the most probable target value,  $V_{MAP}$ , given the attribute values  $\langle \mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n \rangle$  that describe the instance

$$V_{MAP} = \operatorname{argmax}_{v_j \in V} P(v_j | a_1, a_2, \dots, a_n)$$

$$V_{MAP} = \operatorname{argmax}_{v_j \in V} \frac{P(a_1, a_2, \dots, a_n | v_j) P(v_j)}{P(a_1, a_2, \dots, a_n)}$$

$$= \operatorname{argmax}_{v_j \in V} P(a_1, a_2, \dots, a_n | v_j) P(v_j)$$



# Naive Bayes Classifier

- $V_{MAP} = \operatorname{argmax}_{v_j \in V} P(a_1, a_2, \dots, a_n | v_j) P(v_j)$ 
  - $P(\mathbf{v}_j)$  computes by counting the frequency with which each target value  $\mathbf{v}_j$  occurs in the training data
  - Estimating  $P(a_1, a_2, \dots, a_n | v_j)$  is not feasible
- The naive Bayes classifier is based on the simplifying **assumption** that the attribute values are conditionally independent given the target value

$$P(a_1, a_2, \dots, a_n | v_j) = P(a_1 | v_j) \times P(a_2 | v_j) \dots P(a_n | v_j) = \prod_i^n P(a_i | v_j)$$

$$V_{NB} = \operatorname{argmax}_{v_j \in V} P(v_j) \prod_i^n P(a_i | v_j)$$

- If assumption holds, NB is optimal classifier!

# Naive Bayes Algorithm

Naïve\_Bayes\_Learn(*example*)

For each target value  $\mathbf{v}_j$

$$\hat{P}(v_j) \leftarrow \text{estimate } P(v_j)$$

For each attribute value  $\mathbf{a}_i$  of each attribute  $\mathbf{a}$

$$\hat{P}(a_i|v_j) \leftarrow \text{estimate } P(a_i|v_j)$$

Classify\_New\_Instance( $\mathbf{x}$ )

$$V_{NB} = \underset{v_j \in V}{\operatorname{argmax}} (\hat{P}(v_j) \prod_i^n \hat{P}(a_i|v_j))$$

# Example: car theft problem

Example No.	Color	Type	Origin	Stolen?
1	Red	Sports	Domestic	Yes
2	Red	Sports	Domestic	No
3	Red	Sports	Domestic	Yes
4	Yellow	Sports	Domestic	No
5	Yellow	Sports	Imported	Yes
6	Yellow	SUV	Imported	No
7	Yellow	SUV	Imported	Yes
8	Yellow	SUV	Domestic	No
9	Red	SUV	Imported	No
10	Red	Sports	Imported	Yes

classify a **Red Domestic SUV** is getting stolen or not?

# Example: car theft problem

- $P(\text{Stolen}=\text{Yes}) = 5/10$
- $P(\text{Stolen}=\text{No}) = 5/10$
- $P(\text{Color}=\text{Red} \mid \text{Stolen}=\text{Yes}) = 3/5$
- $P(\text{Color}=\text{Red} \mid \text{Stolen}=\text{No}) = 2/5$
- $P(\text{Type}=\text{SUV} \mid \text{Stolen}=\text{Yes}) = 1/5$
- $P(\text{Type}=\text{SUV} \mid \text{Stolen}=\text{No}) = 3/5$
- $P(\text{Origin}=\text{Domestic} \mid \text{Stolen}=\text{Yes}) = 2/5$
- $P(\text{Origin}=\text{Domestic} \mid \text{Stolen}=\text{No}) = 3/5$

Prior

Likelihood

- $P(\text{Stolen}=\text{Yes} \mid \text{Color}=\text{Red}, \text{Type}=\text{SUV}, \text{Origin}=\text{Domestic}) =$   
 $P(\text{Color}=\text{Red} \mid \text{Stolen}=\text{Yes}) \times P(\text{Type}=\text{SUV} \mid \text{Stolen}=\text{Yes}) \times P(\text{Origin}=\text{Domestic} \mid \text{Stolen}=\text{Yes}) \times P(\text{Stolen}=\text{Yes}) = ?$   
 $3/5 \times 1/5 \times 2/5 \times 1/2 = 0.024$
- $P(\text{Stolen}=\text{No} \mid \text{Color}=\text{Red}, \text{Type}=\text{SUV}, \text{Origin}=\text{Domestic}) =$   
 $P(\text{Color}=\text{Red} \mid \text{Stolen}=\text{No}) \times P(\text{Type}=\text{SUV} \mid \text{Stolen}=\text{No}) \times P(\text{Origin}=\text{Domestic} \mid \text{Stolen}=\text{No}) \times P(\text{Stolen}=\text{No}) = ?$   
 $2/5 \times 3/5 \times 3/5 \times 1/2 = 0.072$

$$\frac{0.024}{0.024 + 0.072} = 0.25$$

$$\frac{0.072}{0.024 + 0.072} = 0.75$$

**Answer: the car is not stolen!**

# Subtleties of NB classifier 1 – Violating the NB assumption

- ➡ Usually, features are not conditionally independent:

$$P(X_1 \dots X_n | Y) \neq \prod_i P(X_i | Y)$$

- ➡ Probabilities  $P(Y | \mathbf{X})$  often biased towards 0 or 1
- ➡ Nonetheless, NB is a very popular classifier
  - NB often performs well, even when assumption is violated
  - [Domingos & Pazzani '96] discuss some conditions for good performance

## Subtleties of NB classifier 2 – Insufficient training data

- ➡ What if you never see a training instance where  $X_1=a$  when  $Y=c$ ?
  - e.g.,  $Y=\{\text{NonSpamEmail}\}$ ,  $X_1=\{\text{'Nigeria'}\}$
  - $P(X_1=a \mid Y=c) = 0$
- ➡ Thus, no matter what the values  $X_2, \dots, X_d$  take:
  - $P(Y=c \mid X_1=a, X_2, \dots, X_d) = 0$
- ➡ What now???

## Estimating probabilities (Laplace estimation)

- The probabilities estimate by the fraction of times the event is observed to occur over the total number of opportunities.  $\frac{n_c}{n}$
- It may provides poor/incorrect estimates when  $n_c$ , is very small or zero.
- To avoid this difficulty we can adopt a Bayesian approach to estimating the probability, using the **m-estimate** defined as follows.

$$\frac{n_c + mp}{n + m}$$

assume uniform priors; if an attribute has  $k$  possible values

$$p = \frac{1}{k}$$

$m$  is a constant called the **equivalent sample size**

# Text classification

- ➡ Classify e-mails
  - $Y = \{\text{Spam}, \text{NotSpam}\}$
- ➡ Classify news articles
  - $Y = \{\text{what is the topic of the article?}\}$
- ➡ Classify webpages
  - $Y = \{\text{Student, professor, project, ...}\}$
- ➡ What about the features **X**?
  - The text!



# Features $X$ are entire document – $X_i$ for $i^{\text{th}}$ word in article

## Article from rec.sport.hockey

---

Path: cantaloupe.srv.cs.cmu.edu!das-news.harvard.e  
From: xxx@yyy.zzz.edu (John Doe)  
Subject: Re: This year's biggest and worst (opinion)  
Date: 5 Apr 93 09:53:39 GMT

I can only comment on the Kings, but the most obvious candidate for pleasant surprise is Alex Zhitnik. He came highly touted as a defensive defenseman, but he's clearly much more than that. Great skater and hard shot (though wish he were more accurate). In fact, he pretty much allowed the Kings to trade away that huge defensive liability Paul Coffey. Kelly Hrudey is only the biggest disappointment if you thought he was any good to begin with. But, at best, he's only a mediocre goaltender. A better choice would be Tomas Sandstrom, though not through any fault of his own, but because some thugs in Toronto decided

# NB for Text classification

➡  $P(\mathbf{X} | Y)$  is huge!!!

- Article at least 1000 words,  $\mathbf{X}=\{X_1, \dots, X_{1000}\}$
- $X_i$  represents  $i^{\text{th}}$  word in document, i.e., the domain of  $X_i$  is entire vocabulary, e.g., Webster Dictionary (or more), 10,000 words, etc.

➡ NB assumption helps a lot!!!

- $P(X_i=x_i | Y=y)$  is just the probability of observing word  $x_i$  in a document on topic  $y$

$$h_{NB}(\mathbf{x}) = \arg \max_y P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

# Bag of Words model

- ➡ Typical additional assumption:  
**Position in document doesn't matter:**

$$P(X_i=a \mid Y=y) = P(X_k=a \mid Y=y)$$

- “Bag of words” model – order of words on the page ignored
- Sounds really silly, but often works very well!

$$P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

**When the lecture is over, remember to wake up the person sitting next to you in the lecture room.**

# Bag of Words model

► Typical additional assumption:

**Position in document doesn't matter:**

$$P(X_i=x_i | Y=y) = P(X_k=x_i | Y=y)$$

- “Bag of words” model – order of words on the page ignored
- Sounds really silly, but often works very well!

$$P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

in is lecture lecture next over person remember room  
sitting the the the to to up wake when you

# Bag of Words model

the world of  
**TOTAL**



**all about the company**

Our energy exploration, production, and distribution operations span the globe, with activities in more than 100 countries.

At TOTAL, we draw our greatest strength from our fast-growing oil and gas reserves. Our strategic emphasis on natural gas provides a strong position in a rapidly expanding market.

Our expanding refining and marketing operations in Asia and the Mediterranean Rim complement already solid positions in Europe, Africa, and the U.S.

Our growing specialty chemicals sector adds balance and profit to the core energy business.

► All About The Company

- Global Activities
- Corporate Structure
- TOTAL's Story
- Upstream Strategy
- Downstream Strategy
- Chemicals Strategy
- TOTAL Foundation
- Homepage

aardvark 0

about2

all 2

Africa 1

apple0

anxious 0

...

gas 1

...

oil 1

...

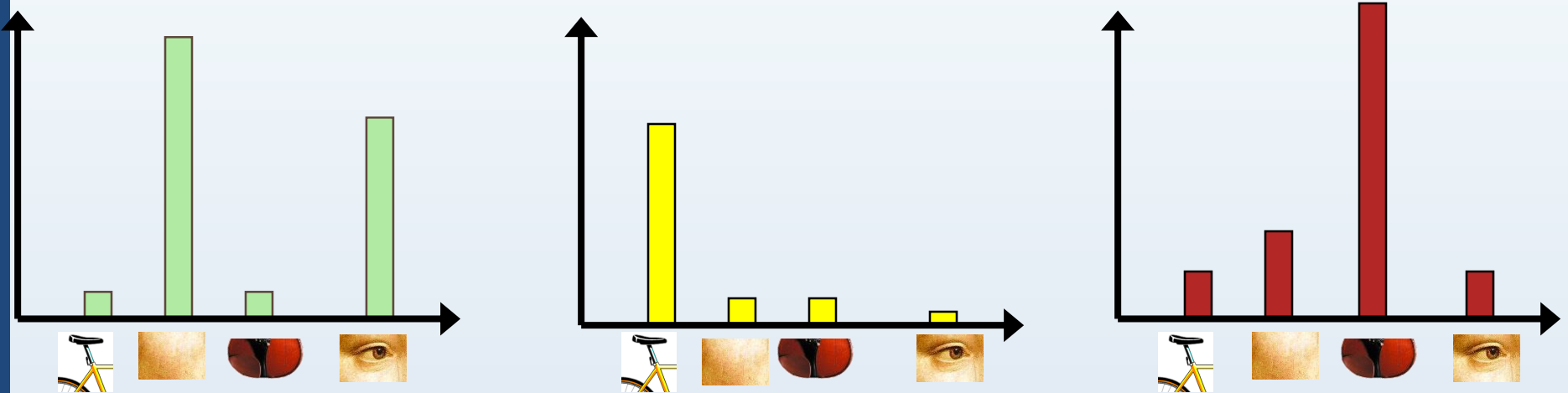
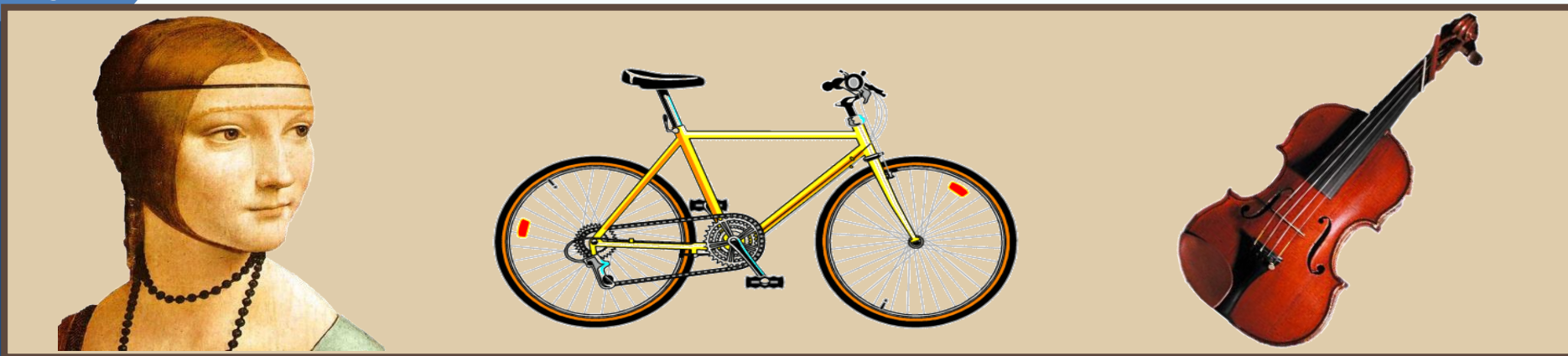
Zaire 0

# Object

# Bag of 'words'



Slide Credit: Fei Fei Li



Slide Credit: Fei Fei Li

# Example

$$P(\text{Class} | \text{Doc}_5) = ?$$

$$P(\text{Class} | \text{Doc}_5) = P(W_{\text{Doc}_5} | \text{Class}) \times P(\text{Class})$$

**Prior:**

$$P(\text{Class} = c) = \frac{3}{4}$$

$$P(\text{Class} = j) = \frac{1}{4}$$

	Doc	Words	Class
Training	1	Chinese Beijing Chinese	c
	2	Chinese Chinese Shanghai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Chinese Tokyo Japan	?

**Likelihood:**  $P(w_j | v_i) = \frac{\text{count}(w_j, v_i) + 1}{\text{count}(v_i) + m}$

$$P(\text{Chinese} | c) = (5 + 1) / (8 + 6) = 3/7$$

$$P(\text{Tokyo} | c) = (0 + 1) / (8 + 6) = 1/14$$

$$P(\text{Japan} | c) = (0 + 1) / (8 + 6) = 1/14$$

$$P(\text{Chinese} | j) = (1 + 1) / (3 + 6) = 2/9$$

$$P(\text{Tokyo} | j) = (1 + 1) / (3 + 6) = 2/9$$

$$P(\text{Japan} | j) = (1 + 1) / (3 + 6) = 2/9$$

$$\begin{aligned} P(c | d_5) &= P(\text{Chinese} | c) \times P(\text{Chinese} | c) \times P(\text{Chinese} | c) \times P(\text{Tokyo} | c) \times P(\text{Japan} | c) \times P(c) \\ &= (3/7) \times (3/7) \times (3/7) \times (1/14) \times (1/14) \times (3/4) \approx 0.0003 \end{aligned}$$

$$\begin{aligned} P(j | d_5) &= P(\text{Chinese} | j) \times P(\text{Chinese} | j) \times P(\text{Chinese} | j) \times P(\text{Tokyo} | j) \times P(\text{Japan} | j) \times P(j) \\ &= (2/9) \times (2/9) \times (2/9) \times (2/9) \times (2/9) \times (1/4) \approx 0.0001 \end{aligned}$$



# Technical Detail: Underflow

- Multiplying lots of probabilities, which are between 0 and 1 by definition, can result in floating-point underflow.
- Since  $\log(xy) = \log(x) + \log(y)$ , it is better to perform all computations by summing logs of probabilities rather than multiplying probabilities.
- Class with highest final un-normalized log probability
- score is still the most probable

$$C_{NB} = \underset{v_j \in V}{\operatorname{argmax}} (\log P(v_j) + \sum_{i=1}^n \log P(a_i | v_j))$$

# What you need to know about NB

- Optimal decision using Bayes Classifier
  - the attribute values are conditionally independent given the target value
- Along with decision tree, neural networks, nearest neighbor, one of the most practical learning methods.
- When to use
  - Moderate or large training set available
  - Attributes that describe instances are conditionally independent given classification
- Very fast
  - Learns with one pass over the data
  - Testing linear in the number of attributes and of documents
  - Low storage requirements
- Successful applications
  - Diagnosis
  - Text classification (Bag of words model)
- Gaussian NB
  - Features are still conditionally independent
  - Each feature has a Gaussian distribution given class

# Reading

- T. Mitchel, **Machine learning**, 1998. (ch. 6)
- C. M. Bishop, **Pattern recognition and machine learning**, Springer, 2006. (ch. 1)
- E. Alpaydin, **Introduction to Machine Learning**, 3<sup>rd</sup> ed., The MIT Press, 2014. (ch. 3)