

# Document Classification

- **What is the role of document classification in information retrieval?**

# Classification Methods (3):

---

## Supervised learning

- Given:
  - A document  $d$
  - A fixed set of classes:  
 $C = \{c_1, c_2, \dots, c_J\}$
  - A training set  $D$  of documents each with a label in  $C$
- Determine:
  - A learning method or algorithm which will enable us to learn a classifier  $\gamma$
  - For a test document  $d$ , we assign it the class  
 $\gamma(d) \in C$

# Document classification

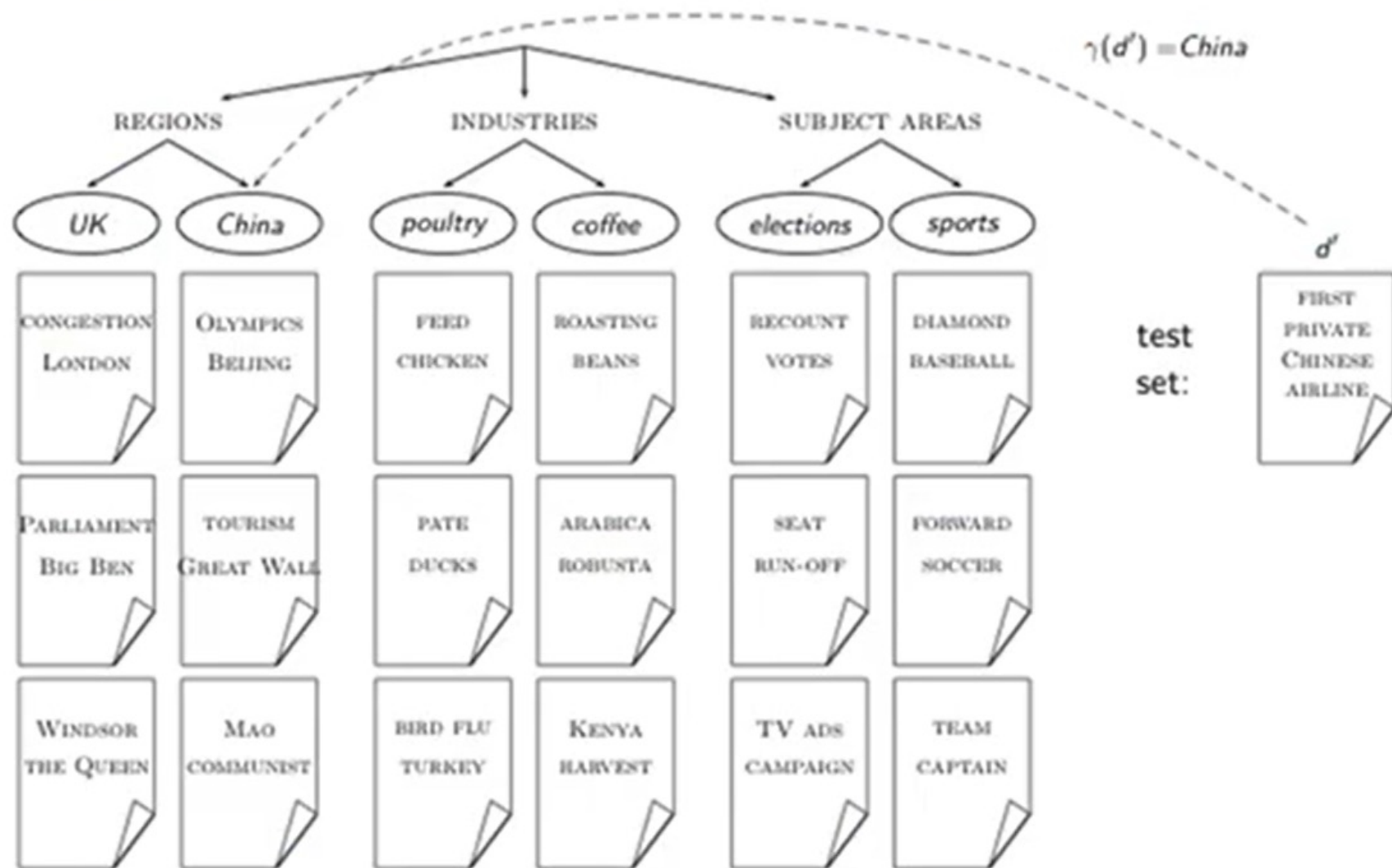
achieve

lead

classes:

training set:

test set:



# The Naive Bayes classifier

- The Naive Bayes classifier is a probabilistic classifier.
- We compute the probability of a document  $d$  being in a class  $c$  as follows:

$$P(c|d) \propto P(c) \prod_{1 \leq k \leq n_d} P(t_k|c)$$

- $n_d$  is the length of the document. (number of tokens)
- $P(t_k|c)$  is the conditional probability of term  $t_k$  occurring in document of class  $c$
- $P(t_k | c)$  as a measure of **how much evidence**  $t_k$  contributes that  $c$  is the correct class.
- $P(c)$  is the prior probability of  $c$ .
- If a document's terms do not provide clear evidence for one class vs. another, we choose the  $c$  with highest  $P(c)$ .

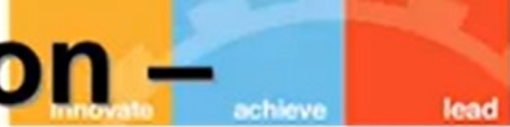
# Maximum a posteriori class

- Our goal in Naive Bayes classification is to find the “best” class.
- The best class is the most likely or maximum a posteriori (MAP) class  $c_{\text{map}}$ :

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} \hat{P}(c|d) = \arg \max_{c \in \mathbb{C}} \hat{P}(c) \prod_{1 \leq k \leq n_d} \hat{P}(t_k|c)$$

# Parameter estimation –

## Maximum likelihood



- Estimate parameters  $\hat{P}(c)$  and  $\hat{P}(t_k|c)$  from train data: How?
- Prior:

$$\hat{P}(c) = \frac{N_c}{N}$$

- $N_c$ : number of docs in class  $c$ ;  $N$ : total number of docs
- Conditional probabilities: 
$$\hat{P}(t|c) = \frac{T_{ct}}{\sum_{t' \in V} T_{ct'}}$$
- $T_{ct}$  is the number of tokens of  $t$  in training documents from class  $c$  (includes multiple occurrences)
- We've made a **Naive Bayes positional independence assumption** here:

$$\hat{P}(t_{k_1}|c) = \hat{P}(t_{k_2}|c)$$



# Second independence assumption



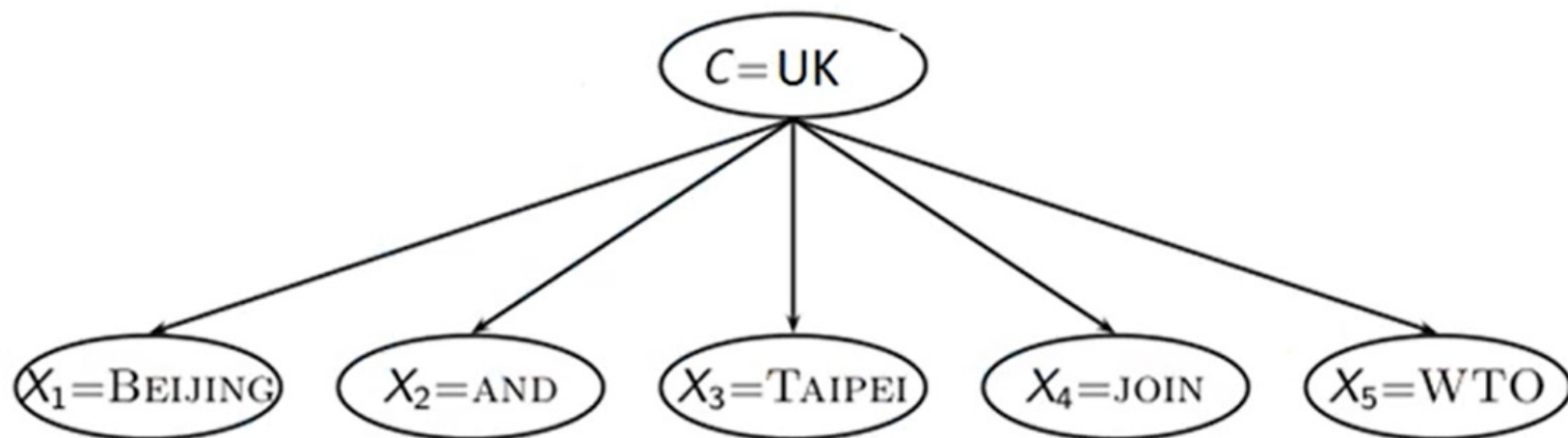
---

$$\hat{P}(t_{k_1}|c) = \hat{P}(t_{k_2}|c)$$

- For example, for a document in the class *UK*, the probability of generating QUEEN in the first position of the document is the same as generating it in the last position.
- The two independence assumptions amount to the **bag of words** model.



# The problem with maximum likelihood estimates: Zeros



- **If WTO never occurs in class China in the train set:**

$$\hat{P}(\text{WTO}|\text{China}) = \frac{T_{\text{China},\text{WTO}}}{\sum_{t' \in V} T_{\text{China},t'}} = \frac{0}{\sum_{t' \in V} T_{\text{China},t'}} = 0$$

# The problem with maximum likelihood estimates: Zeros



- If there were no occurrences of WTO in documents in class China, we'd get a zero estimate:

$$\hat{P}(\text{WTO}|\text{China}) = \frac{T_{\text{China}, \text{WTO}}}{\sum_{t' \in V} T_{\text{China}, t'}} = 0$$

- We will get  $P(\text{China}|d) = 0$  for any document that contains WTO!
- Zero probabilities cannot be conditioned away.

# To avoid zeros: Add-one smoothing

innovate

achieve

lead

- Before:

$$\hat{P}(t|c) = \frac{T_{ct}}{\sum_{t' \in V} T_{ct'}}$$

- Now: Add one to each count to avoid zeros:

$$\hat{P}(t|c) = \frac{T_{ct} + 1}{\sum_{t' \in V} (T_{ct'} + 1)} = \frac{T_{ct} + 1}{(\sum_{t' \in V} T_{ct'}) + B}$$

- B is the number of different words (in this case the size of the vocabulary:)



$$\hat{P}(c) = \frac{N_c}{N}$$

$$\hat{P}(w|c) = \frac{\text{count}(w,c) + 1}{\text{count}(c) + |V|}$$

	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	?

### Priors:

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

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

### Choosing a class:

$$P(c|d5) \propto \frac{3}{4} * \left(\frac{3}{7}\right)^3 * \frac{1}{14} * \frac{1}{14} \approx 0.0003$$

### Conditional Probabilities:

$$P(\text{Chinese}|c) = \frac{(5+1)}{(8+6)} = \frac{6}{14} = \frac{3}{7}$$

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

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

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

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

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

$$P(j|d5) \propto \frac{1}{4} * \left(\frac{2}{9}\right)^3 * \frac{2}{9} * \frac{2}{9} \approx 0.0001$$

# Naive Bayes: Training

TRAINMULTINOMIALNB( $\mathbb{C}, \mathbb{D}$ )

```
1   $V \leftarrow \text{EXTRACTVOCABULARY}(\mathbb{D})$ 
2   $N \leftarrow \text{COUNTDOCS}(\mathbb{D})$ 
3  for each  $c \in \mathbb{C}$ 
4  do  $N_c \leftarrow \text{COUNTDOCSINCLASS}(\mathbb{D}, c)$ 
5       $\text{prior}[c] \leftarrow N_c / N$ 
6       $\text{text}_c \leftarrow \text{CONCATENATETEXTOFALLDOCSINCLASS}(\mathbb{D}, c)$ 
7      for each  $t \in V$ 
8      do  $T_{ct} \leftarrow \text{COUNTTOKENSOFTERM}(\text{text}_c, t)$ 
9      for each  $t \in V$ 
10     do  $\text{condprob}[t][c] \leftarrow \frac{T_{ct}+1}{\sum_{t'} (T_{ct'}+1)}$ 
11 return  $V, \text{prior}, \text{condprob}$ 
```



APPLYMULTINOMIALNB( $\mathbb{C}$ ,  $V$ ,  $prior$ ,  $condprob$ ,  $d$ )

1  $W \leftarrow \text{EXTRACTTOKENSFROMDOC}(V, d)$

2 **for each**  $c \in \mathbb{C}$

3 **do**  $score[c] \leftarrow \log prior[c]$

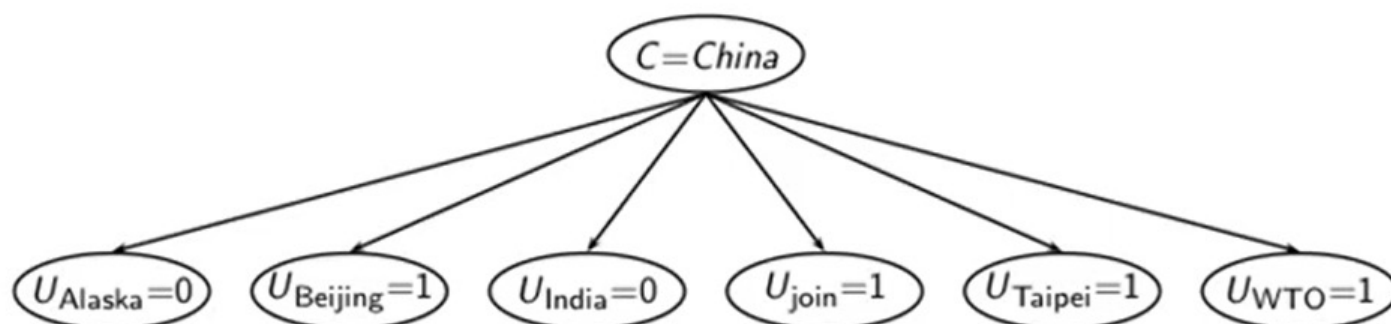
4     **for each**  $t \in W$

5         **do**  $score[c]_+ = \log condprob[t][c]$

6 **return**  $\arg \max_{c \in \mathbb{C}} score[c]$

# A different Naive Bayes model: Bernoulli model

- The Bernoulli model estimates  $P(t|c)$  as the fraction of documents of class  $c$  that contain term  $t$
- The multinomial model estimates  $P(t|c)$  as the fraction of tokens or fraction of positions in documents of class  $c$  that contain term  $t$
- The Bernoulli model considers the binary occurrence information for the terms in the test document ignoring the number of occurrences of the term.
- The probability of non-occurrence of the terms of the vocabulary in the test document is also considered.





# Algorithm

TRAINBERNOULLNB( $\mathbb{C}, \mathbb{ID}$ )

```
1  $V \leftarrow \text{EXTRACTVOCABULARY}(\mathbb{ID})$ 
2  $N \leftarrow \text{COUNTDOCS}(\mathbb{ID})$ 
3 for each  $c \in \mathbb{C}$ 
4 do  $N_c \leftarrow \text{COUNTDOCSINCLASS}(\mathbb{ID}, c)$ 
5    $\text{prior}[c] \leftarrow N_c / N$ 
6   for each  $t \in V$ 
7   do  $N_{ct} \leftarrow \text{COUNTDOCSINCLASSCONTAININGTERM}(\mathbb{ID}, c, t)$ 
8      $\text{condprob}[t][c] \leftarrow (N_{ct} + 1) / (N_c + 2)$ 
9 return  $V, \text{prior}, \text{condprob}$ 
```

APPLYBERNOULLNB( $\mathbb{C}, V, \text{prior}, \text{condprob}, d$ )

```
1  $V_d \leftarrow \text{EXTRACTTERMSFROMDOC}(V, d)$ 
2 for each  $c \in \mathbb{C}$ 
3 do  $\text{score}[c] \leftarrow \log \text{prior}[c]$ 
4   for each  $t \in V$ 
5   do if  $t \in V_d$ 
6     then  $\text{score}[c] += \log \text{condprob}[t][c]$ 
7     else  $\text{score}[c] += \log(1 - \text{condprob}[t][c])$ 
8 return  $\arg \max_{c \in \mathbb{C}} \text{score}[c]$ 
```

# Test doc: “buy cheap dinner”

Doc id	cheap	buy	banking	dinner	the	class
1	0	0	0	0	1	Not spam
2	1	0	1	0	1	spam
3	0	0	0	0	1	Not spam
4	1	0	1	0	1	spam
5	1	1	0	0	1	spam
6	0	0	1	0	1	Not spam
7	0	1	1	0	1	Not spam
8	0	0	0	0	1	Not spam
9	0	0	0	0	1	Not spam
10	1	1	0	1	1	Not spam