

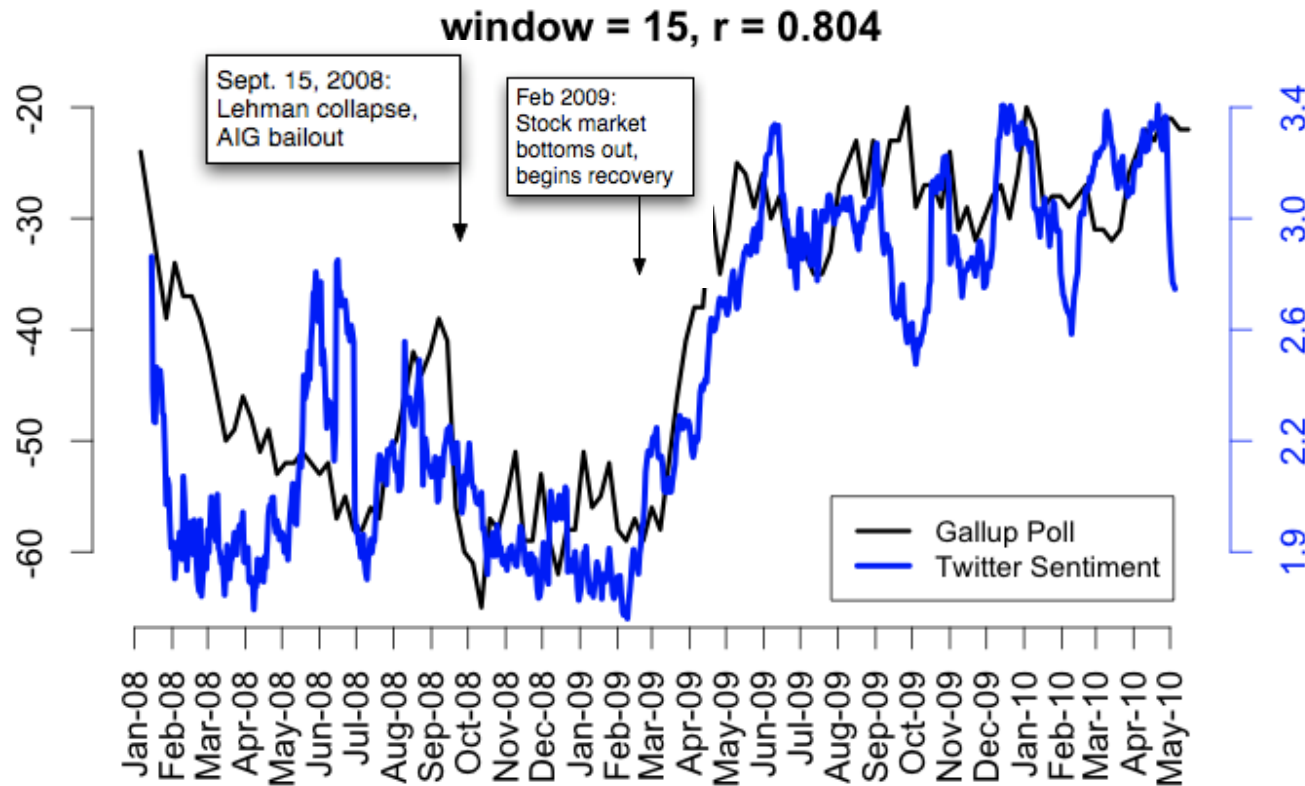
# Clasificación III

---

Rubén Francisco Manrique  
[rf.manrique@uniandes.edu.co](mailto:rf.manrique@uniandes.edu.co)

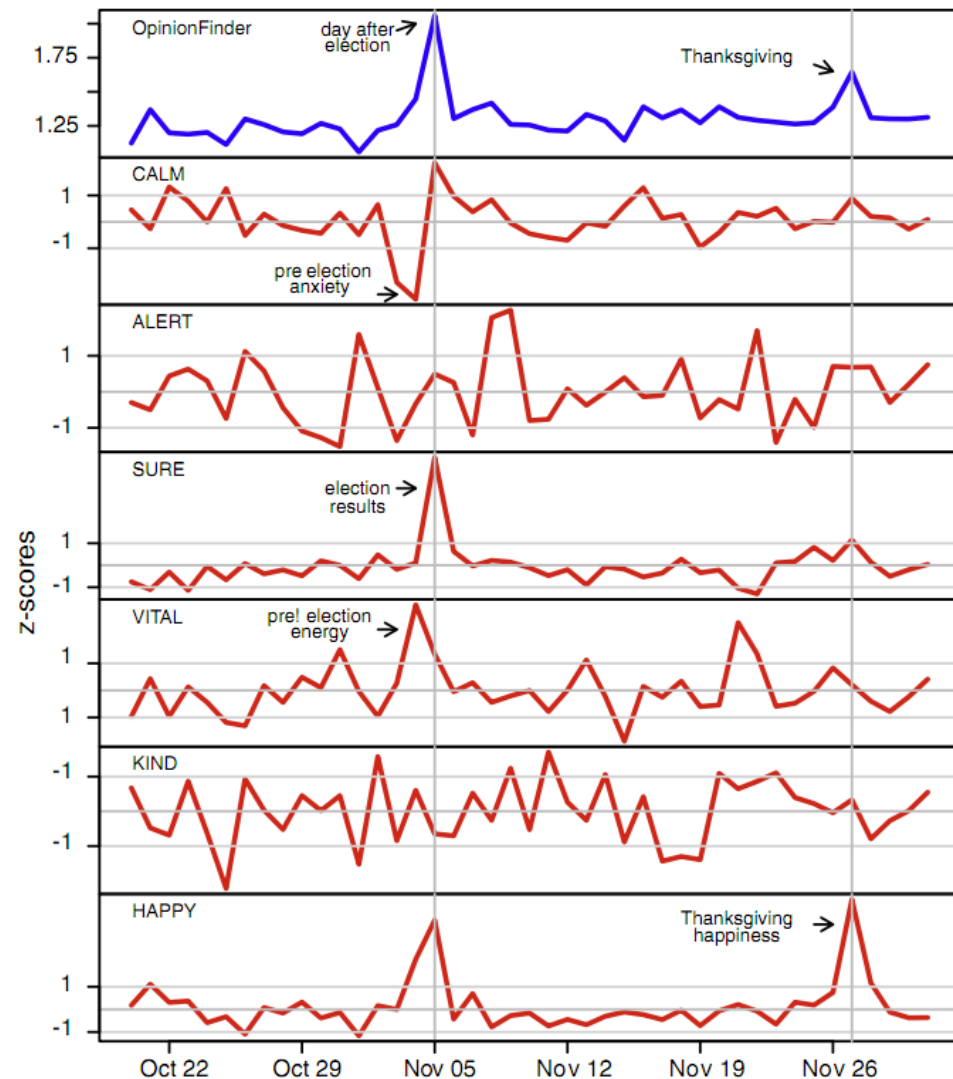
# Twitter sentiment versus Gallup Poll of Consumer Confidence

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



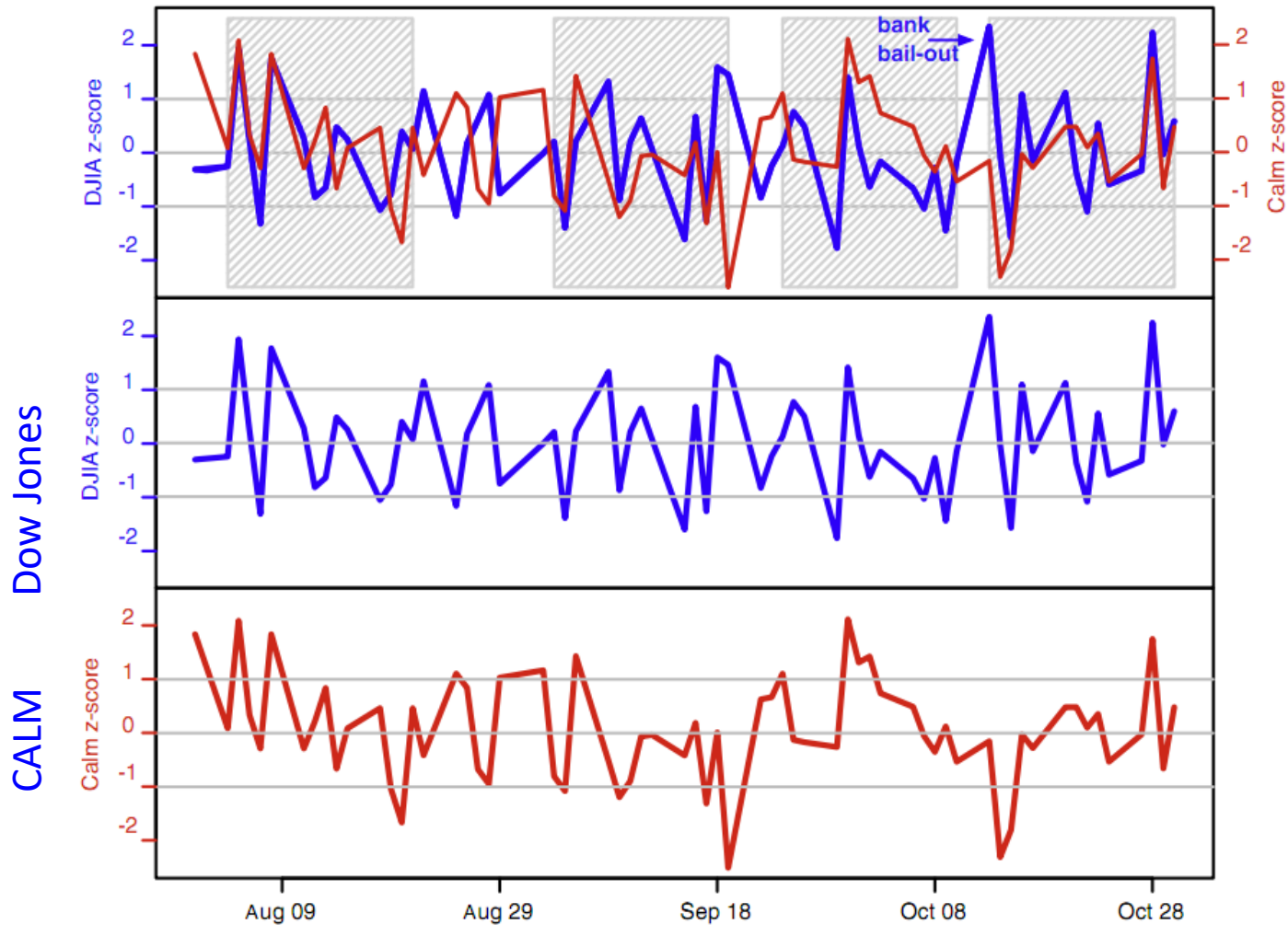
# Twitter sentiment:

Johan Bollen, Huina Mao, Xiaojun Zeng. 2011.  
[Twitter mood predicts the stock market](#),  
Journal of Computational Science 2:1, 1-8.  
10.1016/j.jocs.2010.12.007.



Bollen et al. (2011)

- CALM predicts DJIA 3 days later
- At least one current hedge fund uses this algorithm



# Twitter biz

- [Twitter Sentiment App](#)
- Alec Go, Richa Bhayani, Lei Huang. 2009. Twitter Sentiment Classification using Distant Supervision
- [https://www.csc2.ncsu.edu/faculty/healey/tweet\\_viz/tweet\\_app/](https://www.csc2.ncsu.edu/faculty/healey/tweet_viz/tweet_app/)

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

The figure displays sentiment analysis results for the query "united airlines". It consists of two charts: a pie chart titled "Sentiment by Percent" and a horizontal bar chart titled "Sentiment by Count".

**Sentiment by Percent:** A pie chart showing the distribution of sentiment. The red section represents Negative sentiment at 68%, and the green section represents Positive sentiment at 32%.

**Sentiment by Count:** A horizontal bar chart showing the count of tweets for each sentiment. The green bar represents Positive sentiment with a count of 11, and the red bar represents Negative sentiment with a count of 23. A legend on the right indicates: Positive (11) in green and Negative (23) in red.

Sentiment	Count	Percent
Positive	11	32%
Negative	23	68%

[jlljacobson](#): 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

# 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

# Sentiment Analysis

- Sentiment analysis is the detection of **attitudes**  
“enduring, affectively colored beliefs, dispositions towards objects or persons”
  1. **Holder (source)** of attitude
  2. **Target (aspect)** of attitude
  3. **Type** of attitude
    - From a set of types
      - *Like, love, hate, value, desire, etc.*
    - Or (more commonly) simple weighted **polarity**:
      - *positive, negative, neutral, together with strength*
  4. **Text** containing the attitude
    - Sentence or entire document

# 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 Classification in Movie Reviews

- <https://ai.stanford.edu/~amaas/data/sentiment/>



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

- Tokenization (emoticons?)
- Feature Extraction
- Classification using different classifiers
  - Naïve Bayes

# 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

# 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

# Binarized (Boolean feature) Multinomial Naïve 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 Multinomial Naïve Bayes
    - Clips all the word counts in each document at 1
- First remove all duplicate words from  $d$
- Then compute NB using the same equation:

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

# 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	?

# Binarized (Boolean feature) Multinomial Naïve Bayes

B. Pang, L. Lee, and S. Vaithyanathan. 2002. Thumbs up? Sentiment Classification using Machine Learning Techniques. EMNLP-2002, 79—86.

V. Metsis, I. Androutsopoulos, G. Paliouras. 2006. Spam Filtering with Naive Bayes – Which Naive Bayes? CEAS 2006 - Third Conference on Email and Anti-Spam.

K.-M. Schneider. 2004. On word frequency information and negative evidence in Naive Bayes text classification. ICANLP, 474-485.

JD Rennie, L Shih, J Teevan. 2003. Tackling the poor assumptions of naive bayes text classifiers. ICML 2003

- Binary seems to work better than full word counts
- Other possibility:  $\log(\text{freq}(w))$

# Lexicons

---



# Bing Liu Opinion Lexicon

Minqing Hu and Bing Liu. Mining and Summarizing Customer Reviews. ACM SIGKDD-2004.

- [Bing Liu's Page on Opinion Mining](#)
- <http://www.cs.uic.edu/~liub/FBS/opinion-lexicon-English.rar>
- 6786 words
  - 2006 positive
  - 4783 negative

# SentiWordNet

Stefano Baccianella, Andrea Esuli, and Fabrizio Sebastiani. 2010 SENTIWORDNET 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining. LREC-2010

- All WordNet synsets automatically annotated for degrees of positivity, negativity, and neutrality/objectiveness
- [estimable(J,3)] “may be computed or estimated”  
Pos 0 Neg 0 Obj 1
- [estimable(J,1)] “deserving of respect or high regard”  
Pos .75 Neg 0 Obj .25

# Analyzing the polarity of each word in IMDB

Potts, Christopher. 2011. On the negativity of negation. SALT 20, 636-659.

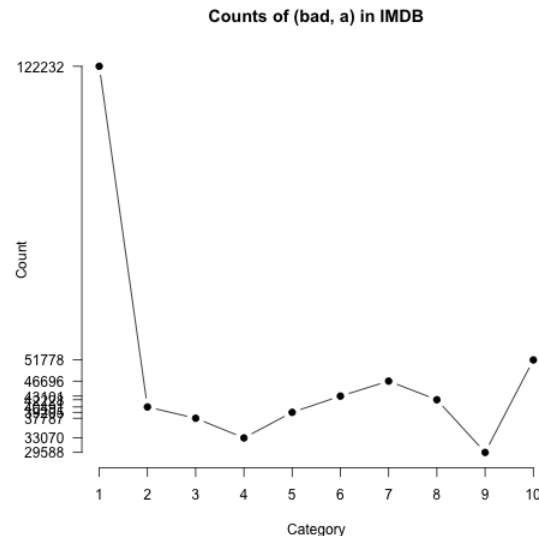
- How likely is each word to appear in each sentiment class?
- Count(“bad”) in 1-star, 2-star, 3-star, etc.
- But can’t use raw counts:

- Instead, **likelihood**:  $P(w|c) = \frac{f(w,c)}{\sum_{w \in \hat{w}} f(w,c)}$

- Make them comparable between words

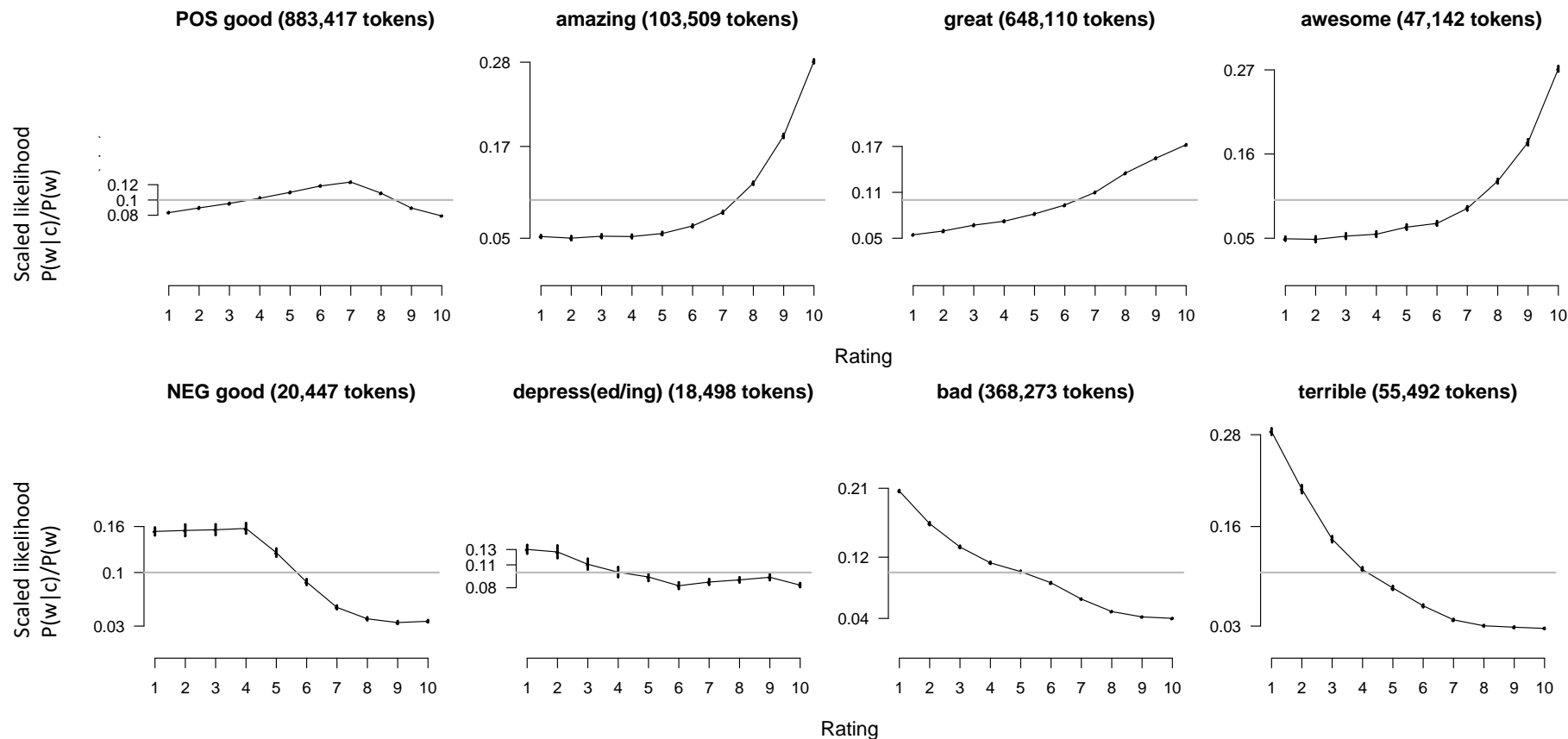
- **Scaled likelihood**:

$$\frac{P(w|c)}{P(w)}$$



# Analyzing the polarity of each word in IMDB

Potts, Christopher. 2011. On the negativity of negation. SALT 20, 636-659.



# Other sentiment feature: Logical negation

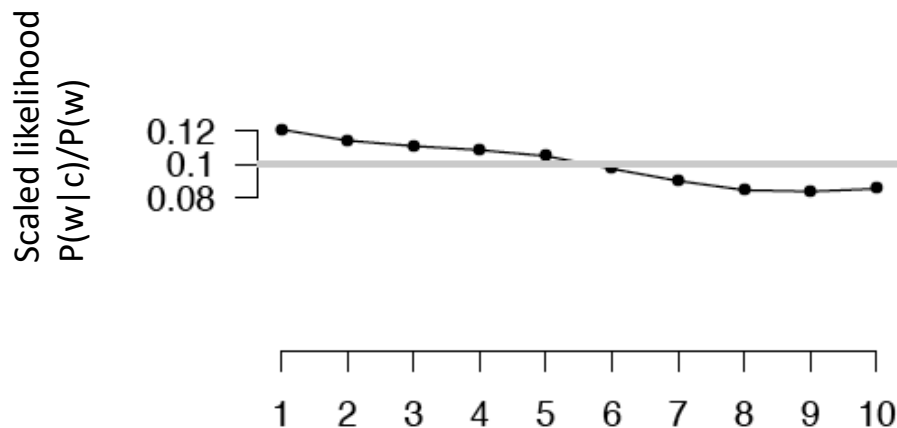
Potts, Christopher. 2011. On the negativity of negation. *SALT* 20, 636-659.

- Is logical negation (*no*, *not*) associated with negative sentiment?
- Potts experiment:
  - Count negation (*not*, *n't*, *no*, *never*) in online reviews
  - Regress against the review rating

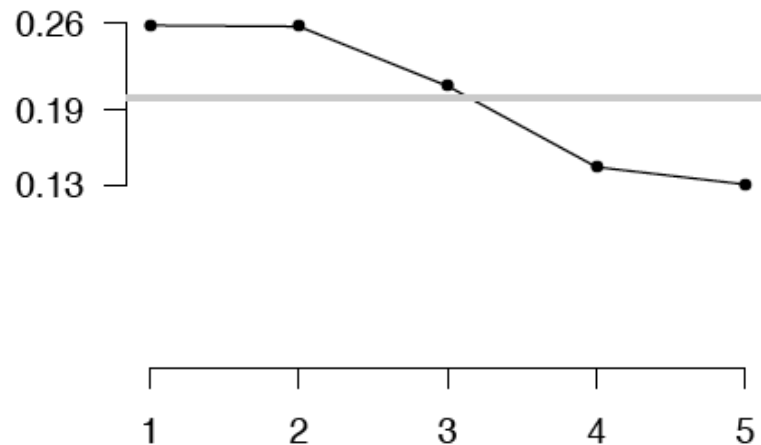
# Potts 2011 Results:

## More negation in negative sentiment

**IMDB (4,073,228 tokens)**



**Five-star reviews (846,444 tokens)**



Gracias por la atención

¿Tiene alguna pregunta?

