



NATURAL LANGUAGE PROCESSING

Text Classification

TOPICS

- **Text Classification**
 - Naïve Bayes
 - Logistic Regression
 - Practical Issues



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.

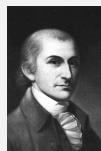
© Stanford University. All Rights Reserved.

YanJun Li

3

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



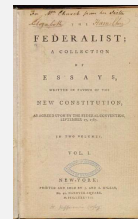
John Jay



James Madison



Alexander
Hamilton



YanJun Li

4

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. Shmuni, 2003. "Gender, Genre, and Writing Style in Formal Written Texts," *Text*, volume 23, number 3, pp. 321–346

YanJun Li

5

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.

YanJun Li

6

WHAT IS THE SUBJECT OF THIS ARTICLE?

MeSH Subject Category Hierarchy



- Yanjun Li

7

TEXT CLASSIFICATION

- Yanjun Li

8

TEXT CLASSIFICATION: DEFINITION

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

YanJun Li

9

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

YanJun Li

10

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 $\gamma: d \rightarrow c$

YanJun Li

11

SUPERVISED MACHINE LEARNING

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

YanJun Li

12

TOPICS

- Text Classification
 - **Naïve Bayes**
 - Logistic Regression
 - Practical Issues

YanJun Li



13

NAÏVE BAYES INTUITION

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

YanJun Li

14

THE BAG OF WORDS REPRESENTATION

$Y(\text{seen sweet whimsical recommend happy} \dots) = c$

seen	2
sweet	1
whimsical	1
recommend	1
happy	1
...	...



YanJun Li

15

BAYES' RULE APPLIED TO DOCUMENTS AND CLASSES

- For a document d and a class c

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

YanJun Li

16

NAÏVE BAYES CLASSIFIER (I)

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

MAP is “maximum a posteriori” = most likely class

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

Bayes Rule

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

Dropping the denominator

YanJun Li

17

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 $x_1 \dots x_n$

YanJun Li

18

NAÏVE BAYES CLASSIFIER (IV)

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

How often does this class occur?

Could only be estimated if a very, very large number of training examples was available.

We can just count the relative frequencies in a corpus

YanJun Li

19

MULTINOMIAL NAÏVE BAYES INDEPENDENCE ASSUMPTIONS

$$P(x_1, x_2, \dots, x_n | c)$$

- **Bag of Words assumption:** Assume position doesn't matter
- **Conditional Independence:** Assume the feature probabilities $P(x_i | c_j)$ are independent given the class c .

$$P(x_1, \dots, x_n | c) = P(x_1 | c) \cdot P(x_2 | c) \cdot P(x_3 | c) \cdot \dots \cdot P(x_n | c)$$

YanJun Li

20

MULTINOMIAL NAÏVE BAYES CLASSIFIER

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

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

YanJun Li

21

APPLYING MULTINOMIAL NAIVE BAYES CLASSIFIERS TO TEXT CLASSIFICATION

positions \leftarrow all word positions in test document

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

YanJun Li

22

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

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

**fraction of times word w_i 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

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

UNKNOWN WORD

- The solution for unknown words is to ignore them.
- Remove them from the test document and not include any probability for them at all in the prediction procedure.

YanJun Li

27

MULTINOMIAL NAÏVE BAYES: LEARNING

- From training corpus, extract *Vocabulary*
- Calculate $P(c_j)$ terms
 - For each c_j in C do
 $docs_j \leftarrow$ all docs with class $= c_j$
- Calculate $P(w_k | c_j)$ terms
 - $Text_j \leftarrow$ single doc containing all $docs_j$
 - For each word w_k in *Vocabulary*
 $n_k \leftarrow$ # of occurrences of w_k in $Text_j$

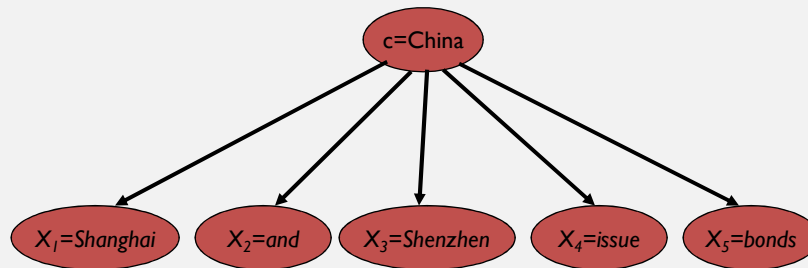
$$P(c_j) \leftarrow \frac{|docs_j|}{|\text{total \# documents}|}$$

$$P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha |Vocabulary|}$$

YanJun Li

28

GENERATIVE MODEL FOR MULTINOMIAL NAÏVE BAYES



YanJun Li

29

NAÏVE BAYES AND LANGUAGE MODELING

- Naïve bayes classifiers can use any sort of feature
 - URL, email address, dictionaries, network features
- But if, as in the previous slides
 - We use **only** word features
 - we use **all** of the words in the text (not a subset)
- Then
 - Naïve bayes has an important similarity to language modeling.

YanJun Li

30

EACH CLASS = A UNIGRAM LANGUAGE MODEL

Sec.13.2.1

- Assigning each word: $P(\text{word} | c)$
- Assigning each sentence: $P(s|c) = \prod P(\text{word}|c)$

Class *pos*

0.1 I

0.1 love

0.01 this

0.05 fun

0.1 film

...

$$P(s | \text{pos}) = 0.0000005$$

YanJun Li

31

NAÏVE BAYES AS A LANGUAGE MODEL

Sec.13.2.1

- Which class assigns the higher probability to *s*?

Model *pos*

0.1 I

0.1 love

0.01 this

0.05 fun

0.1 film

Model *neg*

0.2 I

0.001 love

0.01 this

0.005 fun

0.1 film

I	love	this	fun	film
0.1	0.1	0.01	0.05	0.1
0.2	0.001	0.01	0.005	0.1

$$P(s|\text{pos}) > P(s|\text{neg})$$

YanJun Li

32

Example

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

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

Prior

s:

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

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

Conditional Probabilities:

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

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

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

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

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

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

	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	?

Choosing a class:

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

$$\approx 0.0003$$

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

$$\approx 0.0001$$

YanJun Li

33

SUMMARY: NAIVE BAYES IS NOT SO NAIVE

- Very Fast, low storage requirements
- Very good in domains with many equally important features
- Optimal if the independence assumptions hold: If assumed independence is correct
- A good dependable baseline for text classification
 - **But we will see other classifiers that give better accuracy**

YanJun Li

34

EVALUATION OF NB

		gold standard labels		
		gold positive	gold negative	
system output labels	system positive	true positive	false positive	precision = $\frac{tp}{tp+fp}$
	system negative	false negative	true negative	
		recall = $\frac{tp}{tp+fn}$		accuracy = $\frac{tp+tn}{tp+fp+tn+fn}$

Figure 4.4 Contingency table

YanJun Li

35

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, $\alpha = 1/2$): $F = \frac{2PR}{P+R}$

YanJun Li

36

MORE THAN TWO CLASSES: SETS OF BINARY CLASSIFIERS

- Dealing with **any-of** or **multivalued** classification
 - A document can belong to 0, 1, or >1 classes.
- For each class $c \in C$
 - Build a classifier γ_c to distinguish c from all other classes $c' \in C$
- Given test doc d ,
 - Evaluate it for membership in each class using each γ_c
 - d belongs to **any** class for which γ_c returns true

MORE THAN TWO CLASSES: SETS OF BINARY CLASSIFIERS

- **One-of** or **multinomial** classification (multiclass)
 - Classes are mutually exclusive: each document in exactly one class
- For each class $c \in C$
 - Build a classifier γ_c to distinguish c from all other classes $c' \in C$
- Given test doc d ,
 - Evaluate it for membership in each class using each γ_c
 - d belongs to the **one** class with maximum score

EVALUATION: CLASSIC REUTERS-21578 DATA SET

- Most (over)used data set, 21,578 docs (each 90 types, 200 tokens)
- 9603 training, 3299 test articles (ModApte/Lewis split)
- 118 categories
 - An article can be in more than one category
 - Learn 118 binary category distinctions
- Average document (with at least one category) has 1.24 classes
- Only about 10 out of 118 categories are large

Common categories
(#train, #test)

- | | |
|----------------------------|-----------------------|
| • Earn (2877, 1087) | • Trade (369,119) |
| • Acquisitions (1650, 179) | • Interest (347, 131) |
| • Money-fx (538, 179) | • Ship (197, 89) |
| • Grain (433, 149) | • Wheat (212, 71) |
| • Crude (389, 189) | • Corn (182, 56) |

YanJun Li

39

REUTERS TEXT CATEGORIZATION DATA SET (REUTERS-21578) DOCUMENT

```
<REUTERS TOPICS="YES" LEWISSPLIT="TRAIN" CGISPLIT="TRAINING-SET" OLDID="12981"
NEWID="798">
```

```
<DATE> 2-MAR-1987 16:51:43.42</DATE>
```

```
<TOPICS><D>livestock</D><D>hog</D></TOPICS>
```

```
<TITLE>AMERICAN PORK CONGRESS KICKS OFF TOMORROW</TITLE>
```

```
<DATELINE> CHICAGO, March 2 - </DATELINE><BODY>The American Pork Congress kicks off tomorrow,
March 3, in Indianapolis with 160 of the nations pork producers from 44 member states determining industry positions
on a number of issues, according to the National Pork Producers Council, NPPC.
```

Delegates to the three day Congress will be considering 26 resolutions concerning various issues, including the future direction of farm policy and the tax law as it applies to the agriculture sector. The delegates will also debate whether to endorse concepts of a national PRV (pseudorabies virus) control and eradication program, the NPPC said.

A large trade show, in conjunction with the congress, will feature the latest in technology in all areas of the industry, the NPPC added. Reuter

```
&#3;</BODY></TEXT></REUTERS>
```

YanJun Li

40

CONFUSION MATRIX

- For each pair of classes $\langle c_1, c_2 \rangle$ how many documents from c_1 were incorrectly assigned to c_2 ?
- $c_{3,2}$: 90 wheat documents incorrectly assigned to poultry

Docs in test set	Assigned UK	Assigned poultry	Assigned wheat	Assigned coffee	Assigned interest	Assigned trade
True UK	95	1	13	0	1	0
True poultry	0	1	0	0	0	0
True wheat	10	90	0	1	0	0
True coffee	0	0	0	34	3	7
True interest	-	1	2	13	26	5
True trade	0	0	2	14	5	10

YanJun Li

41

PER CLASS EVALUATION MEASURES

Recall:

Fraction of docs in class i classified correctly:

$$\frac{c_{ii}}{\sum_j c_{ij}}$$

Precision:

Fraction of docs assigned class i that are actually about class i :

$$\frac{c_{ii}}{\sum_j c_{ji}}$$

Accuracy: (1 - error rate)

Fraction of docs classified correctly:

$$\frac{\sum_i c_{ii}}{\sum_j \sum_i c_{ij}}$$

YanJun Li

42

MICRO- VS. MACRO-AVERAGING

- If we have more than one class, how do we combine multiple performance measures into one quantity?
- **Macro-averaging:** Compute performance for each class, then average.
- **Micro-averaging:** Collect decisions for all classes, compute contingency table, evaluate.

MICRO- VS. MACRO-AVERAGING: EXAMPLE

Class 1			Class 2			Micro Ave. Table		
	Truth: yes	Truth: no		Truth: yes	Truth: no		Truth: yes	Truth: no
Classifier: yes	10	10	Classifier: yes	90	10	Classifier: yes	100	20
Classifier: no	10	970	Classifier: no	10	890	Classifier: no	20	1860

- Macro-averaged precision: $(0.5 + 0.9)/2 = 0.7$
- Micro-averaged precision: $100/120 = .83$
- Micro-averaged score is dominated by score on common classes

DEVELOPMENT TEST SETS AND CROSS-VALIDATION

Training set

Development Test Set

Test Set

- Metric: P/R/F1 or Accuracy
- Unseen test set
 - avoid overfitting ('tuning to the test set')
 - more conservative estimate of performance
- Cross-validation over multiple splits
 - Handle sampling errors from different datasets
 - Pool results over each split
 - Compute pooled dev set performance

Training Set Dev Test

Training Set Dev Test

Dev Test Training Set

Test Set

YanJun Li

45