

UNIVERSITY of DELAWARE

## Text Categorization – generative methods

Hui Fang  
Department of Electrical and Computer Engineering  
University of Delaware

1

UNIVERSITY of DELAWARE

### Problem Setup of Text Categorization

- Input:
  - a document  $d$
  - a fixed set of categories  $C = \{c_1, c_2, \dots, c_J\}$

A training set of  $m$  hand-labeled documents  $(d_1, c_1), \dots, (d_m, c_m)$

a learned classifier:  $f(d) \rightarrow c$

- Output: a predicted category  $c \in C$  for the document  $d$

2

UNIVERSITY of DELAWARE

### Generative Classification Method

- Input:
  - a document  $d$
  - a fixed set of categories  $C = \{c_1, c_2, \dots, c_J\}$
  - a training set of  $m$  labeled documents  $(d_1, c_1), \dots, (d_m, c_m)$
- Training
  - For each class, we estimate a probability model
- Testing
  - Given the document  $d$ , we select the category  $c \in C$  which is most likely used to generate  $d$ .

3

UNIVERSITY of DELAWARE

### Naïve Bayes Classification

5

UNIVERSITY of DELAWARE

### Naïve Bayes Intuition

- A simple classification method based on Bayes rule
- Relies on very simple representation of document
  - Bag of words

UNIVERSITY of DELAWARE

### Text Classification with Naive Bayes

- Given a document  $d$ , we need to find  $c$  that is most likely used to generate  $d$ .

$$c^* = \arg \max_{c_i \in C} P(c_i | d) \quad \text{Posterior probability of } c_i$$

$$= \arg \max_{c_i \in C} \frac{P(d | c_i)P(c_i)}{P(d)}$$

$$= \arg \max_{c_i \in C} P(d | c_i)P(c_i)$$

Document model for category  $c_i$

Prior probability of  $c_i$

UNIVERSITY of DELAWARE

### How to estimate $P(d|c)$ ?

- Naïve assumption
  - A document is generated by a unigram language model.
  - Each word is generated *independently* of the others.

$$P(d|c_i) = P(w_1, \dots, w_n | c_i) = \prod_{j=1}^n \underbrace{P(w_j | c_i)}$$

**How to estimate it?**

UNIVERSITY of DELAWARE

$$P(d|c_i) = P(w_1, \dots, w_n | c_i) = \prod_{j=1}^n P(w_j | c_i)$$

Fraction of times word w appears among all words in documents of topic c

$$P(w_j | c_i) = \frac{\sum_{k=1..m, class(d_k)=c_i} count(w_j, d_k)}{\sum_{w \in V} \sum_{k=1..m, class(d_k)=c_i} count(w, d_k)}$$

- Create a mega-document for topic c by concatenating all docs in this topic
- Use frequency of w in the mega-document.

UNIVERSITY of DELAWARE

$$P(d|c_i) = P(w_1, \dots, w_n | c_i) = \prod_{j=1}^n P(w_j | c_i)$$

Fraction of times word w appears among all words in documents of topic c

$$P(w_j | c_i) = \frac{\sum_{k=1..m, class(d_k)=c_i} count(w_j, d_k)}{\sum_{w \in V} \sum_{k=1..m, class(d_k)=c_i} count(w, d_k)}$$

This gives score of zero if d contains a unseen word

$$P(w_j | c_i) = \frac{\sum_{k=1..m, class(d_k)=c_i} count(w_j, d_k) + 1}{\sum_{w \in V} \sum_{k=1..m, class(d_k)=c_i} count(w, d_k) + |V|}$$

UNIVERSITY of DELAWARE

### Text Classification with Naive Bayes

- Given a document d, we need to find c that is most likely used to generate d.

$$c^* = \arg \max_{c_i \in C} P(c_i | d)$$

Posterior probability of  $c_i$

$$= \arg \max_{c_i \in C} \frac{P(d | c_i) P(c_i)}{P(d)}$$

$$= \arg \max_{c_i \in C} P(d | c_i) P(c_i)$$

Document model for category  $c_i$

Prior probability of  $c_i$

UNIVERSITY of DELAWARE

### Implementation of Naïve Bayes

- Calculate  $P(c)$   
For each c
$$P(c) \leftarrow \frac{\# \text{of training documents with label } c}{\# \text{of total training documents}}$$

- Calculate  $P(w_i | c)$ 
  - $Text(c) \leftarrow \text{concatenate all training documents with label } c$
  - For each word  $w_i$  in  $V$
  - $n_i \leftarrow \# \text{ of occurrences of } w_i \text{ in } Text(c)$
  - $n \leftarrow \# \text{ of all word occurrences in } Text(c)$
  - $P(w_i | c) \leftarrow \frac{n_i + 1}{n + |V|}$

UNIVERSITY of DELAWARE

### An Example for Naïve Bayes Classification

- Problem:
  - Decide whether a document talks about China
- Training set
  - Doc 1: Chinese Beijing Chinese; Yes
  - Doc 2 : Chinese Chinese Shanghai; Yes
  - Doc 3: Chinese Macao; Yes
  - Doc 4: Tokyo Japan Chinese; No
- Test set
  - Doc 5: Chinese Chinese Chinese Tokyo Japan; ?

**An Example (Cont.)**

- Prior:  $p(c)=3/4$ ,  $p(\text{not } c) = 1/4$
- Conditional probabilities
  - $P(\text{Chinese}|c) = (5+1)/(8+6) = 6/14=3/7$
  - $P(\text{Tokyo}|c) = P(\text{Japan}|c) = (0+1)/(8+6) = 1/14$
  - $P(\text{Chinese} | \text{not } c) = (1+1)/(3+6) = 2/9$
  - $P(\text{Tokyo}|\text{not } c) = P(\text{Japan}|\text{not } c) = (1+1)/(3+6) = 2/9$
- Posterior probabilities
  - $P(c|d5) = 3/4 * (3/7)^3 * 1/14 * 1/14 = 0.0003$
  - $P(\text{not } c | d5) = 1/4 * (2/9)^3 * 2/9 * 2/9 = 0.0001$
- Yes. Doc5 talks about China.

16

**Summary**

- Pros:
  - Efficient
  - easy-to-implement
  - Very good in domains with many equally important features
- Cons:
  - Seldom gives the very best performance