



Sentiment Analysis

What is Sentiment Analysis?



Normal vs. Boolean Multinomial NB

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

Handwritten annotations for Normal NB:

- Red box highlights "Chinese" in row 1.
- Red bracket groups "Chinese", "Beijing", and "Chinese" in row 1.
- Red bracket groups "Chinese", "Chinese", and "Shanghai" in row 2.
- Red bracket groups "Chinese", "Macao" in row 3.
- Red bracket groups "Tokyo", "Japan", and "Chinese" in row 4.
- Red bracket groups "Chinese", "Chinese", and "Chinese" in row 5.
- Red bracket groups "Tokyo", "Japan", and "Japanese" in row 5.
- Red text "T_{Chinese} = 5" is written below row 5.
- Red text "T_{Chinese} = 3" is written below row 5.
- Red text "T_{Tokyo} = 1" is written below row 5.
- Red text "T_{Japan} = 1" is written below row 5.

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

Handwritten annotations for Boolean NB:

- Red box highlights "Chinese" and "Beijing" in row 1.
- Red box highlights "Chinese" and "Shanghai" in row 2.
- Red box highlights "Chinese" and "Macao" in row 3.
- Red box highlights "Tokyo", "Japan", and "Chinese" in row 4.
- Red box highlights "Chinese", "Tokyo", and "Japan" in row 5.
- Red bracket groups "Chinese", "Beijing", and "Chinese" in row 1.
- Red bracket groups "Chinese", "Shanghai", and "Chinese" in row 2.
- Red bracket groups "Chinese", "Macao", and "Chinese" in row 3.
- Red bracket groups "Tokyo", "Japan", and "Chinese" in row 4.
- Red bracket groups "Chinese", "Tokyo", and "Japan" in row 5.
- Red text "T_{Chinese} = 3" is written below row 5.
- Red text "T_{Tokyo} = 1" is written below row 5.
- Red text "T_{Japan} = 1" is written below row 5.
- Red text "Beijing - not present" is written below row 5.
- Red text "P(B₀ | y) = 1/6" is written below row 5.

TRAINMULTINOMIALNB(C, D)

```
1  $V \leftarrow \text{EXTRACTVOCABULARY}(D)$ 
2  $N \leftarrow \text{COUNTDOCS}(D)$ 
3 for each  $c \in C$ 
4 do  $N_c \leftarrow \text{COUNTDOCSINCLASS}(D, c)$ 
5    $prior[c] \leftarrow N_c/N$   $\leftarrow$ 
6    $text_c \leftarrow \text{CONCATENATETEXTOFAALLDOCSINCLASS}(D, c)$ 
7   for each  $t \in V$ 
8   do  $T_{ct} \leftarrow \text{COUNTTOKENSOFTERM}(text_c, t)$ 
9   for each  $t \in V$   $\xrightarrow{\text{if } T_{ct} > 0}$ 
10   do  $condprob[t][c] \leftarrow \frac{T_{ct}+1}{\sum_{t'}(T_{ct'}+1)} = \frac{T_{ct}}{N+v}$   $T_{ct} = 1$ 
11 return  $V, prior, condprob$ 
```

APPLYMULTINOMIALNB($C, V, prior, condprob, d$)

```
1  $W \leftarrow \text{EXTRACTTOKENSFROMDOC}(V, d)$ 
2 for each  $c \in C$ 
3 do  $score[c] \leftarrow \log prior[c]$ 
4   for each  $t \in W$ 
5   do  $score[c] += \log condprob[t][c]$ 
6 return  $\arg \max_{c \in C} score[c]$ 
```

$|V|$

N

$$P(C) = \frac{N_c}{N}$$

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

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

$$score[c] = \log prior[c] + \left\{ \log condprob[t][c] \right\}$$

	w_0	w_1	w_2	w_3	w_4
Doc 1	41	71	0	21	0
Doc 2					
Doc 3					
.					
Doc p					

1	1	0	1	0
0				
1				
1				
0				

$$P(w|c) = \frac{\text{count of word in all the documents of class}}{\text{total no. of words}}$$

if word is present then prob is p
otherwise it is $(1-p)$

TRAINBERNOULLINB(C, D)

```
1  $V \leftarrow \text{EXTRACTVOCABULARY}(D)$ 
2  $N \leftarrow \text{COUNTDOCS}(D)$ 
3 for each  $c \in C$ 
4 do  $N_c \leftarrow \text{COUNTDOCSINCLASS}(D, c)$ 
5  $prior[c] \leftarrow N_c/N$ 
6 for each  $t \in V$ 
7 do  $N_{ct} \leftarrow \text{COUNTDOCSINCLASSCONTAININGTERM}(D, c, t)$ 
8  $condprob[t][c] \leftarrow (N_{ct} + 1)/(N_c + 2)$ 
9 return  $V, prior, condprob$ 
```

Bernoulli(p)

if $X = 1$ $\Pr(X=1) = p$
else $\Pr(X=0) = 1-p$

APPLYBERNOULLINB($C, V, prior, condprob, d$)

```
1  $V_d \leftarrow \text{EXTRACTTERMSFROMDOC}(V, d)$ 
2 for each  $c \in C$ 
3 do  $score[c] \leftarrow \log prior[c]$ 
4 for each  $t \in V$ 
5 do if  $t \in V_d$ 
6 then  $score[c] += \log condprob[t][c]$ 
7 else  $score[c] += \log(1 - condprob[t][c])$ 
8 return  $\arg \max_{c \in C} score[c]$ 
```

generates the probability
of words that are not
present in a document.

Bernoulli

- ① $P(t|c)$ = fraction of documents of class c that contain term t
- ② Binary occurrence model
- ③ Not performs well for long documents
- ④ Prob. of term not present in the test document is non-zero and $= 1 - P(t|c)$

Multinomial

- ① $P(t|c)$ = fraction of tokens in class c that contains term t
- ② Multiple occurrence model
- ③ Perform well for long documents.
- ④ Prob. of term not present in the test document is not taken into consideration.

$$\hat{P}(\text{Chinese}|c) = (3+1)/(3+2) = 4/5$$

$$\hat{P}(\text{Japan}|c) = \hat{P}(\text{Tokyo}|c) = (0+1)/(3+2) = 1/5$$

$$\hat{P}(\text{Beijing}|c) = \hat{P}(\text{Macao}|c) = \hat{P}(\text{Shanghai}|c) = (1+1)/(3+2) = 2/5$$

$$\hat{P}(\text{Chinese}|\bar{c}) = (1+1)/(1+2) = 2/3$$

$$\hat{P}(\text{Japan}|\bar{c}) = \hat{P}(\text{Tokyo}|\bar{c}) = (1+1)/(1+2) = 2/3$$

$$\hat{P}(\text{Beijing}|\bar{c}) = \hat{P}(\text{Macao}|\bar{c}) = \hat{P}(\text{Shanghai}|\bar{c}) = (0+1)/(1+2) = 1/3$$

$$\begin{aligned}\hat{P}(c|d_5) &\propto \hat{P}(c) \cdot \hat{P}(\text{Chinese}|c) \cdot \hat{P}(\text{Japan}|c) \cdot \hat{P}(\text{Tokyo}|c) \\ &\quad \cdot (1 - \hat{P}(\text{Beijing}|\bar{c})) \cdot (1 - \hat{P}(\text{Shanghai}|\bar{c})) \cdot (1 - \hat{P}(\text{Macao}|\bar{c})) \\ &= 3/4 \cdot 4/5 \cdot 1/5 \cdot 1/5 \cdot (1 - 2/5) \cdot (1 - 2/5) \cdot (1 - 2/5) \\ &\approx 0.005\end{aligned}$$

$$\begin{aligned}\hat{P}(\bar{c}|d_5) &\propto 1/4 \cdot 2/3 \cdot 2/3 \cdot 2/3 \cdot (1 - 1/3) \cdot (1 - 1/3) \cdot (1 - 1/3) \\ &\approx 0.022\end{aligned}$$

prob of
those words
which are not
present in class c



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
 - This is **not** the same as Multivariate Bernoulli Naïve Bayes
 - MBNB doesn't work well for sentiment or other text tasks
 - Other possibility: $\log(\text{freq}(w))$



The General Inquirer

Philip J. Stone, Dexter C Dunphy, Marshall S. Smith, Daniel M. Ogilvie. 1966. The General Inquirer: A Computer Approach to Content Analysis. MIT Press

- Home page: <http://www.wjh.harvard.edu/~inquirer>
- List of Categories: <http://www.wjh.harvard.edu/~inquirer/homecat.htm>
- Spreadsheet: <http://www.wjh.harvard.edu/~inquirer/inquirerbasic.xls>
- Categories:
 - Positiv (1915 words) and Negativ (2291 words)
 - Strong vs Weak, Active vs Passive, Overstated versus Understated
 - Pleasure, Pain, Virtue, Vice, Motivation, Cognitive Orientation, etc
- Free for Research Use



LIWC (Linguistic Inquiry and Word Count)

Pennebaker, J.W., Booth, R.J., & Francis, M.E. (2007). *Linguistic Inquiry and Word Count: LIWC 2007*. Austin, TX

- Home page: <http://www.liwc.net/>
- 2300 words, >70 classes
- **Affective Processes**
 - negative emotion (*bad, weird, hate, problem, tough*)
 - positive emotion (*love, nice, sweet*)
- **Cognitive Processes**
 - Tentative (*maybe, perhaps, guess*), Inhibition (*block, constraint*)
- **Pronouns, Negation** (*no, never*), **Quantifiers** (*few, many*)
- \$30 or \$90 fee



MPQA Subjectivity Cues Lexicon

Theresa Wilson, Janyce Wiebe, and Paul Hoffmann (2005). Recognizing Contextual Polarity in Phrase-Level Sentiment Analysis. Proc. of HLT-EMNLP-2005.

Riloff and Wiebe (2003). Learning extraction patterns for subjective expressions. EMNLP-2003.

- Home page: http://www.cs.pitt.edu/mpqa/subj_lexicon.html
- 6885 words from 8221 lemmas
 - 2718 positive
 - 4912 negative
- Each word annotated for intensity (strong, weak)
- GNU GPL



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

- Home page: <http://sentiwordnet.isti.cnr.it/>
- 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



Disagreements between polarity lexicons

Christopher Potts, [Sentiment Tutorial](#), 2011

	Opinion Lexicon	General Inquirer	SentiWordNet	LIWC
MPQA	33/5402 (0.6%)	49/2867 (2%)	1127/4214 (27%)	12/363 (3%)
Opinion Lexicon		32/2411 (1%)	1004/3994 (25%)	9/403 (2%)
General Inquirer			520/2306 (23%)	1/204 (0.5%)
SentiWordNet				174/694 (25%)
LIWC				



Analyzing the polarity of each word in IMDB

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

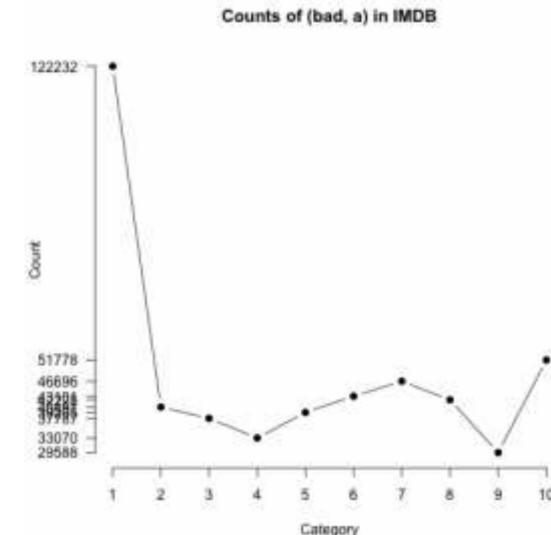
- How likely is each word to appear in each sentiment class?
- Count("bad") in 1-star, 2-star, 3-star, etc.
- But can't use raw counts:

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

poor — bad

- Make them comparable between words

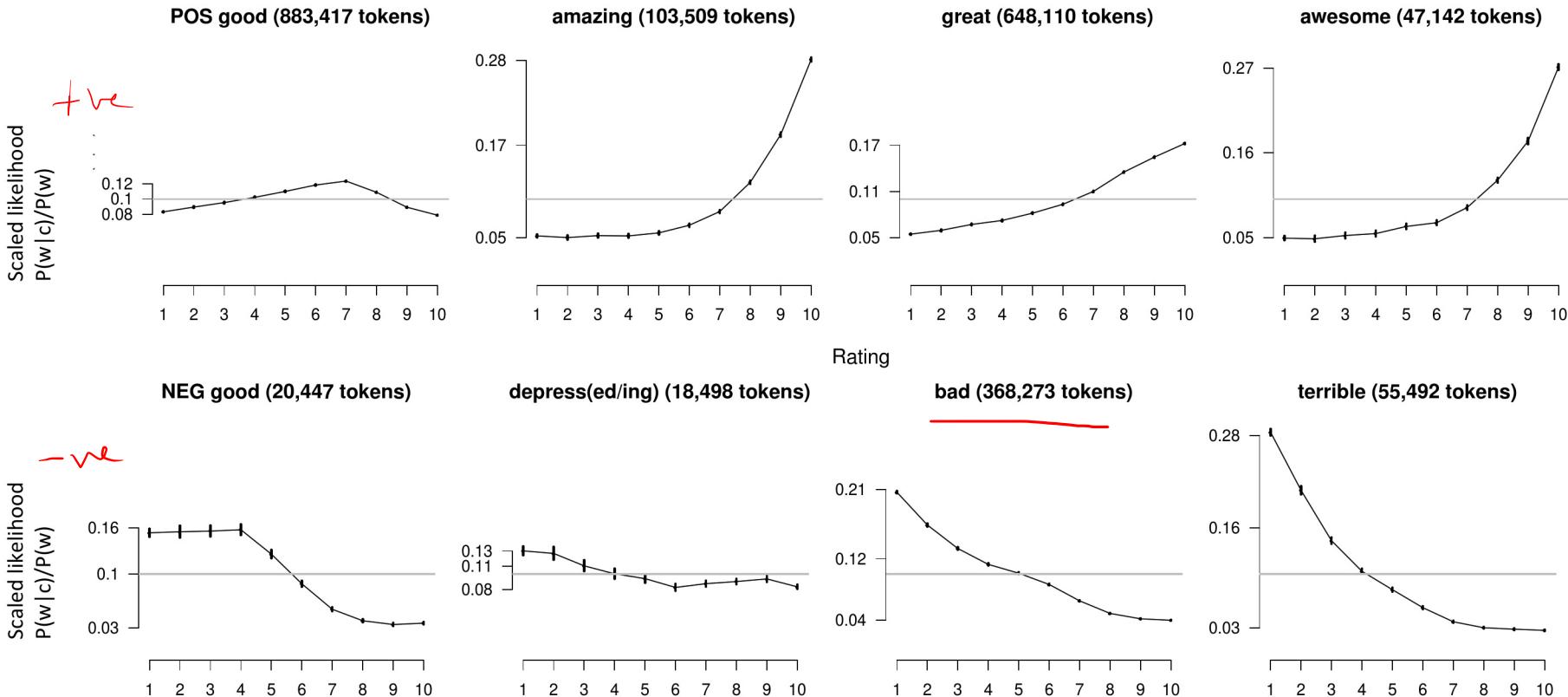
- **Scaled likelihood**:
$$\frac{P(w|c)}{P(w)}$$





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

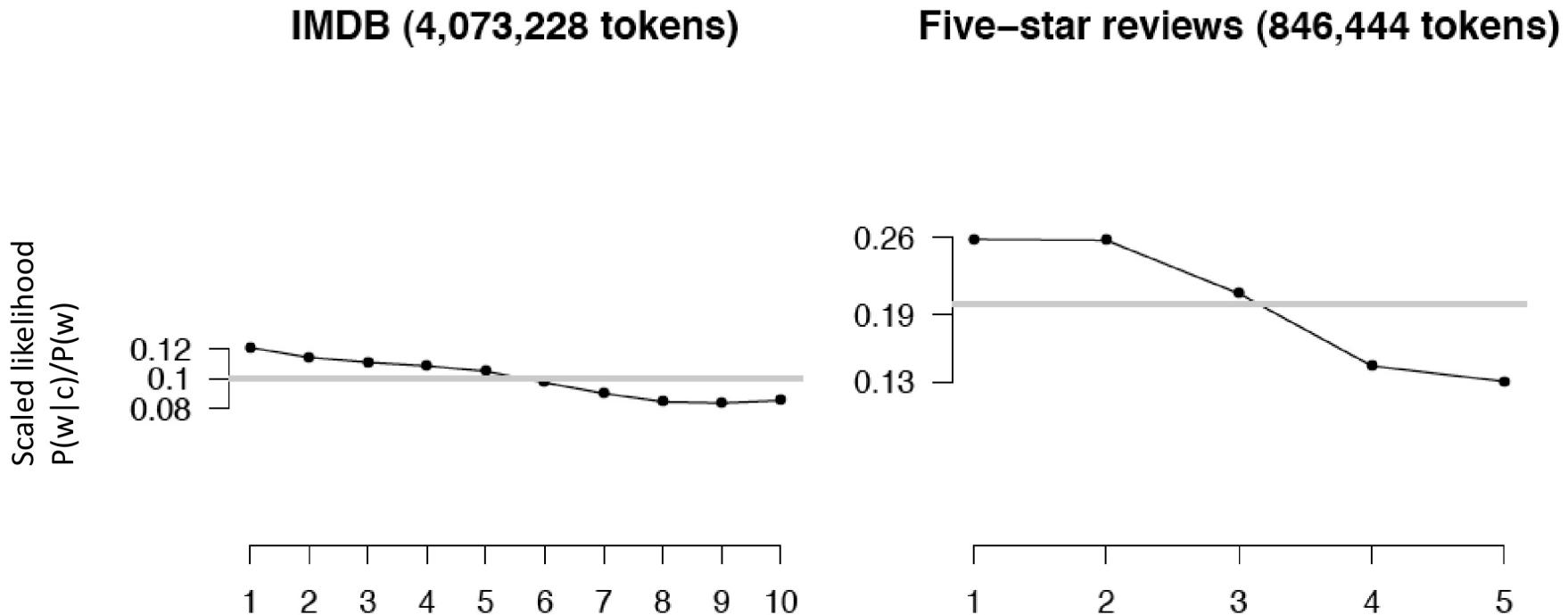
Potts, Christopher. 2011. On the negativity of negation. *SALT* 20, 636-659.

- 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





Sentiment Analysis

Sentiment Lexicons