

Natural Language Processing

Lecture 11—9/29/2015

Jim Martin

Today

- Text classification
 - Tasks
 - Naïve Bayes

9/29/15 Speech and Language Processing - Jurafsky and Martin 2

Is this spam?

Dr.fun
To: jeinaddy@gmail.com
Collect your welcome treat

September 12, 2015 10:54 PM [Hide Details](#)

Hey there,
Hello

A 200% Match Bonus to help triple your first deposit!
That's what you'll score if you simply sign up at Ruby Fortune on your mobile NOW!
<http://rfsb.com/procedures.htm>

Plus, there's games galore, awesome promos, and excellent prizes too.
Have fun

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







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 Cochin-China; 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. Shmioni, 2003. "Gender, Genre, and Writing Style in Formal Written Texts," Text, volume 23, number 3, pp. 321-346

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.

6

What is the subject of this article?

MEDLINE Article



MeSH Subject Category Hierarchy



- Antagonists and Inhibitors
- Blood Supply
- Chemistry
- Drug Therapy
- Embryology
- Epidemiology
- ...

7

Text Classification

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

Text Classification

- **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

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 $y: d \rightarrow c$

11

Supervised Machine Learning

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

Supervised Machine Learning

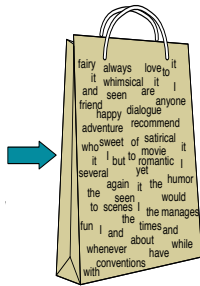
- Any kind of classifier
 - Naïve Bayes
 - Logistic regression
 - Support-vector machines
 - k-Nearest Neighbors
 - ...

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!



it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

Bag of Words Classifier

Y
(

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

) = C



Bayes' Rule

- For a document d and a class c

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

Naïve Bayes Classifier (I)

$$\begin{aligned} c_{MAP} &= \operatorname{argmax}_{c \in C} P(c|d) && \text{MAP is "maximum a posteriori" = most likely class} \\ &= \operatorname{argmax}_{c \in C} \frac{P(d|c)P(c)}{P(d)} && \text{Bayes Rule} \\ &= \operatorname{argmax}_{c \in C} P(d|c)P(c) && \text{Dropping the denominator} \end{aligned}$$

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 $x_1 \dots x_n$

Naïve Bayes Classifier (IV)

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

$O(|X|^n \cdot |C|)$ parameters

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

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 $P(x_i | c)$ are independent given the class c .

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

Multinomial Naïve Bayes Classifier

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

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

Applying Naive Bayes to Text Classification

positions ← all word positions in test document

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

Learning the Naïve Bayes Model

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

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

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

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 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
 - $docs_j \leftarrow$ 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 $Text_j$

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

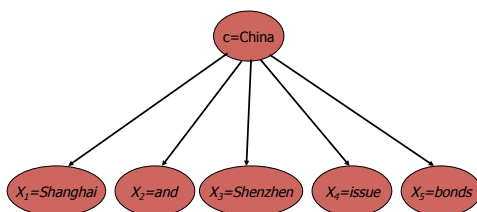
$$P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha |Vocabulary|}$$

Naïve Bayes and Language Modeling

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

29

Generative Model for Naïve Bayes



30

Each class = a unigram language model

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

Class pos			I	love	this	fun	film
0.1	I						
0.1	love						
0.01	this	0.1	0.1	.05	0.01	0.1	
0.05	fun						
0.1	film						
...							

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

Naïve Bayes as a Language Model

Sec. 13.2.1

- Which class assigns the higher probability to s?

Model pos		Model neg			I	love	this	fun	film
0.1	I	0.2	I						
0.1	love	0.001	love						
0.01	this	0.01	this	0.1	0.1	0.01	0.05	0.1	
0.05	fun	0.005	fun	0.2	0.001	0.01	0.005	0.1	
0.1	film	0.1	film						

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

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

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

Priors:
 $P(c) =$
 $P(j) =$

Conditional Probabilities:
 $P(\text{Chinese}|c) = \frac{(5+1)}{(8+6)} = 6/14 = 3/7$
 $P(\text{Tokyo}|c) = \frac{(0+1)}{(8+6)} = 1/14$
 $P(\text{Japan}|c) = \frac{(0+1)}{(8+6)} = 1/14$
 $P(\text{Chinese}|j) = \frac{(1+1)}{(3+6)} = 2/9$
 $P(\text{Tokyo}|j) = \frac{(1+1)}{(3+6)} = 2/9$
 $P(\text{Japan}|j) = \frac{(1+1)}{(3+6)} = 2/9$

	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:
 $P(c|d5) \propto 3/4 * (3/7)^3 * 1/14 * 1/14 = 0.0003$

$P(j|d5) \propto 1/4 * (2/9)^2 * 2/9 * 2/9 = 0.0001$

Practical Issues

- Preventing underflow
 - Multiplying lots of probabilities can result in floating-point underflow.
 - $\log(xy) = \log(x) + \log(y)$
 - Sum logs of probabilities instead of multiplying probabilities.
 - Class with highest un-normalized log probability score is still most probable.

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

- Model is now just max of sum of weights

Practical Issues

- What about new words in the test set?

	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 Boulder Tokyo Japan	?

9/29/15

Speech and Language Processing - Jurafsky and Martin

35

Machine Learning

- This model is now essentially a sum of weighted features + a bias term

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

- Moving on we can move away from just word-based features, and find better ways to set the weights

9/29/15

Speech and Language Processing - Jurafsky and Martin

36
