

# NATURAL LANGUAGE PROCESSING

---

Gayane Petrosyan  
McGill University

Luys DAP 2012

# Text classification

- Spam classification
- Authorship classification
- Age/Gender classification
- Sentiment analysis
- Subject identification
- Language Identification

# Text Classification

## Supervised learning

- Input
  - Document  $\mathbf{d}$
  - A fixed set of classes  $\mathcal{C} = \{c_1, c_2, \dots, c_n\}$
  - A training set  $\mathbf{m}$  hand-labeled documents  $(\mathbf{d}_1, c_1), \dots, (\mathbf{d}_m, c_m)$
- Output
  - A learned classifier  $\mathbf{y}: \mathbf{d} \rightarrow \mathbf{c}$

# Naïve Bayes

- Document is bag of words(only count of words is considered)
- Bayes, because Bayes rule is used
  - $P(c/d) = \frac{P(d/c)P(c)}{P(d)}$
  - $C = \operatorname{argmax} P(c/d) = \operatorname{argmax} \frac{P(d/c)P(c)}{P(d)} = \operatorname{argmax} P(d/c)P(c)$
- Naïve, because
  - $C = \operatorname{argmax} P(d/c)P(c) = \operatorname{argmax} P(x_1 \dots x_n/c)P(c) \approx \operatorname{argmax} P(x_1/c) \dots P(x_n/c)P(c)$
  - $P(c) = \frac{\text{count}(c)}{N}$
  - $P(w/c) = \frac{\text{count}(w,c)+1}{\text{count}(c)+|V|}$

# Google search engine

- Spiders/Googlebot/Google's Web Crawler
- Building index
- PageRank

# Before Indexing

- Tokenize (split words)
- Stop words removal (a, the, ...)
- Stemming (friends, friendship->friend, )
- Normalization (U.S.A->USA)

# Search Index

- Term-document matrix

	D1	D2	...	Dn
W1				
W2				
...				
Wm				

- Posting list
  - Word1 -> d1, d2, ... dj
- Positional indexes
  - Word1 -> d1(p11,...pi1), d2(p12...pi2), ... dj(p1j...pij)

# Ranking

- The more the term appear in the doc the higher the score
- Best known: TF-IDF
  - TF: term frequency of the document
  - Df: document frequency of the term
- How long the Web page has existed?
- Is the term in the title
- Is the term in the URL
- Google uses: The number of other Web pages that link to the page in question (PageRank)
- <http://www.youtube.com/watch?v=BNHR6IQJGZs>



# Personalization

- “Zuckerberg Said,  
‘A Squirrel Dying In Your Front Yard May Be More Relevant To Your Interests Right Now Than People Dying In Africa’”
- Google uses 57 signals like
  - What computer you are using?
  - What browser you are using?
  - Where are you located ?
  - ...
- [http://www.ted.com/talks/lang/en/eli\\_pariser\\_beware\\_online\\_filter\\_bubbles.html](http://www.ted.com/talks/lang/en/eli_pariser_beware_online_filter_bubbles.html)