

# **Naive Bayes Classification (and Sentiment)**

Chapter 4

# **Text Classification**

## Is this spam?

RE: INVESTMENT/BUSINESS PARTNERSHIP.

DOES YOUR BUSINESS/PROJECT STILL NEED FUNDING?

DO YOU HAVE PROJECT/INVESTMENT CAPABLE OF GENERATING 15% AROI?

IF THE ANSWER IS YES, I represent a group of company based in the Gulf Region. We are seeking means of expanding and relocating our business interest abroad in the following sector, Oil & Gas, Energy, Mining, Construction, Real Estate, Communication, Agriculture, Health Sector or any other VIABLE business/project capable of generating 15% AROI.

If you have a solid background and idea of making good profit in any of the following SECTORS, please write me for possible business co-operation. More so, we are ready to facilitate and fund any business that is capable of generating 15% Annual Return on Investment (AROI) Joint Venture partnership and Hard loan funding can also be considered.

I look forward to discussing this opportunity further with you.

## Who wrote which Federalist papers?

- 1787-8: anonymous essays try to convince New York to ratify U.S Constitution.
- 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

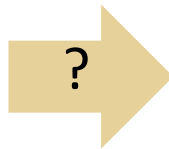
## What is the subject of this article?

### MEDLINE Article



### MeSH Subject Category Hierarchy

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



6

## Text Classification

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

## 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  $\gamma: d \rightarrow c$

## **Classification Methods: Supervised Machine Learning**

- Any kind of classifier
  - Naive Bayes
  - Logistic regression
  - Support-vector machines
  - k-Nearest Neighbors
  - Neural Nets
- ...

## **Naive Bayes**

Intuition and  
Formalization





## The bag of words representation

$Y$  (

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

) =  $C$

## Bayes' Rule Applied to Documents and Classes

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

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



## Naive 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 Rule

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

Dropping the denominator

## Naive 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  $x_1 \dots x_n$

## Naive Bayes Classifier

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

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

## Naive 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  $c$ .

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

## Naive Bayes Classifier

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

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

## Applying Naive Bayes to Text Classification

positions ← all word positions in test document

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

# Naive Bayes

## Learning

### Learning the Naive 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 (add-1) smoothing for Naive Bayes

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

## Naive Bayes: Learning

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

# Naive Bayes

## Relationship to Language Modeling

### Generative vs Discriminative Classifiers

Naive Bayes is the prototypical generative classifier.

- It describes a probabilistic process – “generative story” for a text input  $X$
- But why model  $X$ ? It's always observed.

Discriminative models instead:

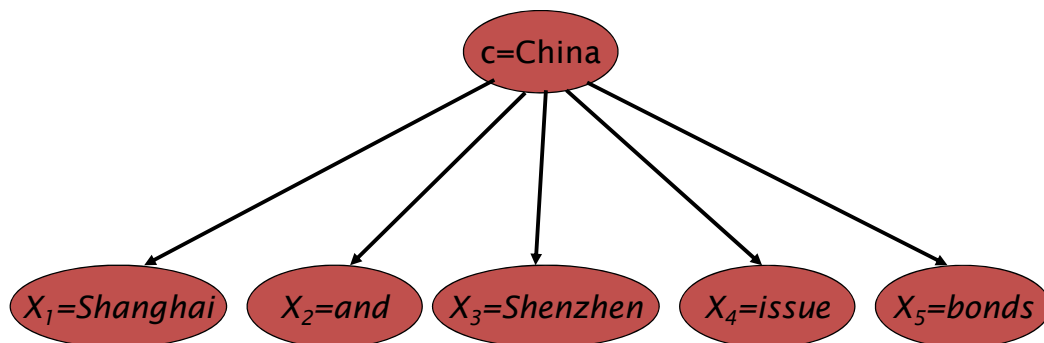
- seek to optimize a performance measure, like accuracy
- do not worry about  $p(X)$

## Generative vs Discriminative Classifiers

- Generative
  - Probabilistic
  - Specify a joint probability distribution over observations and targets:  $P(c,d)$
  - Bayes rule enables a conditional distribution
- Discriminative
  - Provide a model for the target variable
  - Use analysis of observed variables
  - Learn boundaries between classes
  - Infer outputs based on inputs:  $P(c|d)$

31

## Generative Model for Naive Bayes



32



## Naive Bayes and Language Modeling

- Naive 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
  - Naive bayes has an important similarity to language modeling.

33

## 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*

0.1	<u>I</u>					
0.1	<u>love</u>					
0.01	<u>this</u>	0.1	0.1	.01	.05	0.1
0.05	<u>fun</u>					
0.1	<u>film</u>					

$$P(s | pos) = 0.0000005$$

## Naïve Bayes as a Language Model

- Which class assigns the higher probability to s?

Model pos		Model neg			<u>l</u>	<u>love</u>	<u>this</u>	<u>fun</u>	<u>film</u>
0.1	l	0.2	l	0.1	0.1	0.01	0.05	0.1	
0.1	love	0.001	love	0.2	0.001	0.01	0.005	0.1	
0.01	this	0.01	this						
0.05	fun	0.005	fun						
0.1	film	0.1	film						

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

## Naive Bayes

### A Worked Example

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

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

**Priors:**

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

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

**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:**

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

$$\approx 0.0003$$
  

$$P(j|d5) \propto 1/4 * (2/9)^3 * 2/9 * 2/9$$

$$\approx 0.0001$$

# 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'

## Summary: Naive Bayes is Not So Naive

- Very fast, low storage requirements
- Robust to irrelevant features
- Very good in domains with many equally important features
- Optimal if the independence assumptions hold
- A good dependable baseline for text classification
  - **But often other classifiers give better accuracy**

## Naive Bayes

Evaluation: Precision,  
Recall, F-measure

## The 2-by-2 contingency table

	correct	not correct
selected	tp	fp
not selected	fn	tn

## 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 conservative average
- People usually use balanced F1 measure
  - i.e., with  $\beta = 1$  (that is,  $\alpha = \frac{1}{2}$ ):  $F = 2PR/(P+R)$

## Classification Methods: Review

- *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  $\gamma: d \rightarrow c$

## Naïve Bayes: Review

- What type of classifier?
- Two simplifying assumptions (one specific to text classification)
- Two types of probabilities

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

- Learning

45

## 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  $c \in \mathcal{C}$ 
  - Build a classifier  $\gamma_c$  to distinguish  $c$  from all other classes  $c' \in \mathcal{C}$
- Given test doc  $d$ ,
  - Evaluate it for membership in each class using each  $\gamma_c$
  - $d$  belongs to **any** class for which  $\gamma_c$  returns true

46

## More Than Two Classes: Sets of binary classifiers

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

47

## 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)

48

- |                            |                       |
|----------------------------|-----------------------|
| • 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)      |



## 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

<sup>49</sup>  
&#3;</BODY></TEXT></REUTERS>

## Confusion matrix c

- 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

50

## 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:

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

51

## Micro- vs. Macro-Averaging

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

52

## Micro- vs. Macro-Averaging: Example

Class 1

	Truth: yes	Truth: no
Classifier: yes	10	10
Classifier: no	10	970

Class 2

	Truth: yes	Truth: no
Classifier: yes	90	10
Classifier: no	10	890

Micro Ave. Table

	Truth: yes	Truth: no
Classifier: yes	100	20
Classifier: no	20	1860

- Macroaveraged precision:
- Microaveraged precision:
- Microaveraged score is dominated by score on common classes

53

## Micro- vs. Macro-Averaging: Example

Class 1

	Truth: yes	Truth: no
Classifier: yes	10	10
Classifier: no	10	970

Class 2

	Truth: yes	Truth: no
Classifier: yes	90	10
Classifier: no	10	890

Micro Ave. Table

	Truth: yes	Truth: no
Classifier: yes	100	20
Classifier: no	20	1860

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

54

## 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

## Statistical Significance

- Suppose we have two classifiers,  $\text{classify}_1$  and  $\text{classify}_2$ .
- Is  $\text{classify}_1$  better? The "null hypothesis," denoted  $H_0$ , is that it isn't. But if  $\text{Accuracy}_1 \gg \text{Accuracy}_2$  (or whatever your evaluation metric is instead of accuracy) we are tempted to believe otherwise.
- How much larger must  $A_1$  be than  $A_2$  to reject  $H_0$ ?
- Frequentist view: how (im)probable is the observed difference, given  $H_0 = \text{true}$ ?

# Text Classification

## Practical Issues

### The Real World

- Gee, I'm building a text classifier for real, now!
- What should I do?

## No training data? Manually written rules

If (wheat or grain) and not (whole or bread) then  
Categorize as grain

- Need careful crafting
  - Human tuning on development data
  - Time-consuming: 2 days per class

59

## Very little data?

- Use Naive Bayes
  - On Discriminative vs. Generative classifiers: A comparison of logistic regression and naive Bayes (Ng and Jordan 2002 NIPS)
- Get more labeled data
  - Find clever ways to get humans to label data for you
- Try semi-supervised machine learning methods

60

## **A reasonable amount of data?**

- Try more clever classifiers

61

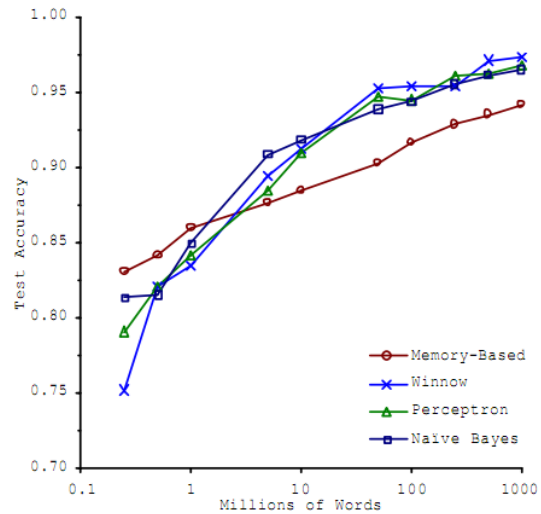
## **A huge amount of data?**

- Can achieve high accuracy!
- At a cost (high train or test time for some methods)
- So Naive Bayes can come back into its own again!

62

## Accuracy as a function of data size

- With enough data
  - Classifier may not matter



Brill and Banko on spelling correction

63

## Underflow Prevention: log space

- Multiplying lots of probabilities can result in floating-point underflow.
- Since  $\log(xy) = \log(x) + \log(y)$ 
  - Better to 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



## How to tweak performance

- Domain-specific features and weights: *very* important in real performance
- Sometimes need to collapse terms:
  - Part numbers, chemical formulas, ...
  - But stemming generally doesn't help
- Upweighting: Counting a word as if it occurred twice:
  - title words (Cohen & Singer 1996)
  - first sentence of each paragraph (Murata, 1999)
  - In sentences that contain title words (Ko *et al*, 2002)

65

## Sentiment Analysis

## 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.

67

## Google Product Search



HP Officejet 6500A Plus e-All-in-One Color Ink-jet - Fax / copier / printer / scanner  
\$89 online, \$100 nearby ★★★★★ 377 reviews  
September 2010 - Printer - HP - Inkjet - Office - Copier - Color - Scanner - Fax - 250 sh

### Reviews

Summary - Based on 377 reviews



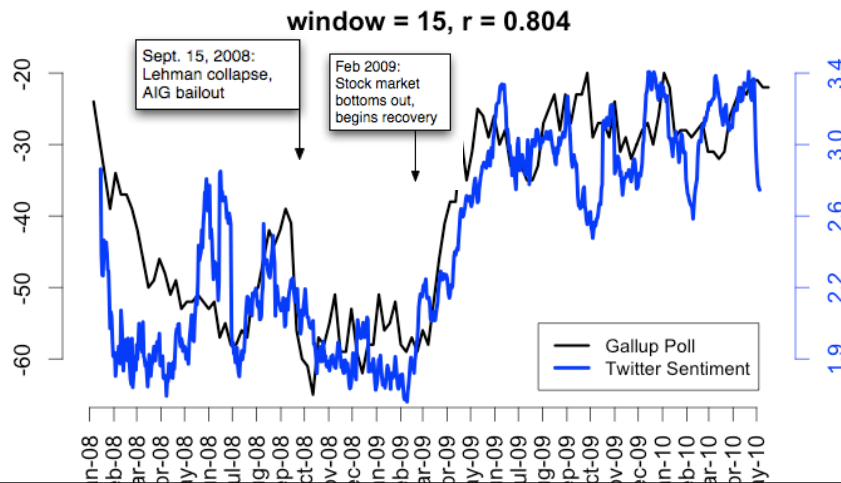
What people are saying

Category	Rating	Quote
ease of use	★★★★★	"This was very easy to setup to four computers."
value	★★★★★	"Appreciate good quality at a fair price."
setup	★★★★★	"Overall pretty easy setup."
customer service	★★★☆☆	"I DO like honest tech support people."
size	★★★★★	"Pretty Paper weight."
mode	★★★☆☆	"Photos were fair on the high quality mode."
colors	★★★★★	"Full color prints came out with great quality."

68

# Twitter sentiment versus Gallup Poll of Consumer Confidence

Brendan O'Connor, Ramnath Balasubramanian, Bryan R. Routledge, and Noah A. Smith. 2010. From Tweets to Polls: Linking Text Sentiment to Public Opinion Time Series. In ICWSM-2010



# Target Sentiment on Twitter

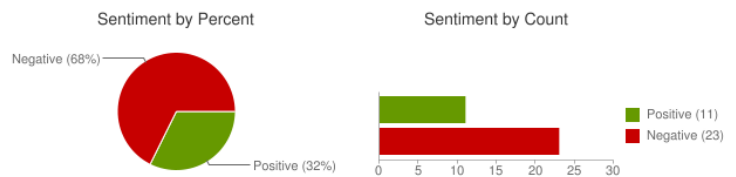
Type in a word and we'll highlight the good and the bad

- [Twitter Sentiment App](#)

[Save this search](#)

- Alec Go, Richa Bhayani, Lei Huang. 2009. Twitter Sentiment Classification using Distant Supervision

## Sentiment analysis for "united airlines"



**jjacobson:** OMG... Could @United airlines have worse customer service? W8g now 15 minutes  
Posted 2 hours ago

**12345clumsy6789:** I hate United Airlines Ceiling!!! Fukn impossible to get my conduit in this d  
Posted 2 hours ago

**EMLandPRGbelgiu:** EML/PRG fly with Q8 united airlines and 24seven to an exotic destination  
Posted 2 hours ago

**CountAdam:** FANTASTIC customer service from United Airlines at XNA today. Is tweet more  
Posted 4 hours ago

## Sentiment analysis has many other names

- Opinion extraction
- Opinion mining
- Sentiment mining
- Subjectivity analysis

71

## Why sentiment analysis?

- *Movie*: is this review positive or negative?
- *Products*: what do people think about the new iPhone?
- *Public sentiment*: how is consumer confidence? Is despair increasing?
- *Politics*: what do people think about this candidate or issue?
- *Prediction*: predict election outcomes or market trends from sentiment

72

## Sentiment Analysis

- Simplest task:
  - Is the attitude of this text positive or negative?
- More complex:
  - Rank the attitude of this text from 1 to 5
- Advanced:
  - Detect the target, source, or complex attitude types

## Sentiment Analysis

- Simplest task:
  - Is the attitude of this text positive or negative?
- More complex:
  - Rank the attitude of this text from 1 to 5
- Advanced:
  - Detect the target, source, or complex attitude types

# Sentiment Analysis

## Using Naive Bayes

### Sentiment Classification in Movie Reviews

Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? Sentiment Classification using Machine Learning Techniques. EMNLP-2002, 79—86.

Bo Pang and Lillian Lee. 2004. A Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts. ACL, 271-278

- Polarity detection:
  - Is an IMDB movie review positive or negative?
- Data: *Polarity Data 2.0*:
  - <http://www.cs.cornell.edu/people/pabo/movie-review-data>

## IMDB data in the Pang and Lee database



when `_star wars_` came out some twenty years ago , the image of traveling throughout the stars has become a commonplace image . [...]

when han solo goes light speed , the stars change to bright lines , going towards the viewer in lines that converge at an invisible point .

cool .

`_october sky_` offers a much simpler image—that of a single white dot , traveling horizontally across the night sky . [ . . . ]



“ snake eyes ” is the most aggravating kind of movie : the kind that shows so much potential then becomes unbelievably disappointing .

it’s not just because this is a brian depalma film , and since he’s a great director and one who’s films are always greeted with at least some fanfare .

and it’s not even because this was a film starring nicolas cage and since he gives a brauvara performance , this film is hardly worth his talents .

## Baseline Algorithm (adapted from Pang and Lee)

- Tokenization
- Feature Extraction
- Classification using different classifiers
  - Naive Bayes
  - ...

## Sentiment Tokenization Issues

- Deal with HTML and XML markup
- Twitter mark-up (names, hash tags)
- Capitalization (preserve for words in all caps)
- Phone numbers, dates
- Emoticons

79

## Extracting Features for Sentiment Classification

- How to handle negation
  - I **didn't** like this movie
  - vs
  - I really like this movie
- Which words to use?
  - Only adjectives
  - All words
    - All words turns out to work better, at least on this data

80



## Negation

Das, Sanjiv and Mike Chen. 2001. Yahoo! for Amazon: Extracting market sentiment from stock message boards. In Proceedings of the Asia Pacific Finance Association Annual Conference (APFA).  
Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? Sentiment Classification using Machine Learning Techniques. EMNLP-2002, 79–86.

Add NOT\_ to every word between negation and following punctuation:

didn't like this movie , but I



didn't NOT\_like NOT\_this NOT\_movie but I

## Reminder: Naïve Bayes

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

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

## Binarized (Boolean feature) Naive Bayes

- Intuition:
  - For sentiment (and probably for other text classification domains)...
  - Word occurrence may matter more than word frequency
    - The occurrence of the word *fantastic* tells us a lot
    - The fact that it occurs 5 times may not tell us much more.
  - Boolean Naive Bayes
    - Clips all the word counts in each document at 1

83

## Boolean Naive 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
  - Remove duplicates in each doc
  - For each word  $w_k$  in vocabulary
    - Retain only a single instance of  $w_k$
    - $n_k \leftarrow$  # of occurrences of  $w_k$  in  $Text_j$
$$P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha |Vocabulary|}$$

## Boolean Multinomial Naive Bayes on a test document $d$

- First remove all duplicate words from  $d$
- Then compute NB using the same equation:

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

85

## Normal vs. Boolean Multinomial NB

Normal	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	?

Boolean	Doc	Words	Class
Training	1	Chinese Beijing	c
	2	Chinese Shanghai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Tokyo Japan	?

## Problems:

### What makes reviews hard to classify?

- Subtlety:
  - Perfume review in *Perfumes: the Guide*:
    - “If you are reading this because it is your darling fragrance, please wear it at home exclusively, and tape the windows shut.”
  - Dorothy Parker on Katherine Hepburn
    - “She runs the gamut of emotions from A to B”

87

## Thwarted Expectations and Ordering Effects

- “This film should be **brilliant**. It sounds like a **great** plot, the actors are **first grade**, and the supporting cast is **good** as well, and Stallone is attempting to deliver a good performance. However, it **can’t hold up.**”
- Well as usual Keanu Reeves is nothing special, but surprisingly, the **very talented** Laurence Fishbourne is **not so good** either, I was surprised.

88