

Text Classification and Naïve Bayes

These slides are based on:

Dan Jurafsky and James H. Martin, Speech and Language Processing (3rd ed. draft)
<https://web.stanford.edu/~jurafsky/slp3/>
(Chapter 7)

These slides are an edited version of Jurafsky's slides:
<https://web.stanford.edu/~jurafsky/NLPCourseraSlides.html>

Examples of text classification problems

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.

Is this spam?

From Jack <huixinsoft17@foxmail.com>
Subject **lab furniture info from Kerric**
To [REDACTED]

Reply Reply All Forward Archive Junk Delete More
05:48

Dear My Friends,

Good day!

Kerric Laboratory Equipment Research & Development Manufacturer Co.,Ltd is established since 1999, is the leading manufacturer & supplier of the LABORATORY FURNITURE and relevant accessory for school, college, university and chemical or biology industry.

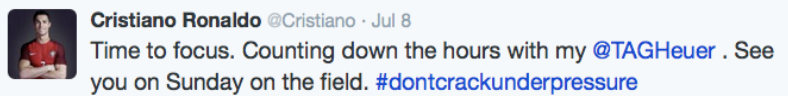
Our products Range Includes:
Laboratory Bench (All Steel, Steel Wood)
Laboratory Cabinet(All Steel, Aluminum Wood)
Fume Hood(All steel, PP)
Laboratory Rack/Shelf(All Steel, Steel Wood)

What is the topic of this post?



?

- Agriculture
- Robotics
- Sport
- Religion
- Psychology
- ...



Text Classification

Assigning subject categories:

- Sentiment analysis
- Spam detection
- Language identification
- Authorship identification
- Topic identification
- ...

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 and dictionary methods

- 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

Text Classification: Supervised Machine Learning

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 $\gamma: d \rightarrow c$

Classification Methods:

Supervised Machine Learning

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

The bag of words (BOW)

Y (

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.

) = C



BOW: important words

Y (

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.

) = C

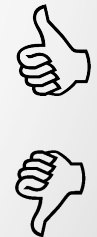


BOW: using a subset of words

Y (

```
x love xxxxxxxxxxxxxxxxxxxx sweet
xxxxxxxx satirical xxxxxxxxxxxx
xxxxxxxxxxxx great xxxxxxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxx fun xxxx
xxxxxxxxxxxxxxxxxxxx whimsical xxxx
romantic xxxx laughing
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxxxxxxxxxxxxxxxxx recommend xxxxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xx several xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxx happy xxxxxxxxxxxx again
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
```

) = C





BOW: word counts

Y (

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

) = C

Text Classification: Naïve Bayes

Naïve Bayes Intuition

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

Bayes' Rule applied to documents and classes

For a document d and a class c

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

likelihood

prior

posterior

Bayes theorem

Naïve Bayes Classifier

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

MAP is “maximum a posteriori” = most likely class

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

Bayes theorem

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

Dropping the denominator

Naïve Bayes Classifier

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

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

Document d
represented as
features
 $x_1 \dots x_n$

Features could be:

- words (binary value)
- word counts
- word frequencies (tf)
- tf-idf

Naïve Bayes Classifier

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

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

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

How often does this class occur?

We can just count the relative frequencies in a corpus

Multinomial NB: Independence Assumptions

- **Bag of words assumption:** Assume that position of words doesn't matter (the **exchangeability** of random variables)

$$P(x_1, x_2, \dots, x_n | c) = P(x_{\delta(1)}, x_{\delta(2)}, \dots, x_{\delta(n)} | c)$$

- **Conditional Independence:** Assume the feature probabilities $P(x_i | c_j)$ are independent given the class c .

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

Multinomial Naïve Bayes Classifier

Features: counts

Likelihood: $P(\mathbf{x}) = \text{Multinomial}(\mathbf{x})$

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c) \prod_{i \in V} P(x_i | c)$$

Boolean Multinomial Naïve Bayes

Features: binarized counts (0/1 values)

Likelihood: $P(\mathbf{x}) = \text{Multinomial}(\mathbf{x})$

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c) \prod_{i \in V} P(x_i | c)$$

- Boolean (aka binarized) multinomial NB good for polarity prediction
- Different from Bernoulli Naïve Bayes classifier

Bernoulli Naïve Bayes Classifier

Features: binarized counts (0/1 values)

Likelihood: $P(\mathbf{x}) = \text{MultivariateBernoulli}(\mathbf{x})$

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c) \prod_{i \in V} (P(x_i | c))^{x_i} (1 - P(x_i | c))^{1-x_i}$$

Learning parameters of Naïve Bayes

Learning parameters of Multinomial NB

- 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 the 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*)?

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

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

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

Laplace (add-1) smoothing for Naïve Bayes

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

Additive
smoothing:

$$\hat{P}(w_i | c) = \frac{\text{count}(w_i, c) + \alpha}{\left(\sum_{w \in V} \text{count}(w, c) \right) + \alpha |V|}$$

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

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

- 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 $Text_j$

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

Summary: Naive Bayes surprisingly good

- 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**

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

**Evaluation metrics:
precision, recall,
accuracy,**

The 2-by-2 contingency table

Actual class:

		Actual class:	
		positive	negative
Predicted class:	positive	tp	fp
	negative	fn	tn

true negative



Precision and recall

- **Precision:** % of predicted positive items that are correctly classified
- **Recall:** % of actual positive items that are correctly classified
- **Accuracy:** % of all items that are correctly classified

$$P = \frac{tp}{tp + fp}$$

$$R = \frac{tp}{tp + fn}$$

$$A = \frac{tp + tn}{n + p}$$

Actual class:

	positive	negative
Predicted class: positive	tp	fp
negative	fn	tn

Confusion matrix c_{ij}

		Predicted		
		Cat	Dog	Rabbit
Actual class	Cat	5	3	0
	Dog	2	3	1
	Rabbit	0	2	11

$c_{33} = 11$
i.e., 11 rabbits correctly
classified as rabbits

$c_{32} = 2$
i.e., 2 rabbits incorrectly
classified as dogs

Per class evaluation measures

Recall:

Fraction of docs in class i classified correctly:

$$R = \frac{C_{ii}}{\sum_j C_{ij}}$$

		Predicted		
		Cat	Dog	Rabbit
Actual class	Cat	5	3	0
	Dog	2	3	1
	Rabbit	0	2	11

Precision:

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

$$P = \frac{C_{ii}}{\sum_j C_{ji}}$$

		Predicted		
		Cat	Dog	Rabbit
Actual class	Cat	5	3	0
	Dog	2	3	1
	Rabbit	0	2	11

Accuracy: (1 - error rate)

Fraction of docs classified correctly:

$$Accuracy = \frac{\sum_i C_{ii}}{\sum_j \sum_i C_{ij}}$$

		Predicted		
		Cat	Dog	Rabbit
Actual class	Cat	5	3	0
	Dog	2	3	1
	Rabbit	0	2	11

A combined measure: F

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

$$F = \frac{1}{\alpha \frac{1}{P} + (1 - \alpha) \frac{1}{R}}$$

- The harmonic mean is a very conservative average
- We typically use balanced F measure with $\alpha = 0.5$
 - Namely, $F_1 = 2PR/(P+R)$

Cross-validation

Data:

Training set

Development Test
(Validation) Set

Test Set

- Score 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
 - compute score metric for each split
 - average score over all splits
- Grid search over hyperparameters
 - repeat previous step for various values of hyperparameters to find the best ones
 - choose hyperparameters that maximize score

CV, 3 splits (folds):

Training Set Dev Test

Training Set Dev Test

Dev Test Training Set

Test Set

Text Classification and Naïve Bayes

These slides are based on:

Dan Jurafsky and James H. Martin, Speech and Language Processing (3rd ed. draft)
<https://web.stanford.edu/~jurafsky/slp3/>
(Chapter 7)

These slides are an edited version of Jurafsky's slides:
<https://web.stanford.edu/~jurafsky/NLPCourseraSlides.html>