



## Text Classification and Naïve Bayes

### The Task of Text Classification



#### Male or female author?

1. By 1925 present-day Vietnam was divided into three parts under French colonial rule. The southern region embracing Saigon and the Mekong delta was the colony of Cochinchina; the central area with its imperial capital at Hue was the protectorate of Annam...
2. Clara never failed to be astonished by the extraordinary felicity of her own name. She found it hard to trust herself to the mercy of fate, which had managed over the years to convert her greatest shame into one of her greatest assets...

S. Argamon, M. Koppel, J. Fine, A. R. Shimoni. 2003. "Gender, Genre, and Writing Style in Formal Written Texts," *Text*, volume 23, number 3, pp. 321–346.



#### Text Classification

- Assigning subject categories, topics, or genres
- Spam detection
- Authorship identification
- Age/gender identification
- Language Identification
- Sentiment analysis
- ...



#### Is this spam?

Subject: Important notice!  
From: Stanford University <newsforum@stanford.edu>  
Date: October 28, 2011 12:34:16 PM PDT  
To: undisclosed-recipients:;

Greats News!

You can now access the latest news by using the link below to login to Stanford University News Forum.  
<http://www.123contactform.com/contact-form-StanfordNew1-236335.html>

Click on the above link to login for more information about this new exciting forum. You can also copy the above link to your browser bar and login for more information about the new services.

© Stanford University. All Rights Reserved.



#### Who wrote which Federalist papers?

- 1787-8: anonymous essays try to convince New York to ratify U.S Constitution: Jay, Madison, Hamilton.
- Authorship of 12 of the letters in dispute
- 1963: solved by Mosteller and Wallace using Bayesian methods



James Madison

Alexander Hamilton



#### Positive or negative movie review?



- unbelievably disappointing
- Full of zany characters and richly applied satire, and some great plot twists
- this is the greatest screwball comedy ever filmed
- It was pathetic. The worst part about it was the boxing scenes.

5



#### What is the subject of this article?



6

- MeSH Subject Category Hierarchy**
- Antagonists and Inhibitors
  - Blood Supply
  - Chemistry
  - Drug Therapy
  - Embryology
  - Epidemiology
  - ...



#### Text Classification: definition

- **Input:**
  - a document  $d$
  - a fixed set of classes  $C = \{c_1, c_2, \dots, c_J\}$
- **Output:** a predicted class  $c \in C$



#### Classification Methods: Hand-coded rules

- Rules based on combinations of words or other features
  - spam: black-list-address OR ("dollars" AND "have been selected")
- Accuracy can be high
  - If rules carefully refined by expert
- But building and maintaining these rules is expensive



## Classification Methods: Supervised Machine Learning

- Input:
  - a document  $d$
  - a fixed set of classes  $C = \{c_1, c_2, \dots, c_j\}$
  - A training set of  $m$  hand-labeled documents  $(d_1, c_1), \dots, (d_m, c_m)$
- Output:
  - a learned classifier  $y: d \rightarrow c$

10



## Classification Methods: Supervised Machine Learning

- Any kind of classifier
  - Naïve Bayes
  - Logistic regression
  - Support-vector machines
  - k-Nearest Neighbors
  - ...



## Text Classification and Naïve Bayes

### Naïve Bayes (I)



## Naïve Bayes Intuition

- Simple ("naïve") classification method based on Bayes rule
- Relies on very simple representation of document
  - Bag of words



## The bag of words representation

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet.

$$Y( \quad ) = C$$


## The bag of words representation: using a subset of words



```
x love XXXXXXXXXXXXXXXX sweet
XXXXXX satirical XXXXXXXXXX
XXXXXXX great XXXXXX
XXXXXXXXXXXXXX fun XXXX
XXXXXXXXXXXXXX whimsical XXXX
XXXXXXXXXXXXXX recommend XXXX
XXXXXXXXXXXXXX several XXXXXXXXXX
XX several XXXXXXXXXX
XXXX happy XXXXXXXX again
XXXXXXXXXXXXXX XXXXXXXXXX
XXXXXXXXXXXXXX
```

$$Y( \quad ) = C$$


## The bag of words representation

great	2
love	2
recommend	1
laugh	1
happy	1
...	...

$$Y( \quad ) = C$$


## Bag of words for document classification



Test document  
parser  
language  
label  
translation  
...

Machine Learning  
learning  
training  
algorithm  
shrinkage  
network...

NLP  
parser  
tag  
algorithm  
translation  
language...

Garbage Collection  
garbage  
collection  
memory  
optimization  
plan  
temporal  
reasoning  
region...

Planning  
planning  
temporal  
reasoning  
optimization  
plan  
language...

GUI

?



## Text Classification and Naïve Bayes

### Formalizing the Naïve Bayes Classifier



#### Naïve Bayes Classifier (II)

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(d | c)P(c)$$

$$= \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c)P(c)$$

Document d represented as features  $x_1..x_n$



#### Bayes' Rule Applied to Documents and Classes

- For a document  $d$  and a class  $c$

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



#### Naïve Bayes Classifier (I)

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(c | d)$$

$$= \operatorname{argmax}_{c \in C} \frac{P(d | c)P(c)}{P(d)}$$

$$= \operatorname{argmax}_{c \in C} P(d | c)P(c)$$

MAP is "maximum a posteriori" = most likely class

Bayes Rule

Dropping the denominator



#### Multinomial Naïve Bayes Classifier

$$c_{MAP} = \operatorname{argmax} P(x_1, x_2, \dots, x_n | c)P(c)$$

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c_j) \prod_{x \in X} P(x | c)$$



#### Applying Multinomial Naïve Bayes Classifiers to Text Classification

positions  $\leftarrow$  all word positions in test document

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i \in \text{positions}} P(x_i | c_j)$$



## Text Classification and Naïve Bayes

Naïve Bayes:  
Learning



## Learning the Multinomial Naïve Bayes Model

- First attempt: maximum likelihood estimates
  - simply use the frequencies in the data

$$\hat{P}(c_j) = \frac{\text{doccount}(C=c_j)}{N_{\text{doc}}}$$

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$



## Parameter estimation

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)} \quad \text{fraction of times word } w_i \text{ appears among all words in documents of topic } c_j$$

- Create mega-document for topic  $j$  by concatenating all docs in this topic
- Use frequency of  $w$  in mega-document



## Problem with Maximum Likelihood

- What if we have seen no training documents with the word **fantastic** and classified in the topic **positive (thumbs-up)**?

$$\hat{P}(\text{"fantastic"} | \text{positive}) = \frac{\text{count}(\text{"fantastic"}, \text{positive})}{\sum_{w \in V} \text{count}(w, \text{positive})} = 0$$

- Zero probabilities cannot be conditioned away, no matter the other evidence!

$$c_{MAP} = \operatorname{argmax}_c \hat{P}(c) \prod_i \hat{P}(x_i | c)$$



## Laplace (add-1) smoothing for Naïve Bayes

$$\begin{aligned} \hat{P}(w_i | c) &= \frac{\text{count}(w_i, c) + 1}{\sum_{w \in V} (\text{count}(w, c) + 1)} \\ &= \frac{\text{count}(w_i, c) + 1}{\left( \sum_{w \in V} \text{count}(w, c) \right) + |V|} \end{aligned}$$



## Multinomial Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*
- Calculate  $P(c_j)$  terms
  - For each  $c_j$  in  $C$  do
    - $\text{docs}_{c_j} \leftarrow$  all docs with class  $= c_j$
    - $P(c_j) \leftarrow \frac{|\text{docs}_{c_j}|}{\text{total # documents}}$
  - Calculate  $P(w_k | c_j)$  terms
    - $\text{Text}_{c_j} \leftarrow$  single doc containing all  $\text{docs}_{c_j}$
    - For each word  $w_k$  in *Vocabulary*
      - $n_k \leftarrow$  # of occurrences of  $w_k$  in  $\text{Text}_{c_j}$
      - $P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha |Vocabulary|}$



## Laplace (add-1) smoothing: unknown words

Add one extra word to the vocabulary, the "unknown word"  $w_u$

$$\begin{aligned} \hat{P}(w_u | c) &= \frac{\text{count}(w_u, c) + 1}{\left( \sum_{w \in V} \text{count}(w, c) \right) + |V + 1|} \\ &= \frac{1}{\left( \sum_{w \in V} \text{count}(w, c) \right) + |V + 1|} \end{aligned}$$



## Text Classification and Naïve Bayes

### Multinomial Naïve Bayes: A Worked Example



$$\begin{aligned} \hat{P}(c) &= \frac{N_c}{N} \\ \hat{P}(w | c) &= \frac{\text{count}(w, c) + 1}{\text{count}(c) + |V|} \end{aligned}$$

	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:

$$P(c | d5) \propto \frac{3}{4} * \frac{1}{3} * \frac{1}{14} * \frac{1}{14} = 0.0003$$

#### Conditional Probabilities:

$$\begin{aligned} P(\text{Chinese} | c) &= (5+1) / (8+6) = 6/14 = 3/7 \\ P(\text{Tokyo} | c) &= (0+1) / (8+6) = 1/14 \\ P(\text{Japan} | c) &= (0+1) / (8+6) = 1/14 \\ P(\text{Chinese} | j) &= (1+1) / (3+6) = 2/9 \\ P(\text{Tokyo} | j) &= (1+1) / (3+6) = 2/9 \\ P(\text{Japan} | j) &= (1+1) / (3+6) = 2/9 \end{aligned}$$

45



## Naïve Bayes in Spam Filtering

#### SpamAssassin Features:

- Mentions Generic Viagra
- Online Pharmacy
- Mentions millions of (dollar) ((dollar) NN,NNN,NNN.NN)
- Phrase: impress ... girl
- From: starts with many numbers
- Subject is all capitals
- HTML has a low ratio of text to image area
- One hundred percent guaranteed
- Claims you can be removed from the list
- Prestigious Non-Accredited Universities'
- [http://spamassassin.apache.org/tests\\_3\\_3\\_x.html](http://spamassassin.apache.org/tests_3_3_x.html)



## Summary: Naive Bayes is Not So Naive

- Very Fast, low storage requirements
- Robust to Irrelevant Features  
Irrelevant Features cancel each other without affecting results
- Very good in domains with many equally important features  
Decision Trees suffer from *fragmentation* in such cases – especially if little data
- Optimal if the independence assumptions hold: if assumed independence is correct, then it is the Bayes Optimal Classifier for problem
- A good dependable baseline for text classification
  - **But we will see other classifiers that give better accuracy**



## Precision and recall

- **Precision:** % of selected items that are correct
- **Recall:** % of correct items that are selected

	correct	not correct
selected	tp	fp
not selected	fn	tn



## A combined measure: F

- A combined measure that assesses the P/R tradeoff is F measure (weighted harmonic mean):

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

- The harmonic mean is a very conservative average; see *IIR* § 8.3
- People usually use balanced F1 measure
  - i.e., with  $\beta = 1$  (that is,  $\alpha = \frac{1}{2}$ ):

$$F = 2PR/(P+R)$$



## Text Classification and Naïve Bayes

Precision, Recall, and the F measure



## The 2-by-2 contingency table

	correct	not correct
selected	tp	fp
not selected	fn	tn