



UPPSALA  
UNIVERSITET

# Information Retrieval (5LN712)

## Text classification and naïve Bayes

2020-05-25

Ali Basirat

Department of Linguistics and Philology





# Today

- The text classification problem
- The naïve Bayesian text classification
- Multinomial model
- The Bernoulli model
- A comparison between the two models
- Summary



# The text classification problem

- Queries are generated by probabilistic models from some topic (or document class)
- For example: the terms Taipei and Beijing are generated from a broader document class, China
- In order to classify a document, we should find a class that is more likely to generate the document



# The text classification problem

- For a given document space  $X$  and a set of classes  $C$ , a document classifier maps each document description in  $X$  to a class in  $C$

$$\gamma: X \rightarrow C$$

- The classification function  $\gamma$  is learnt from a training set of labelled documents  $(d, c) \in X \times C$



# The Naïve Bayesian Approach

- The document-class probabilities can be estimated by the Bayesian equation:

$$P(c|d) = \frac{P(c)P(d|c)}{P(d)}$$

- The denominator  $P(d)$  is a constant for a document and can be ignored



# The Naïve Bayesian Approach

- The prior probability  $P(c)$  can be estimate from the data and other information about the classes
- We can make a naïve assumption about the independence occurrence of words in a document
- This assumption is not correct but helps us to have an estimation of the document probabilities



# The Naïve Bayesian Approach

- If we consider a multinomial probability distribution over positions

$$P(d|c) \approx \prod_{t \in d} P(t|c)$$

- If we consider a Bernoulli distribution over the occurrence of words

$$P(d|c) \approx \prod_{e \in V} P(e|c)$$



# Multinomial model

- The probability of a document  $d$  being in class  $c$

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

- The best class for a document in NB is the most likely one (maximum a posteriori):

$$c_{map} = \arg \max_{c \in C} \hat{P}(c|d)$$





# Multinomial model

- If we expand  $\hat{P}(c|d)$ , we will have

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

- We can also take a logarithm from the right side

$$c_{map} = \arg \max_{c \in C} [\log \hat{P}(c) + \sum_{1 \leq k \leq n_d} \log \hat{P}(t_k | c)]$$



# Multinomial model

- We may have zero probabilities for some terms
- It can be resolved by an add-one or Laplace smoothing

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

- $|c|$  is the total number of terms in class  $c$
- $|V|$  is the vocabulary size of the collection



# Multinomial model

	Doc ID	Words in documents	class
Training set	1	Chinese Beijing Chinese	1
	2	Chinese Chinese Shanghai	1
	3	Chinese Macao	1
	4	Tokyo Japan Chinese	0
Test set	5	Chinese Chinese Chinese Tokyo Japan	?



# Multinomial model

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

```
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    $\text{text}_c \leftarrow \text{CONCATENATETEXTOFALLDOCSINCLASS}(\mathbb{ID}, c)$ 
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( $\mathbb{C}, V, \text{prior}, \text{condprob}, d$ )

```
1  $W \leftarrow \text{EXTRACTTOKENSFROMDOC}(V, d)$ 
2 for each  $c \in \mathbb{C}$ 
3 do  $\text{score}[c] \leftarrow \log \text{prior}[c]$ 
4   for each  $t \in W$ 
5   do  $\text{score}[c] += \log \text{condprob}[t][c]$ 
6 return  $\arg \max_{c \in \mathbb{C}} \text{score}[c]$ 
```



# Multinomial model

- The multinomial NB (A) is identical to the multinomial unigram language model (B)
- A:  $P(c|d) \propto p(c) \prod_{1 \leq k \leq n_d} P(t_k|c)$
- B:  $P(d|q) \propto P(d) \prod_{t \in q} P(t|M_d)$
- The document in A takes the role of a query in B
- The class in A takes the role of a document in B



# The Bernoulli model

- A Bernoulli random variable takes one of two values
- A multivariate Bernoulli variable is a vector of Bernoulli variables which is used as an indicator
- When using in text classification, a Bernoulli model only considers the presence of a term in a class
- It ignores the number of occurrences
- It is equivalent to the binary independence model



# Bernoulli model

$$P(c|d) \propto P(c)P(d|c)$$

$$\begin{aligned} P(d = [e_1, \dots, e_M] | c; V) &= \prod_{e \in V} P(e|c) \\ &= \prod_{e \in V_d} P(e|c) \prod_{e \notin V_d} (1 - P(e|c)) \end{aligned}$$

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

$$P(d|c) = \frac{N_{ct} + 1}{N_c + 2}$$

- $N_c$  is the number of documents in class  $c$
- $N_{ct}$  is the number of document in class  $c$  containing term  $t$
- $N$  is the total number of documents



# Bernoulli model

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

```
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  $V, \text{prior}, \text{condprob}$ 
```

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

```
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  $\arg \max_{c \in \mathbb{C}} \text{score}[c]$ 
```





# Bernoulli model

	Doc ID	Words in documents	class
Training set	1	Chinese Beijing Chinese	1
	2	Chinese Chinese Shanghai	1
	3	Chinese Macao	1
	4	Tokyo Japan Chinese	0
Test set	5	Chinese Chinese Chinese Tokyo Japan	?



# Multinomial vs Bernoulli model

	Multinomial model	Bernoulli model
Event model	Generation of token	Generation of document
Random variable	$X=t$ iff $t$ occurs at given position	$U_t = 1$ iff $t$ occurs in the document
Document representation	$d = \langle t_1, \dots, t_k, \dots, t_{nd} \rangle, t_k \in V$	$d = \langle e, \dots, e_i, \dots, e_M \rangle, e_i \in \{0,1\}$
Parameter estimation	$\hat{P}(X = t c)$	$\hat{P}(U_i = e c)$
Multiple occurrences	Taken into account	Ignored
Length of documents	Can handle longer documents	Works best for short documents
Number of features	Can handle more	Works best with fewer
Estimate for term 'the'	$\hat{P}(X = 'the' c) \approx 0.05$	$\hat{P}(U_{the} = 1 c) \approx 1.0$



# Summary

- Documents can be classified into multiple classes
- A text classifier is a function from document space to a label set
- A naïve Bayes classifier can be used for document classification
- It is similar to a language model information retrieval
- A Bernoulli model can also be used for document classification