



Sentiment Analysis

What is 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."



Bing Shopping

HP Officejet 6500A E710N Multifunction Printer

[Product summary](#) [Find best price](#) **Customer reviews** [Specifications](#) [Related items](#)



\$121.53 - \$242.39 (14 stores)

Compare

Average rating (144)

(55)

(54)

(10)

(6)

(23)

(0)

Most mentioned

Performance (57)

Ease of Use (43)

Print Speed (39)

Connectivity (31)

More ▾

Show reviews by source

[Best Buy \(140\)](#)

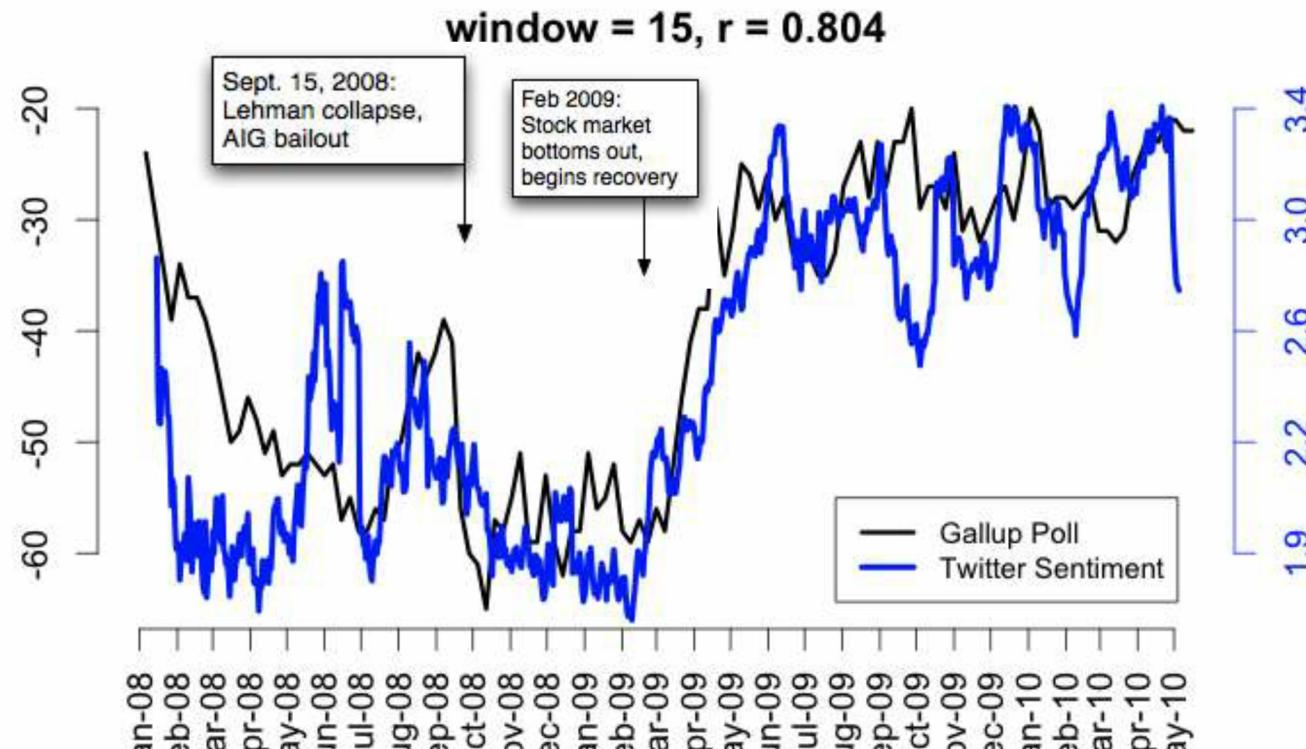
[CNET \(5\)](#)

[Amazon.com \(3\)](#)



Twitter sentiment versus Gallup Poll of Consumer Confidence

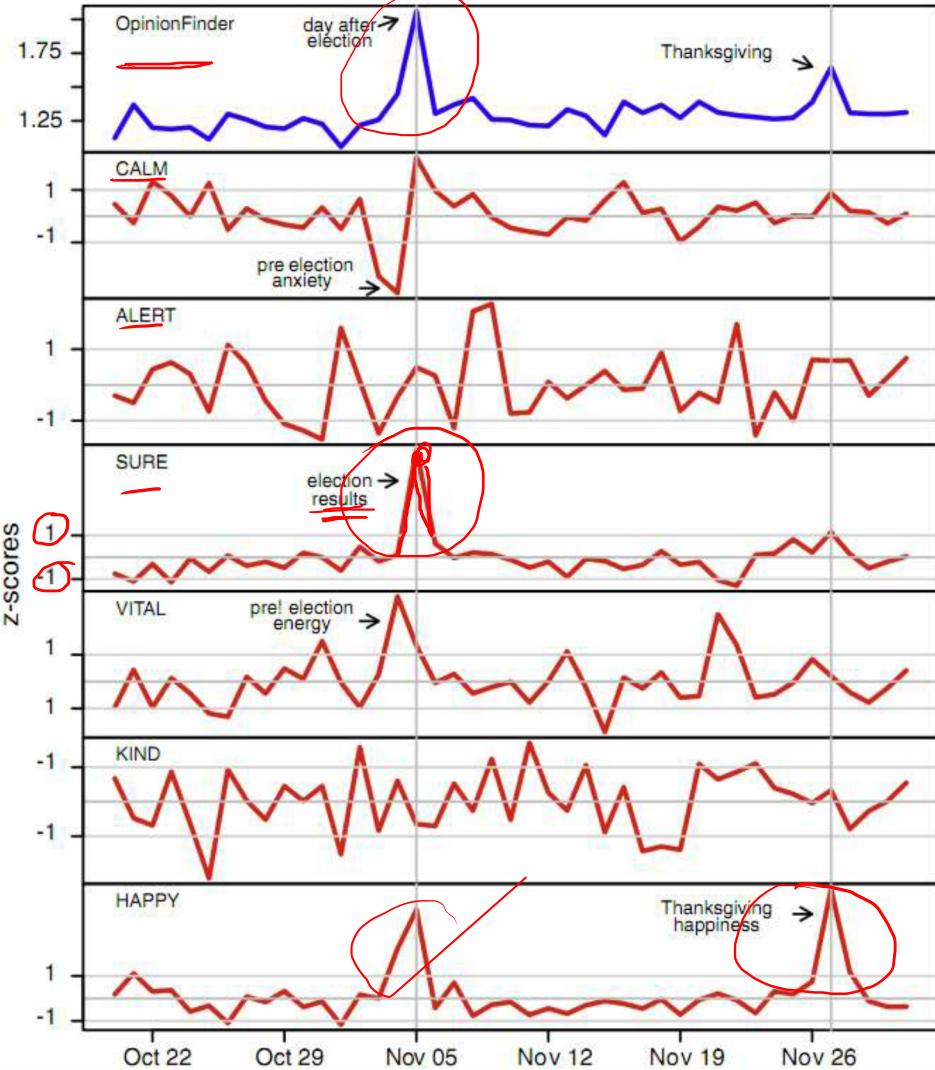
Brendan O'Connor, Ramnath Balasubramanyan, Bryan R. Routledge, and Noah A. Smith. 2010.
From Tweets to Polls: Linking Text Sentiment to Public Opinion Time Series. In ICWSM-2010





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.





Target Sentiment on Twitter

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

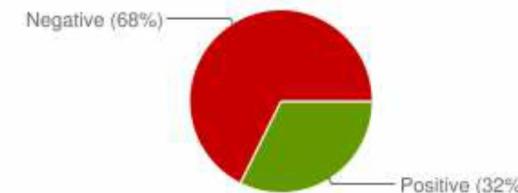
- Twitter Sentiment App

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

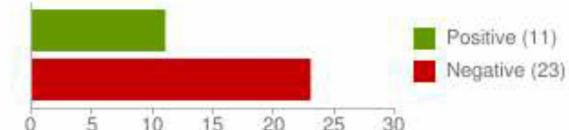
"united airlines" [Save this search](#)

Sentiment analysis for "united airlines"

Sentiment by Percent



Sentiment by Count



✓ lljacobson: 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



Sentiment analysis has many other names

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



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



Scherer Typology of Affective States

- **Emotion:** brief organically synchronized ... evaluation of a major event
 - *angry, sad, joyful, fearful, ashamed, proud, elated*
- **Mood:** diffuse non-caused low-intensity long-duration change in subjective feeling
 - *cheerful, gloomy, irritable, listless, depressed, buoyant*
- **Interpersonal stances:** affective stance toward another person in a specific interaction
 - *friendly, flirtatious, distant, cold, warm, supportive, contemptuous*
- **Attitudes:** enduring, affectively colored beliefs, dispositions towards objects or persons
 - *liking, loving, hating, valuing, desiring*
- **Personality traits:** stable personality dispositions and typical behavior tendencies
 - *nervous, anxious, reckless, morose, hostile, jealous*



Scherer Typology of Affective States

- **Emotion:** brief organically synchronized ... evaluation of a major event
 - *angry, sad, joyful, fearful, ashamed, proud, elated*
- **Mood:** diffuse non-caused low-intensity long-duration change in subjective feeling
 - *cheerful, gloomy, irritable, listless, depressed, buoyant*
- **Interpersonal stances:** affective stance toward another person in a specific interaction
 - *friendly, flirtatious, distant, cold, warm, supportive, contemptuous*
- **Attitudes: enduring, affectively colored beliefs, dispositions towards objects or persons**
 - *liking, loving, hating, valuing, desiring*
- **Personality traits:** stable personality dispositions and typical behavior tendencies
 - *nervous, anxious, reckless, morose, hostile, jealous*



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

more +ve
+ ve
neutral
- ve
more -ve

fine grained
labels to
sentiments



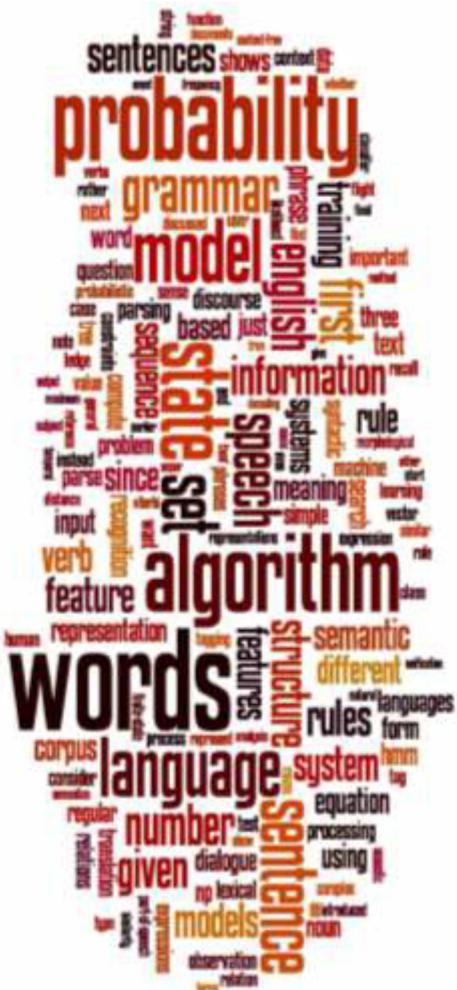
Sentiment Analysis

- 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



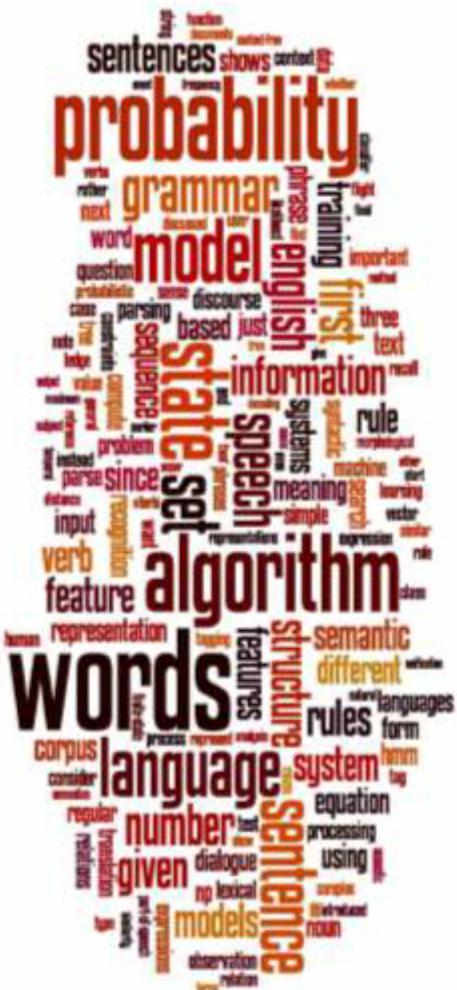
Sentiment Analysis

- 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

What is Sentiment Analysis?



Sentiment Analysis

A Baseline Algorithm



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 ✓
 - 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](#)

@ABC → NAME
 @ →
 # → TAG

Potts emoticons

```

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



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



Reminder: Naïve Bayes

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \cdot P(w_i | c_j)$$

i positions

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



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



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
 - Calculate $P(w_k | c_j)$ terms
 - Remove single doc containing all $docs_j$
 - For each word w_k in *Vocabulary*
 $n_k \leftarrow$ # of occurrences of w_k in $Text_j$
- $$P(c_j) \leftarrow \frac{|docs_j|}{|\text{total \# documents}|}$$
- $$P(w_k | c_j) \leftarrow \frac{n_k + 1}{n + |\text{Vocabulary}|}$$



Boolean Multinomial Naïve Bayes on a test document d

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

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



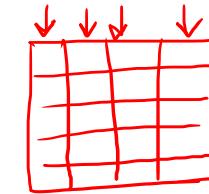
Normal vs. Boolean Multinomial NB

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	?

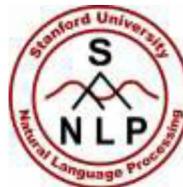


Binarized (Boolean feature) Multinomial Naïve Bayes



- 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

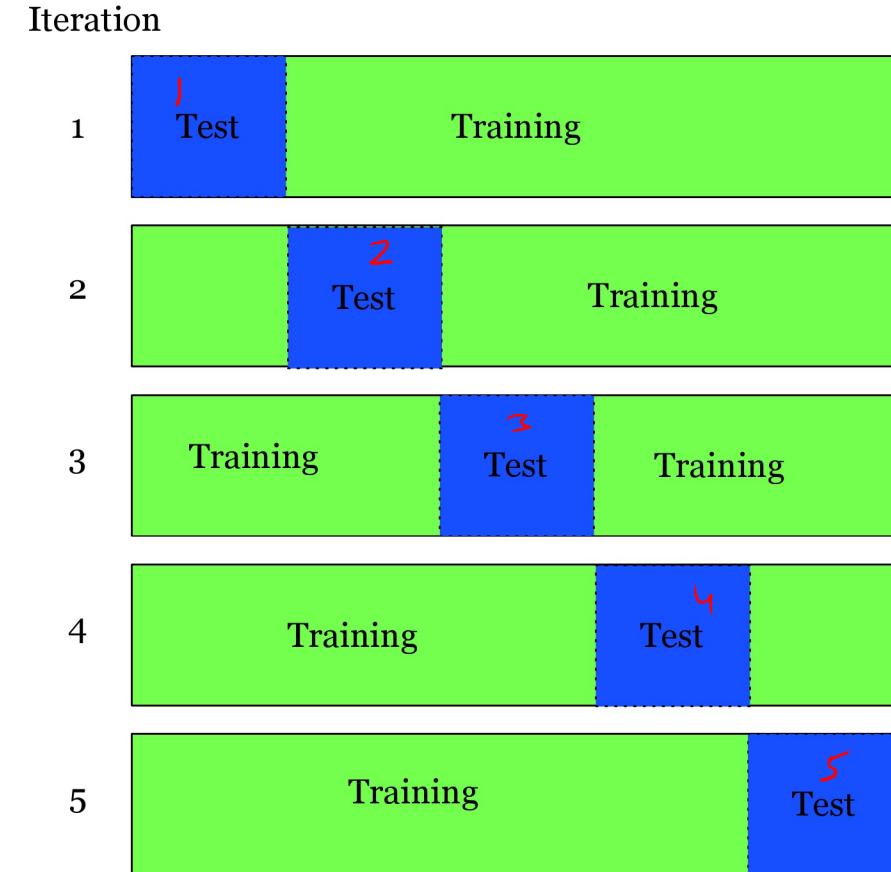
- Binary seems to work better than full word counts
 - This is **not** the same as Multivariate Bernoulli Naïve Bayes
 - MBNB doesn't work well for sentiment or other text tasks
 - Other possibility: $\log(\text{freq}(w))$



Cross-Validation

n folds

- 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





Other issues in Classification

- MaxEnt and SVM tend to do better than Naïve Bayes