

# 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.

# 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



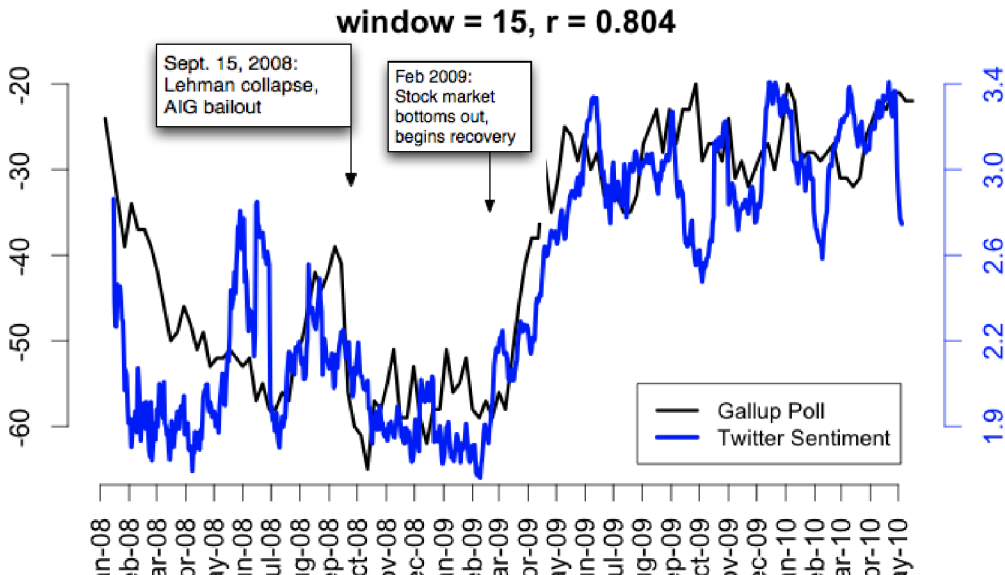
What people are saying

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

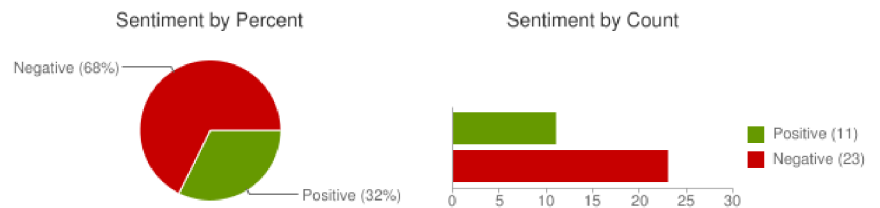
- [Twitter Sentiment App](#)

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

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


[Save this search](#)

### Sentiment analysis for "united airlines"



[ijacobson](#): OMG... Could @United airlines have worse customer service? W8g now 15 minut  
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

70

## 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
  - Classify for 6 basic emotions Happiness, Sadness, Fear, Disgust, Anger, Surprise
- 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 containing all  $docs_j$
  - For each word type  $w$  in doc
    - Retain only a single instance of  $w$
    - $n_k \leftarrow$  # of occurrences of  $w_k$  in  $Text_j$
  - For each word  $w_k$  in *Vocabulary*
    - $$P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha |Vocabulary|}$$

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

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

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

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	?

86

## 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