# Sentiment Analysis

## A Baseline Algorithm

Dan Jurafsky

# 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 .

Dan Jurafsky

# Baseline Algorithm (adapted from Pang and Lee)

- Tokenization

- Feature Extraction

- Classification using different classifiers
  - Naïve Bayes
  - MaxEnt
  - SVM

# 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

- Useful code:
    - Christopher Potts sentiment tokenizer
    - Brendan O'Connor twitter tokenizer

21

Potts emoticons

```
[<>]?                        # optional hat/brow
[:;=8]                       # eyes
[\-o\*\']?                   # optional nose
[\)\]\(\[dDpP/\:\}\{@\|\\]   # mouth
|                            #### reverse orientation
[\)\]\(\[dDpP/\:\}\{@\|\\]   # mouth
[\-o\*\']?                   # optional nose
[:;=8]                       # eyes
[<>]?                        # optional hat/brow
```

Dan Jurafsky

# 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

22

Dan Jurafsky

# 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

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

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

# 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

25

# Boolean 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 $P(w_k \mid c_j)$ terms
  - Remove duplicates in each doc:
  - $Text_j \leftarrow$ single doc containing all $docs_j$
    - For each word type $w$ in doc
  - For each word $w_k$ in *Vocabulary*
    - Retain only a single instance of $w$
    $n_k \leftarrow$ \# of occurrences of $w_k$ in $Text_j$

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

# Boolean Multinomial Naïve Bayes on a test document *d*

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

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

27

# 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 | ? |

Dan Jurafsky

# 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
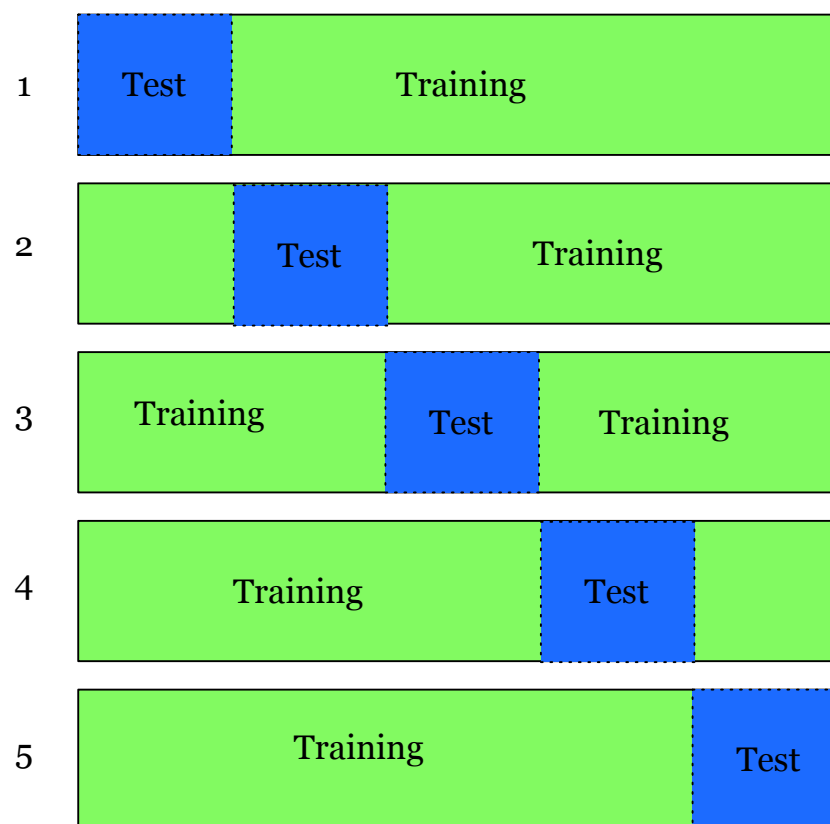
- Binary seems to work better than full word counts
  - This is **not** the same as Multivariate Bernoulli Naïve Bayes
    - MBNB doesn't work well for sentiment or other text tasks
- Other possibility: log(freq($w$))

29

Dan Jurafsky

# Cross-Validation

Iteration

- **Break up data into 10 folds**
  - (Equal positive and negative inside each fold?)
- **For each fold**
  - Choose the fold as a temporary test set
  - Train on 9 folds, compute performance on the test fold
- **Report average performance of the 10 runs**

| | | |
|---|---|---|
| 1 | Test | Training |
| 2 | Training | Test | Training |
| 3 | Training | Test | Training |
| 4 | Training | Test |
| 5 | Training | Test |

Dan Jurafsky

# **Other issues in Classification**

- MaxEnt and SVM tend to do better than Naïve Bayes

31

# 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"

Dan Jurafsky

# 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.

33

# Sentiment Analysis

## A Baseline Algorithm