



POLITÉCNICA

"Ingeniamos el futuro"

CAMPUS  
DE EXCELENCIA  
INTERNACIONAL

# Text Mining 4

# Text Classification

Madrid Summer School 2014  
Advanced Statistics and Data Mining

Florian Leitner  
florian.leitner@upm.es

# Incentive and Applications

- Assign one or more “labels” to a collection of “texts”.
- Spam filtering
- Marketing and politics (“**Opinion/Sentiment Analysis**”)
- Grouping similar items (e.g., “**Recommendation Engines**”)
- Ranking/searching for [topic- or query-specific] documents
- Today’s topics:
  - Document Similarity, Text Classification, Sentiment Analysis

# Document Similarity

- Similarity Measures
  - Cosine Similarity
  - Correlation Coefficients
- Word Vector Normalization
  - TF-IDF
- Latent Semantic Indexing
  - Dimensionality Reduction/Clustering

# Information Retrieval (IR)

## Text Vectorization: Inverted Index

Text 1: He that not wills to the end neither  
wills to the means.

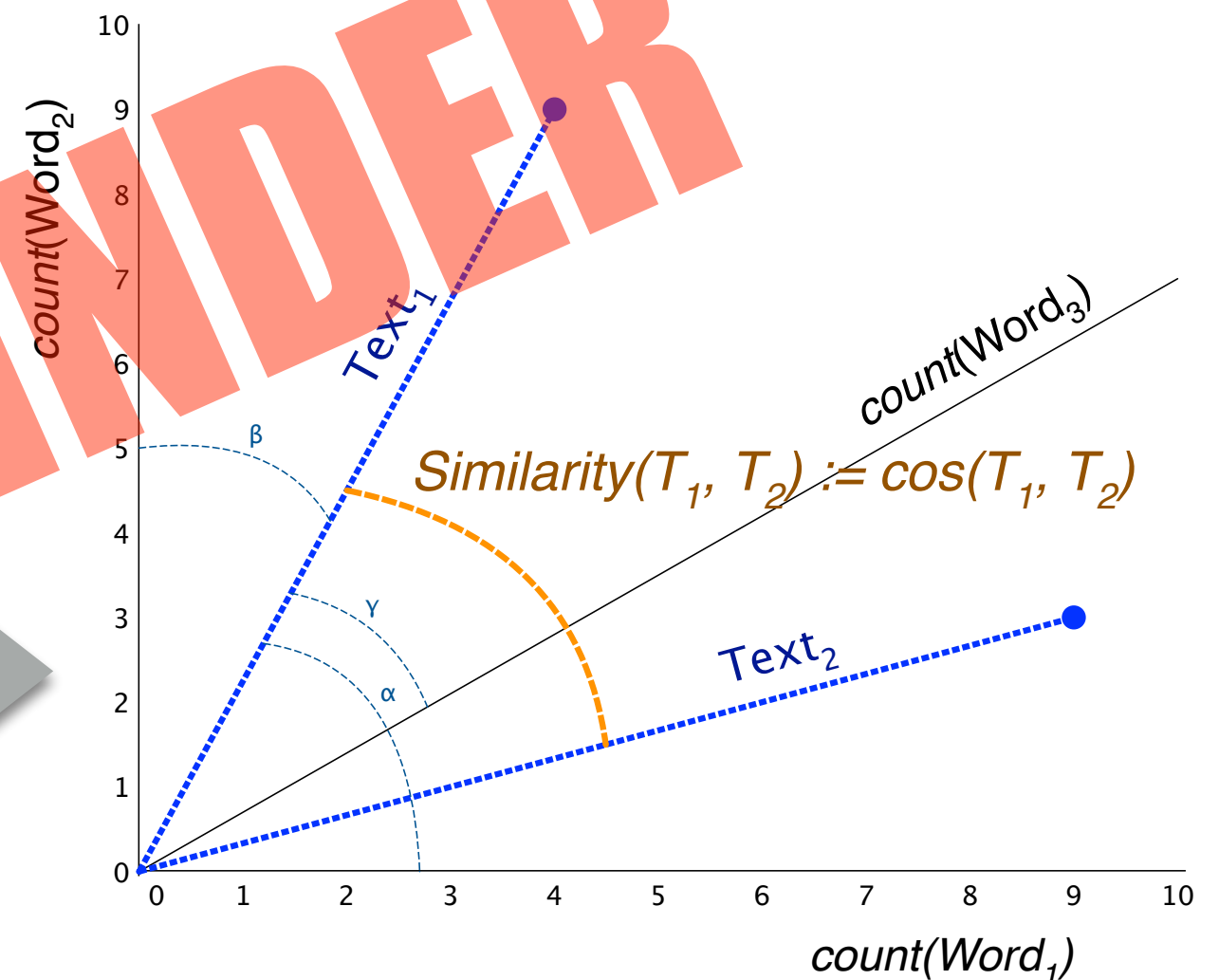
Text 2: If the mountain will not go to Moses,  
then Moses must go to the mountain.

Terms	Doc 1	Doc 2
end	1	0
go	0	2
he	1	0
if	0	1
means	1	0
Moses	0	2
mountain	0	2
must	0	1
not	1	1
that	1	0
the	2	2
then	0	1
to	2	2
will	2	1

document vector

term/word vector

## Comparing Word Vectors: Cosine Similarity



# Cosine Similarity

- Define a similarity score between document vectors (and/or query vectors in Information Retrieval)
- **Euclidian** vector **distance** is **length dependent**
- **Cosine** [angle] between two vectors **is not**

$$\text{sim}(\vec{x}, \vec{y}) = \cos(\vec{x}, \vec{y}) = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| |\vec{y}|} = \frac{\sum x_i y_i}{\sqrt{\sum x_i^2} \sqrt{\sum y_i^2}}$$

*can be dropped if using unit ("length-normalized") vectors  
a.k.a. "cosine normalization"*

# Alternative Similarity Coefficients

- **Spearman's rank correlation coefficient**  $\rho$  ( $r[ho]$ )
  - Ranking is done by term frequency (**TF**; count)
  - Critique: sensitive to ranking differences that are likely to occur with high-frequency words (e.g., "the", "a", ...) → use the **log** of the term count, rounded to two significant digits
    - NB that this is not relevant when only short documents (e.g. titles) with low TF counts are compared
- **Pearson's chi-square test**  $\chi^2$ 
  - Directly on the TFs (counts) - Intuition:  
Are the TFs "random samples" from the same base distribution?
  - Usually,  $\chi^2$  should be preferred over  $\rho$  (Kilgarrieff & Rose, 1998)
    - NB that both measures have no inherent normalization of document size
      - preprocessing might be necessary!

# Term Frequency times Inverse Document Frequency (TF-IDF)

- **Motivation and Background**

- The Problem

- ▶ **Frequent terms** contribute most to a document vector's direction, but **not all** terms are **relevant** ("the", "a", ...).

- The Goal

- ▶ **Separate** important terms from frequent, but **irrelevant terms** in the collection.

- The Idea

- ▶ Frequent **terms** appearing **in all documents** tend to be less important **versus** frequent terms in just a **few documents**. → *Zipf's Law!*

- also **dampens** the effect of **topic-specific noun phrases** or an **author's bias** for a specific set of adjectives

# Term Frequency times Inverse Document Frequency (TF-IDF)

► **tf.idf**( $w$ ) :=  $\text{tf}(w) \times \text{idf}(w)$

- $\text{tf}$ : term frequency
- $\text{idf}$ : inverse document frequency

► **tf<sub>natural</sub>**( $w$ ) :=  $\text{count}(w)$

- $\text{tf}_{\text{natural}}$ : total count of a term in all documents

► **tf<sub>log</sub>**( $w$ ) :=  $\log(\text{count}(w) + 1)$

- $\text{tf}_{\text{log}}$ : the TF is smoothed by taking its log

► **idf<sub>natural</sub>**( $w$ ) :=  $N / \sum^N \{w_i > 0\}$

- $\text{idf}_{\text{natural}}$ : # of documents in which a term occurs

► **idf<sub>log</sub>**( $w$ ) :=  $\log(N / \sum^N \{w_i > 0\})$

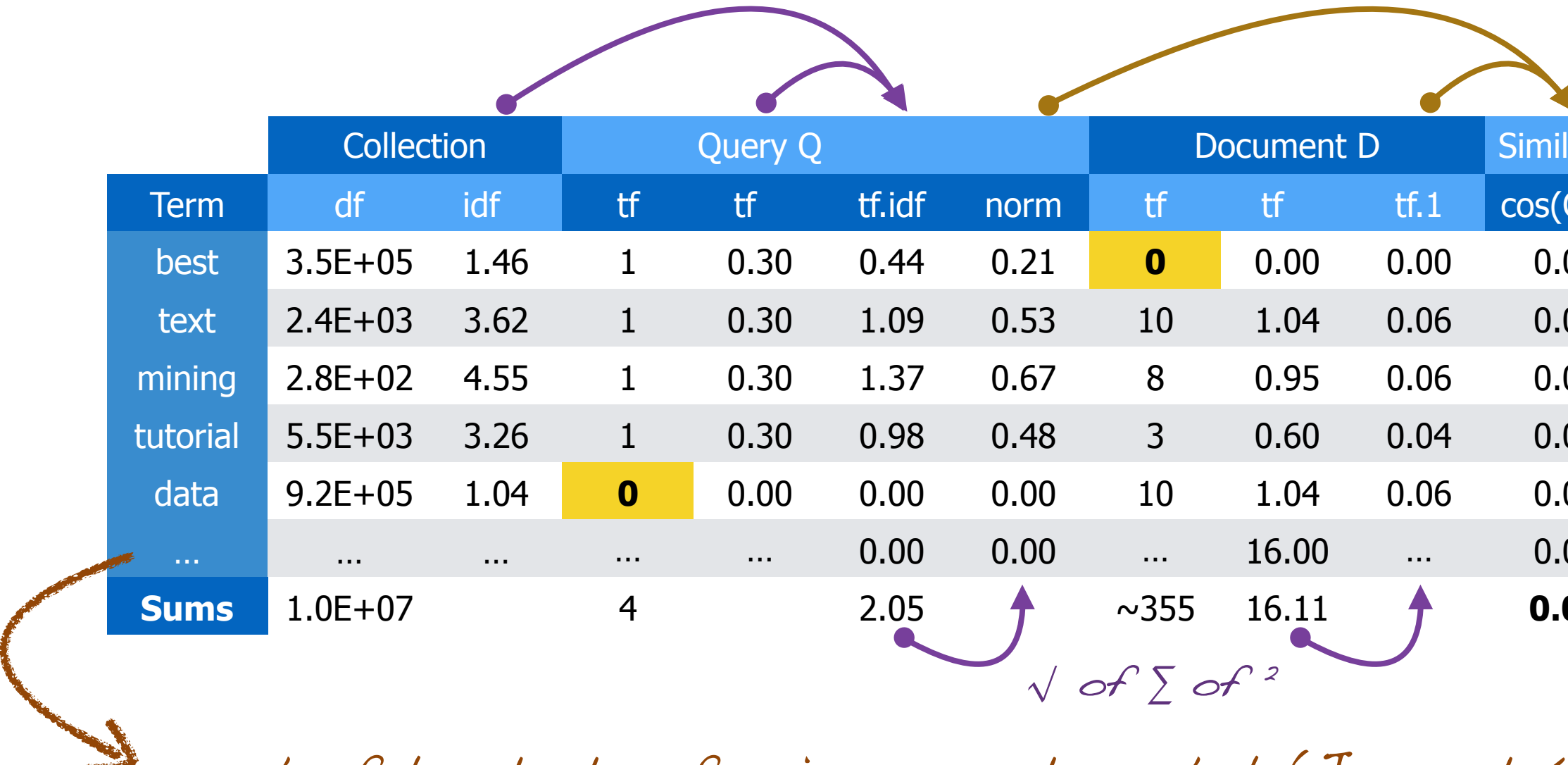
- $\text{idf}_{\text{log}}$ : smoothed IDF by taking its log
- where  $N$  is the **number of documents** and  $w_i$  the **count of word**  $w$  in document  $i$



# TF-IDF in Information Retrieval

- Document Vector =  $tf_{log} \ 1 \rightarrow$  *i.e. the DVs do not use any IDF weighting (simply for efficiency; QV has IDF and gets multiplied with DV values)*
- Query Vector =  $tf_{log} \ idf_{log}$
- terms are counted on each individual document/the query
- Cosine vector length normalization for tf.idf scores:
  - Document  $W$  normalization  $\sqrt{\sum_{w \in W} tf_{log}(w)^2}$
  - Query  $Q$  normalization  $\sqrt{\sum_{q \in Q} (tf_{log}(q) \times idf_{log}(q))^2}$
- IDF is calculated over the indexed collection of all documents

# TF-IDF Query Similarity Calculation Example



	Collection		Query Q				Document D			Similarity
Term	df	idf	tf	tf	tf.idf	norm	tf	tf	tf.1	cos(Q,D)
best	3.5E+05	1.46	1	0.30	0.44	0.21	<b>0</b>	0.00	0.00	0.00
text	2.4E+03	3.62	1	0.30	1.09	0.53	10	1.04	0.06	0.03
mining	2.8E+02	4.55	1	0.30	1.37	0.67	8	0.95	0.06	0.04
tutorial	5.5E+03	3.26	1	0.30	0.98	0.48	3	0.60	0.04	0.02
data	9.2E+05	1.04	<b>0</b>	0.00	0.00	0.00	10	1.04	0.06	0.00
...	...	...	...	...	0.00	0.00	...	16.00	...	0.00
<b>Sums</b>	1.0E+07		4		2.05		~355	16.11		<b>0.09</b>

$\sqrt{\text{of } \sum \text{of }^2}$

*3 out of hundreds of unique words match (Jaccard < 0.03)*

Example idea from: Manning et al. Introduction to Information Retrieval. 2009 *free PDF!*

# From Syntactic to Semantic Similarity

Cosine Similarity,  $\chi^2$ , or Spearman's  $\rho$  all only compare tokens.

[or n-grams!]

But what if you are talking about “automobiles” and I am lazy, calling it a “car”?

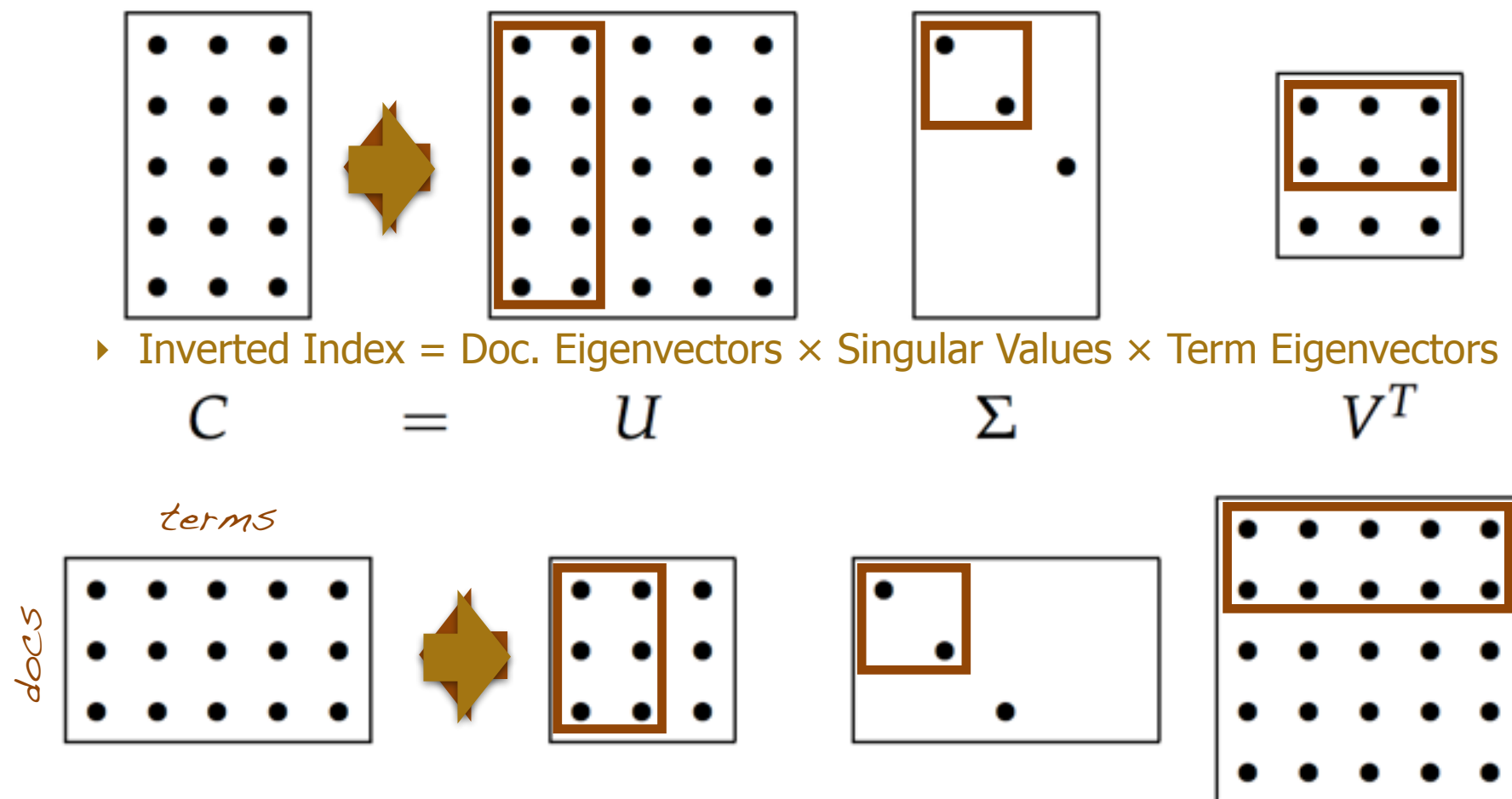
We can solve this with Latent Semantic Indexing!

# Latent Semantic Analysis (LSI 1/3)

- a.k.a. Latent Semantic Indexing (in Text Mining):  
**dimensionality reduction** for **semantic inference**
- Linear Algebra Background
  - ▶ Symmetric Diagonalization of Matrix  $Q$ :  $S = Q\Lambda Q^T$ 
    - symmetric* (pointing to  $Q$ )
    - eigenvectors* (pointing to  $Q$ )
    - eigenvalues (diagonal matrix)* (pointing to  $\Lambda$ )
  - ▶ Singular Value Decomposition of Matrix  $Q$ :  $Q = U\Sigma V^T$ 
    - orthogonal eigenvectors ( $QQ^T$  and  $Q^TQ$ )* (pointing to  $U$  and  $V$ )
    - singular values ( $dm$ )* (pointing to  $\Sigma$ )
- SVD in Text Mining
  - ▶ Inverted Index = Doc. Eigenvectors  $\times$  Singular Values  $\times$  Term Eigenvectors

# Latent Semantic Analysis (LSI 2/3)

$\hat{C}$  = DimRed by selecting only the largest  $n$  eigenvalues



- Image taken from: Manning et al. An Introduction to IR. 2009

# Latent Semantic Analysis (LSI 3/3)

[Spearman's]  $\rho(\text{human}, \text{user}) = -0.38$   
 $\rho(\text{human}, \text{minors}) = -0.29$



**C**

	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	1	0	0	1	0	0	0	0	0
interface	1	0	1	0	0	0	0	0	0
computer	1	1	0	0	0	0	0	0	0
user	0	1	1	0	1	0	0	0	0
system	0	1	1	2	0	0	0	0	0
response	0	1	0	0	1	0	0	0	0
time	0	1	0	0	1	0	0	0	0
EPS	0	0	1	1	0	0	0	0	0
survey	0	1	0	0	0	0	0	0	1
trees	0	0	0	0	0	1	1	1	0
graph	0	0	0	0	0	0	1	1	1
minors	0	0	0	0	0	0	0	1	1

- c1: *Human* machine interface for ABC computer applications  
 c2: A survey of *user* opinion of computer system response time  
 c3: The EPS *user* interface management system  
 c4: System and *human* system engineering testing of EPS  
 c5: Relation of *user* perceived response time to error measurement
- m1: The generation of random, binary, ordered *trees*  
 m2: The intersection *graph* of paths in *trees*  
 m3: *Graph minors* IV: Widths of *trees* and well-quasi-ordering  
 m4: *Graph minors*: A *survey*

**C-hat**

	c1	c2	c3	c4	c5	m1	m2	m3	m4
human	0.16	0.40	0.38	0.47	0.18	-0.05	-0.12	-0.16	-0.09
interface	0.14	0.37	0.33	0.40	0.16	-0.03	-0.07	-0.10	-0.04
computer	0.15	0.51	0.36	0.41	0.24	0.02	0.06	0.09	0.12
user	0.26	0.84	0.61	0.70	0.39	0.03	0.08	0.12	0.19
system	0.45	1.23	1.05	1.27	0.56	-0.07	-0.15	-0.21	-0.05
response	0.16	0.58	0.38	0.42	0.28	0.06	0.13	0.19	0.22
time	0.16	0.58	0.38	0.42	0.28	0.06	0.13	0.19	0.22
EPS	0.22	0.55	0.51	0.63	0.24	-0.07	-0.14	-0.20	-0.11
survey	0.10	0.53	0.23	0.21	0.27	0.14	0.31	0.44	0.42
trees	-0.06	0.23	-0.14	-0.27	0.14	0.24	0.55	0.77	0.66
graph	-0.06	0.34	-0.15	-0.30	0.20	0.31	0.69	0.98	0.85
minors	-0.04	0.25	-0.10	-0.21	0.15	0.22	0.50	0.71	0.62

top 2 dim  
 test # dim  
 to use via  
 synonyms  
 or missing  
 words

From: Landauer et al. An Introduction to LSA. 1998

$\rho(\text{human}, \text{user}) = 0.94$   
 $\rho(\text{human}, \text{minors}) = -0.83$

# Principal Component vs. Latent Semantic Analysis

- **LSA** seeks for the **best linear subspace** in **Frobenius norm**, while **PCA** aims for the **best affine linear subspace**.
- **LSA** (**can**) **use** TF-IDF weighting as **preprocessing** step.
- **PCA requires the** (square) **covariance matrix** of the original matrix as its first step and therefore can only compute term-term or doc-doc similarities.
- **PCA matrices are more dense** (zeros occur only when true independence is detected).

- So far, we have seen how to establish if two documents are syntactically (kNN/LSH) and even semantically (LSI) similar.
- But how do we assign some “label” (or “class”) to a document?
  - ▶ E.g., **relevant**/not relevant; **polarity** (positive, neutral, negative); a **topic** (politics, sport, people, science, healthcare, ...)
  - ▶ We could use the distances (e.g., from LSI) to **cluster** the documents
  - ▶ Instead, let's look at **supervised methods** next.



# Text Classification Approaches

*efficient/fast*

- **Multinomial Naïve Bayes**
- **Nearest Neighbor classification** (ASDM Course 03)
  - ▶ Reminder: **Locality Sensitivity Hashing\*** (see part 3)
- **Latent Semantic Indexing\*** and/or **Clustering\*** (Course 03/10)
- **Maximum Entropy classification**
- **Latent Dirichlet Allocation\***

*your speaker's "favorites"*

*high accuracy*

- **Support Vector Machines** (ASDM Course 11)
- **Random Forests**
- **(Recurrent) Neural Networks** (ASDM Course 05)
- ...

\* (optionally) unsupervised

# Three Text Classifiers

- Multinomial Naïve Bayes
- Maximum Entropy (Multinomial Logistic Regression)
- Latent Dirichlet Allocation → *tomorrow!*

# Bayes' Rule: Diachronic Interpretation

*prior* → *likelihood*

*posterior* →  $P(H|D) = \frac{P(H) \times P(D|H)}{P(D)}$

↑  
"normalizing constant"  
(law of total probability)

REMINDER

$H$  - Hypothesis

$D$  - Data

# Maximum A Posterior (MAP) Estimator

- Issue: Predict the class  $c \in C$  for a given document  $d \in D$
- Solution: MAP, a “perfect” Bayesian estimator:

$$C_{MAP}(d) = \underset{c \in C}{\operatorname{argmax}} \underbrace{P(c|d)}_{\text{the posterior}} = \underset{c \in C}{\operatorname{argmax}} \frac{P(d|c)P(c)}{\cancel{P(d)}_{\text{redundant}}}$$

- Problem:  $d := \{w_1, \dots, w_n\}$  dependent words/features  $W$ 
  - exponential parametrization: one param. for each combination of  $W$  per  $C$

# Multinomial Naïve Bayes Classification

- A simplification of the MAP Estimator
  - ▶  $\text{count}(w)$  is a discrete, **multinomial** variable (unigrams, bigrams, etc.)
  - ▶ Reduce space by making a **strong independence assumption** ("naïve")

*independence assumption*



$$C_{MAP}(d) = \underset{c \in C}{\operatorname{argmax}} P(d|c)P(c) \approx \underset{c \in C}{\operatorname{argmax}} P(c) \prod_{w \in W} P(w|c)$$

*"bag of words/features"*

- Easy **parameter estimation**

$$\hat{P}(w_i|c) = \frac{\text{count}(w_i, c) + 1}{|V| + \sum_{w \in V} \text{count}(w, c)}$$

$\text{count}(w_i, c)$ : the total count of word  $i$  in all documents of class  $c$  [in our training set]

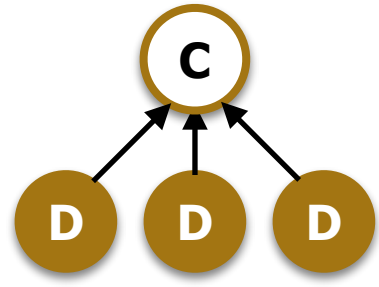
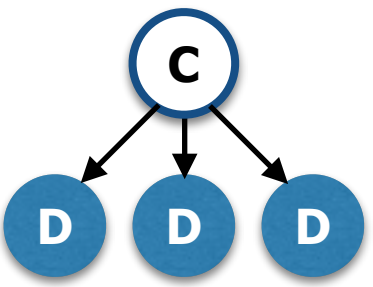
- ▶  $V$  is the entire **vocabulary** (collection of unique words/n-grams/...) in  $D$
- ▶ uses a Laplacian/Add-One Smoothing

# Multinomial Naïve Bayes: Practical Aspects

- Can gracefully handle **unseen words**
- Has **low space requirements**:  $|V| + |C|$  floats *→ sum  $\prod$  using logs!*
- **Irrelevant** (=ubiquitous) **words** cancel each other out
- Opposed to SVM or Nearest Neighbor, it is very **fast**
  - Reminder: the k-shingle LSH approach to NN is fast, too.
  - But Multi-NB will probably result in lower accuracy (⇒ “**baseline**”).
- Each class has its own **n-gram language model**
- **Logarithmic damping** ( $\log(count)$ ) might improve classification

$$\hat{P}(w_i|c) = \frac{\log(count(w_i, c) + 1)}{\log(|V| + \sum_{w \in V} count(w, c))}$$

# Generative vs. Discriminative Models



- **Generative models** describe how the [hidden] labels “generated” the [observed] input as **joint probabilities**:  $P(class, data)$ 
  - Examples: Markov Chain, Naïve Bayes, Latent Dirichlet Allocation, Hidden Markov Model, ...
- **Discriminative models** only predict (“discriminate”) the [hidden] labels **conditioned** on the [observed] input:  $P(class | data)$ 
  - Examples: Logistic Regression, Support Vector Machine, Conditional Random Field, Neural Network, ...
- Both can identify the most likely labels and their likelihoods
- **Only generative models:**
  - Most likely input value[s] and their likelihood[s]
  - Likelihood of input value[s] for some particular label[s]

$$P(H|D) = \frac{P(H) \times P(D|H)}{P(D)}$$

# Maximum Entropy (MaxEnt)

## Intuition

### The Principle of Maximum Entropy

go	{	to	1/5		1/6		??
		up	1/5		1/6		??
		outside	1/5	→	1/6	→	??
		home	1/5		1/4		??
		out	1/5		1/4		??

$w \in W$

$$\sum_{w \in W} P(w | go) = 1$$

$$P(out | go) +$$

$$P(home | go) = 0.5$$

$$P(to | go) +$$

$$P(home | go) = 0.75$$



# Supervised MaxEnt Classification

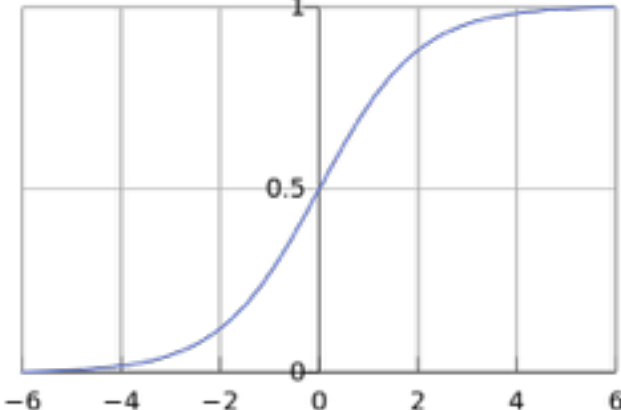
$$p(x) = \frac{1}{1 + \exp(-(\lambda_0 + \lambda_1 x))}$$

*ln* ↗

$$\ln \frac{p(x)}{1 - p(x)} = \lambda_0 + \lambda_1 x$$

↓ *e*

*odds-ratio!* ←  $\frac{p(x)}{1 - p(x)} = \exp(\lambda_0 + \lambda_1 x)$



*logistic function p*

Image Source: WikiMedia Commons, Qef

- a.k.a. Multinomial Logistic Regression
- **Does not assume independence between the features**
- Can **model mixtures of** binary, discrete, and real **features**
- Training data are per-feature-label probabilities:  $P(F, L)$ 
  - ▶ I.e.,  $\text{count}(f_i, l_i) \div \sum_{i=1}^N \text{count}(f_i, l_i)$
- ➔ words → very sparse training data (zero or few examples)
- The model is commonly “learned” using gradient descent
  - ▶ Expensive if compared to Naïve Bayes, but efficient optimizers exist (**L-BFGS**)

# Example Feature Functions for MaxEnt Classifiers

- Examples of **indicator functions** (a.k.a. **feature functions**)
  - Assume we wish to classify the general polarity (positive, negative) of product reviews:
- $f(c, w) := \{c = \text{POSITIVE} \wedge w = \text{"great"}\}$ 
  - Equally, for classifying words in a text, say to detect proper names, we could create a feature:
- $f(c, w) := \{c = \text{NAME} \wedge \text{isCapitalized}(w)\}$
- Note that while we can have multiple classes, we cannot require more than one class in the whole match condition of a single indicator (feature) function.

*NB: typical text mining models can have a million or more features:  
unigrams + bigrams + trigrams + counts + dictionary matches + ...*

# Maximizing Conditional Entropy

- The **conditioned** (on  $X$ ) version of Shannon's **entropy**  $H$ :

*NB: the chain rule*

$$P(x, y) = P(x) P(y|x)$$

$$\begin{aligned} H(Y|X) &= - \sum_{x \in X} P(x) H(Y|X = x) \\ &= - \sum_{x \in X} P(x) \sum_{y \in Y} P(y|x) \log_2 P(y|x) \\ &\quad \text{(swapped nom/denom to remove the minus)} \\ &= \sum_{x, y \in X, Y} P(x, y) \log_2 \frac{P(x)}{P(x, y)} \end{aligned}$$

- MaxEnt **training** then is about selecting the model  $p^*$  that maximizes  $H$ :

$$p^* = \underset{p \in P}{\operatorname{argmax}} H(P) = \underset{p \in P}{\operatorname{argmax}} H(Y|X)$$

# Maximum Entropy (MaxEnt 1/2)

- Some definitions:

- ▶ The observed probability of  $y$  (the class) with  $x$  (the words) is:

$$\hat{P}(x, y) = \text{count}(x, y) \div N$$

- ▶ An **indicator function** ("**feature**") is defined as a binary valued function that returns 1 iff class and data match the **indicated** requirements (**constraints**):

$$f(x, y) = \begin{cases} 1 & \text{if } y = c_i \wedge x = w_i \\ 0 & \text{otherwise} \end{cases}$$

*real/discrete/binary features now are all the same!*

- ▶ The probability of a feature with respect to the observed distribution is:

$$\hat{P}(f_i, X, Y) = E_{\hat{P}}[f_i] = \sum \hat{P}(x, y) f_i(x, y)$$

# Getting lost?

## Reality check:

- I have told you:
  - ▶ MaxEnt is about to maximize “conditional entropy”:
  - ▶ By multiplying binary (0/1) feature functions for observations with the joint (observation, class) probabilities, we can calculate the conditional probability of a class given its observations  $H(Y|X)$
- We will still have to do:
  - ▶ Find weights for each feature [function] that lead to the best model of the [observed] class probabilities.
- And you want to know:
  - ▶ How the heck do we actually classify stuff???

# Maximum Entropy (MaxEnt 2/2)

- ▶ In a **linear** model, we'd use weights ("lambdas") that identify the most relevant features of our model, i.e., we use the following MAP to select a class:

$$\operatorname{argmax}_{y \in Y} \sum \lambda_i f_i(x, y)$$

- ▶ To do **multinomial logistic** regression, expand with a **linear combination**:

$$\operatorname{argmax}_{y \in Y} \frac{\exp(\sum \lambda_i f_i(x, y))}{\sum_{y \in Y} \exp(\sum \lambda_i f_i(x, y))} \quad \text{"exponential model"}$$

- ▶ Next: **Estimate** the  $\lambda$  weights (parameters) that **maximize** the conditional **likelihood** of this logistic model (**MLE**)

# Maximum Entropy (MaxEnt 2/2) [again]

- ▶ In summary, MaxEnt is about selecting the “maximal” model  $p^*$ :

$$p^* = \operatorname{argmax}_{p \in P} - \sum_{x \in X} p(x) \sum_{y \in Y} p(y|x) \log_2 p(y|x)$$

*select some model that maximizes the conditional entropy...*

- ▶ That obeys the following conditional equality constraint:

$$\sum_{x \in X} P(x) \sum_{y \in Y} P(y|x) f(x, y) = \sum_{x \in X, y \in Y} P(x, y) f(x, y)$$

*...using a conditional model that matches the (observed) joint probabilities*

- ▶ Next: Using, e.g., **Lagrange multipliers**, one can establish the optimal  $\lambda$  parameters of the model that maximize the entropy of this probability:

$$p^*(y|x) = \frac{\exp(\sum \lambda_i f_i(x, y))}{\sum_{y \in Y} \exp(\sum \lambda_i f_i(x, y))}$$

# Newton's Method for Parameter Optimization

- Problem: find the  $\lambda$  parameters
  - an “**optimization problem**”
- MaxEnt surface is **concave**
  - one **single maximum**
- Using **Newton's method**
  - iterative, hill-climbing search for max.
  - the **first derivative**  $f'$  is zero at the [global] maximum (the “goal”)
  - the **second derivative**  $f''$  indicates rate of change:  $\Delta\lambda_i$  (search direction)
  - takes the most direct route to the maximum *as opposed to gradient descent, which will follow a possibly curved path to the optimum*
- Using **L-BFGS**
  - a **heuristic** to simplify Newton's method *it is said to be “quasi-Newtonian”*
  - L-BFGS: **limited memory** Broyden–Fletcher–Goldfarb–Shanno
  - normally, the **partial second derivatives** would be stored in the **Hessian**, a matrix that **grows quadratically** with respect to the number of features
  - only uses the last few [partial] gradients to **approximate the search direction**



# MaxEnt vs. Naïve Bayes

MaxEnt model  
(as observed)

Lights Working

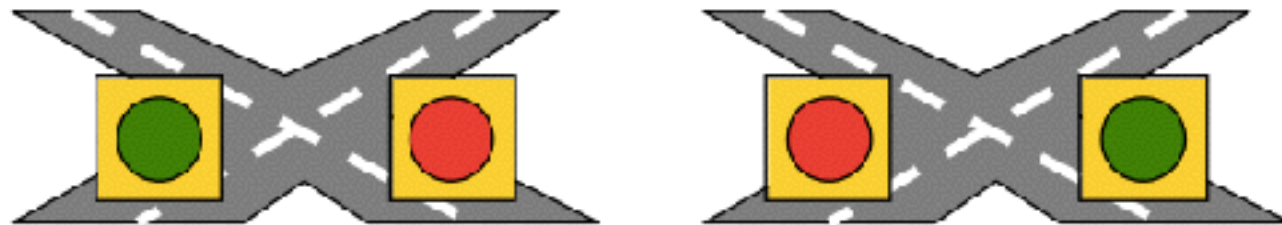
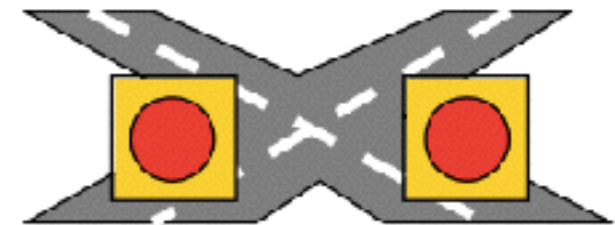


Image Source: Klein & Manning. Maxent Models, Conditional Estimation, and Optimization. ACL 2003 Tutorial

$$P(\text{g}, \text{r}, \text{w}) = 3/7$$

$$P(\text{r}, \text{g}, \text{w}) = 3/7$$

Lights Broken



$$P(\text{r}, \text{r}, \text{b}) = 1/7$$

MaxEnt adjusts the Lagrange multipliers (weights) to **model** the correct (observed) **joint probabilities**.

But even MaxEnt cannot model **feature interaction**!

F	a		F <sub>1</sub> = 2/3 F <sub>2</sub> = 2/3 4/9 2/9 2/9 1/9
	b	b	
b	1	1	
b	1	0	

*NB the example has dependent features: the two stoplights*

Naïve Bayes model

- $P(\text{w}) = 6/7$
  - $P(\text{r}|\text{w}) = 1/2$
  - $P(\text{g}|\text{w}) = 1/2$
  - $P(\text{b}) = 1/7$
  - $P(\text{r}|\text{b}) = 1$
  - $P(\text{g}|\text{b}) = 0$
  - $P(\text{r}, \text{r}, \text{b}) = (1/7)(1)(1) = 4/28$
  - $P(\text{r}, \text{g}, \text{b}) = P(\text{g}, \text{r}, \text{b}) = P(\text{g}, \text{g}, \text{b}) = 0$
  - $P(*, *, \text{w}) = (6/7)(1/2)(1/2) = 6/28$
- $P(\text{g}, \text{g}, \text{w}) = 6/28???$*

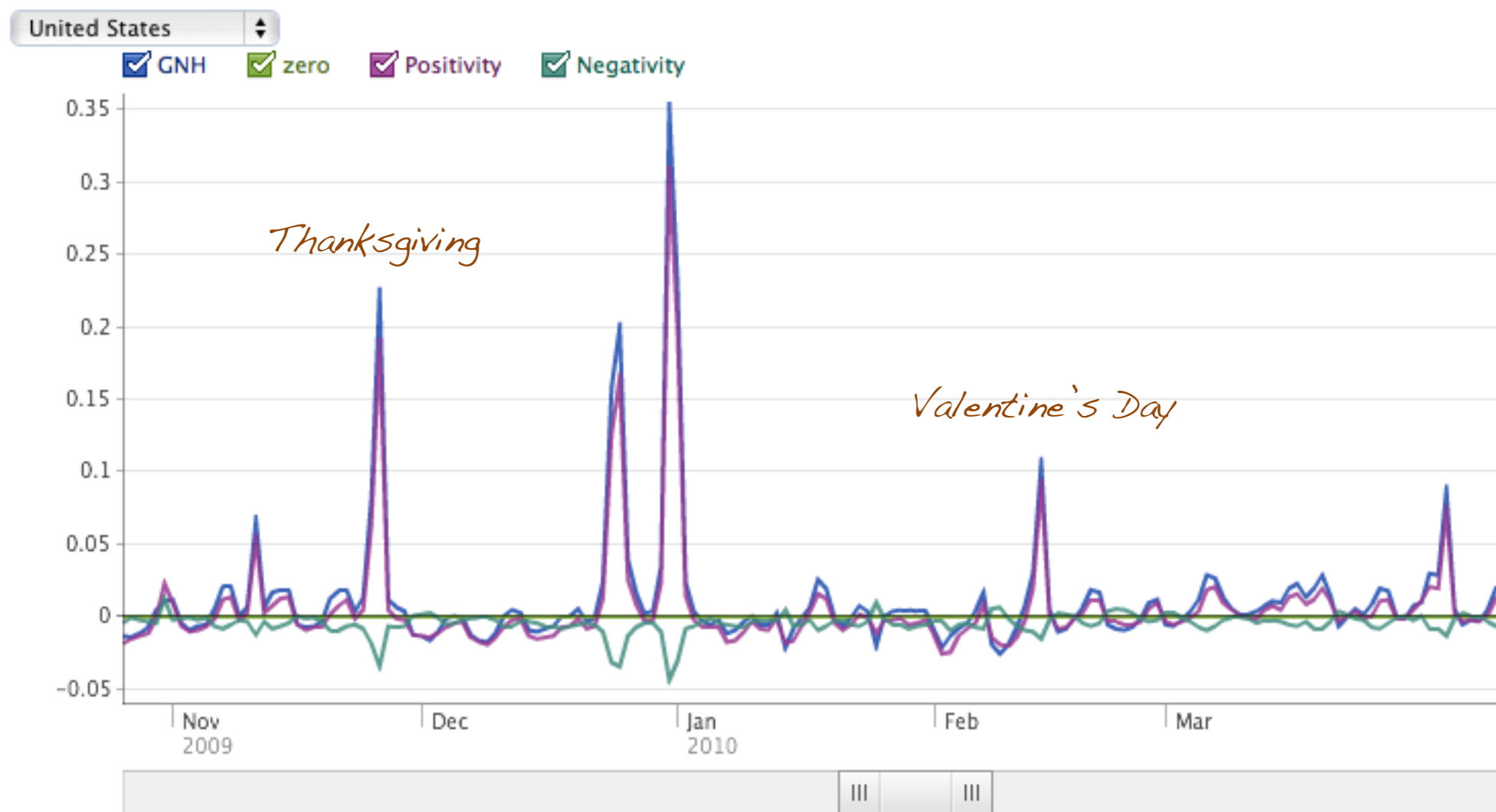
Klein & Manning. MaxEnt Models, Conditional Estimation and Optimization. ACL 2003

# Sentiment Analysis

as an example **domain** for text classification

Cristopher Potts. Sentiment Symposium Tutorial. 2011  
<http://sentiment.christopherpotts.net/index.html>

# Facebook's "Gross National Happiness" Timeline



Source: Cristopher Potts. Sentiment tutorial. 2011

# Opinion/Sentiment Analysis

- Harder than “regular” document classification
  - irony, neutral (“non-polar”) sentiment, negations (“not good”), syntax is used to express emotions (“!”), context dependent
- Confounding polarities from individual aspects (phrases)
  - e.g., a car company’s “customer service” vs. the “safety” of their cars
- Strong commercial interest in this topic
  - “Social” (commercial?) networking sites (FB, G+, ...; advertisement)
  - Reviews (Amazon, Google Maps), blogs, fora, online comments, ...
  - Brand reputation and political opinion analysis

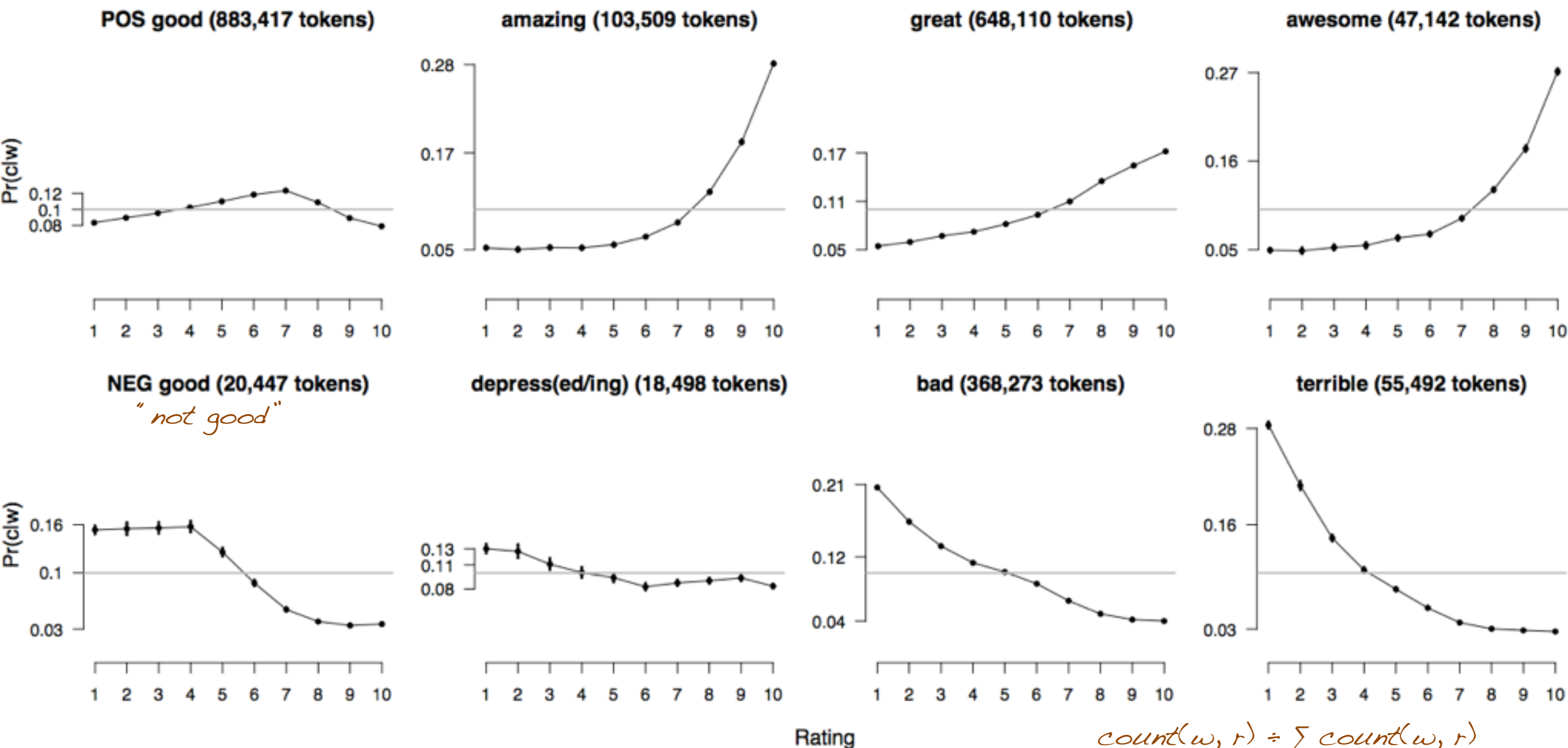
# 5+1 Lexical Resources for Sentiment Analysis

Cristopher Potts. Sentiment Symposium Tutorial. 2011

Disagree- ment	Opinion Lexicon	General Inquirer	SentiWordNet	LIWC
Subjectivity Lexicon	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%)
SentiWord Net				174/694 (25%)

MPQA Subjectivity Lexicon: <http://mpqa.cs.pitt.edu/>  
 Liu's Opinion Lexicon: <http://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html>  
 General Inquirer: <http://www.wjh.harvard.edu/~inquirer/>  
 SentiWordNet: <http://sentiwordnet.isti.cnr.it/>  
 LIWC (commercial, \$90): <http://www.liwc.net/>  
 NRC Emotion Lexicon (+1): <http://www.saifmohammad.com/> (➡Publications & Data)

# Polarity of Sentiment Keywords in IMDB

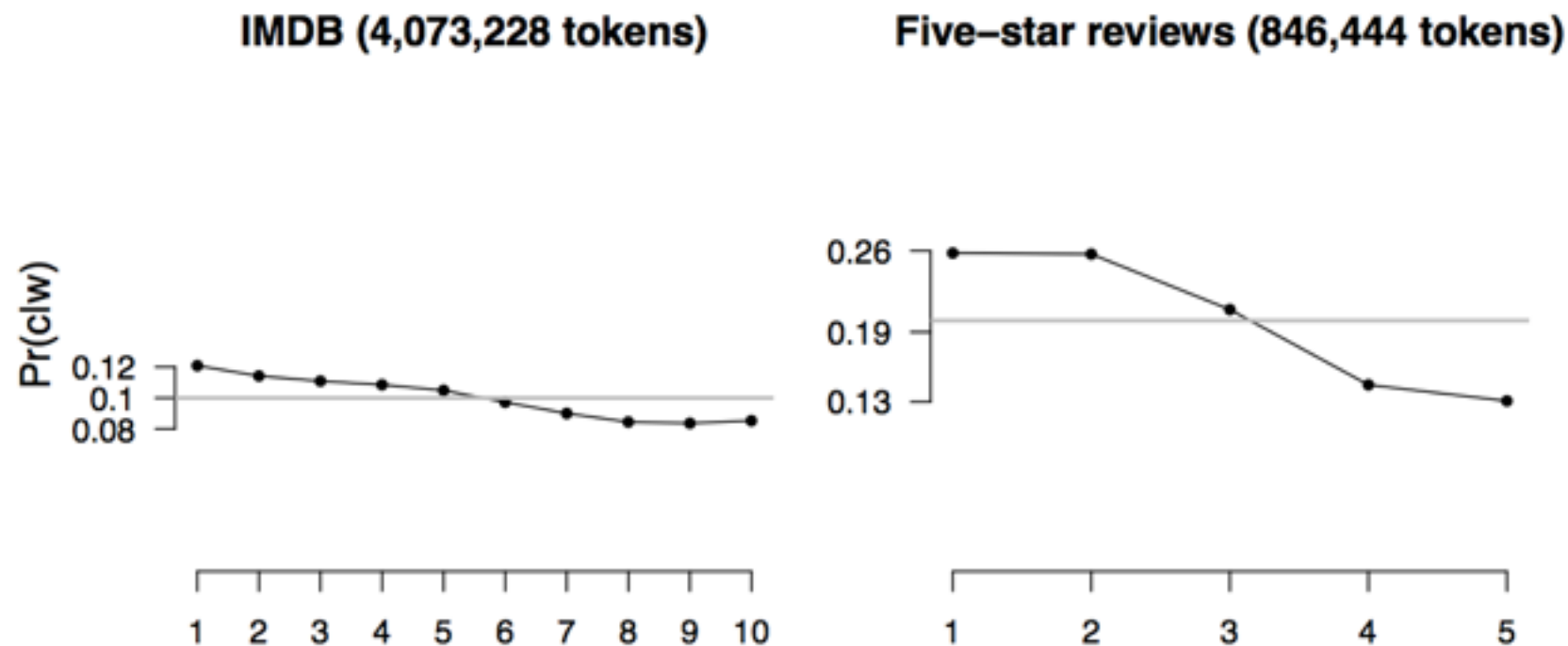


Cristopher Potts. On the negativity of negation. 2011

**Note:**  $P(\text{rating} | \text{word}) = P(\text{word} | \text{rating}) \div P(\text{word})$

# Negations

- Is the use of negation associated to polarity?



Cristopher Potts. On the negativity of negation. 2011

- **Yes, but** there are far more deeper issues going on...

# Detecting the Sentiment of Individual Aspects

- Goal: Determine the sentiment for a particular aspect or establish their polarity.
  - An “aspect” here is a phrase or concept, like “customer service”.
  - “They have a **great**+ customer service team, but the delivery **took ages**-.”
- Solution: Measure the co-occurrence of the aspect with words of distinct sentiment or relative co-occurrence with words of the same polarity.
  - The “sentiment” keywords are taken from some lexical resource.



# Google's Review Summaries


**Google maps** peter lugar nyc Search Maps Show search options  
Find businesses, addresses and places of interest.

[Get Directions](#) [My Maps](#) Edit this place - Business owner? << Print Email Link


### Peter Luger Steakhouse

178 Broadway, Brooklyn, NY  
(718) 387-7400  
[peterluger.com](http://peterluger.com)  
★★★★☆ 402 reviews  
[Directions](#) [Search nearby](#) [Save to...](#) more▼

**Category:** Steak House  
**Hours:** Today 11:30am – 9:45pm  
**Transit:** [Marcy Ave Station](#) (0.3 mi) J M Z



Panoramio



"The steak is great, but the ambience is awesome" - insiderpages.com ... "Don't come looking for ambience, service or value!" - citysearch.com ... "For \$100 plus meal for two I expect much better service and much better food" - tripadvisor.com ... "The generous portions, side dishes, and service are top notch!" - dine.com ... "Steak Heaven" - citysearch.com

**Details**

<b>Cuisine:</b> Specialties, Steakhouse	<b>Prices:</b> Very Expensive
<b>Menu:</b> <a href="#">Menu</a>	<b>Atmosphere:</b> Old City Feel, Power Scene
<b>Parking:</b> Parking lot	<b>Reservations Policy:</b> Required
<b>Attire:</b> Business casual	<b>Year Opened:</b> 1887
<b>Neighborhood:</b> Williamsburg	<b>Meals Served:</b> Dinner, Lunch

[gayot.com](http://gayot.com), [metromix.com](http://metromix.com), [nymag.com](http://nymag.com), [zagat.com](http://zagat.com)  
[More details »](#)

**What people are saying about**

<a href="#">steak</a>	<div><div style="width: 100%;"></div></div>	"The steak is great, but the ambience is awesome." - insiderpages.com
<a href="#">service</a>	<div><div style="width: 75%;"></div></div>	"Steak, Steak or Steak. Service is okay." - citysearch.com
<a href="#">food</a>	<div><div style="width: 100%;"></div></div>	"For \$100 plus meal for two I expect much better service and much better food." - tripadvisor.com
<a href="#">meal</a>	<div><div style="width: 100%;"></div></div>	"My husband would eat steak every meal for the rest of his life if he could." - tripadvisor.com
<a href="#">atmosphere</a>	<div><div style="width: 75%;"></div></div>	"Great food, great atmosphere." - virtualtourist.com
<a href="#">dining_decor_dishes_ambience_ambiance</a>		

**Sponsored Links**

[\*\*Ruth's Chris Steak House\*\*](#)  
Dine on US Prime Steak at Ruth's Chris. Official Site. See Our Menu!  
[www.RuthsChris.com](http://www.RuthsChris.com)

[\*\*Restaurant Supplies\*\*](#)  
Looking For Restaurant Equipment & Supplies? Find Them Here Now.  
[www.ThomasNet.com](http://www.ThomasNet.com)

[\*\*Commercial Kitchen Design\*\*](#)  
Restaurant Kitchen Design  
Design, Equipment and Installation  
[www.BostonShowcase.com](http://www.BostonShowcase.com)

[\*\*Arabic Classes NYC\*\*](#)  
ABC Language Exchange  
Group & Private Arabic Lessons NYC  
[www.abclang.com](http://www.abclang.com)

# Point-wise Mutual Information

- Mutual Information measures the **dependence of two variables**.

$$I(X; Y) = \sum_Y \sum_X P(x, y) \log_2 \frac{P(x, y)}{P(x)P(y)}$$

- Point-wise MI: MI of two **individual** events only
  - e.g., neighboring words, phrases in a document, ...
  - even a mix of two co-occurring events (e.g. a word and a phrase)

*NB: define a maximum distance between  $w_1$  and  $w_2$*

$$PMI(w_1, w_2) = \log_2 \frac{P(w_1, w_2)}{P(w_1)P(w_2)}$$

*$w_1$ : a phrase/aspect  
 $w_2$ : one of a set of pos./neg. sentiment words*

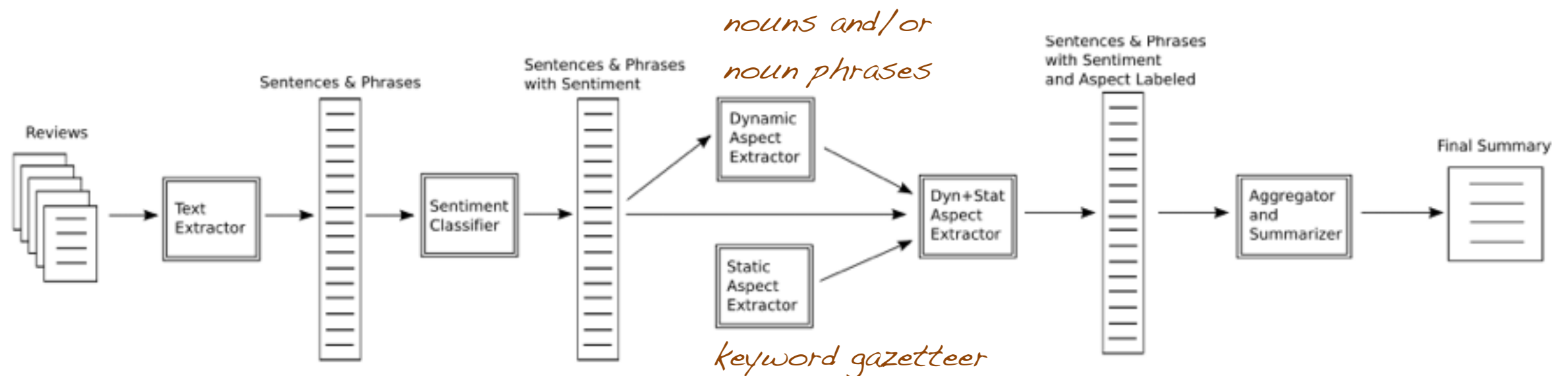
- Can be normalized to a [-1,1] range:

$$\frac{PMI(w_1, w_2)}{-\log_2 P(w_1, w_2)}$$

- -1: the two words/phrases/events do not occur together; 0: the words/phrases/events are independent; +1: the words/phrases/events always co-occur

# Using PMI to Detect Aspect Polarity

- **Polarity(aspect)** :=  $\text{PMI}(\text{aspect}, \text{pos-sent-kwds}) - \text{PMI}(\text{aspect}, \text{neg-sent-kwds})$ 
  - Polarity > 0 = positive sentiment
  - Polarity < 0 = negative sentiment
- Google's approach:



- Blair-Goldensohn et al. Building a Sentiment Summarizer for Local Service Reviews. WWW 2008

# Subjectivity Clues and Polarity Intensifiers

*subjectivity clues*

(4) Philip Clapp, president of the National Environment Trust, sums up well the general thrust of the reaction of environmental movements: “There is no reason at all to believe that the polluters are suddenly going to become reasonable.” *unmodified*

*long-distance negations*

*polarity intensifiers*

- Known as the problem of “**contextual polarity**”:
  - ▶ The **evil** baron was **held in check**. (“evil”=neg. subjective expression, “held in check”=context reverses neg. subjective expression)
- Wilson et al. Recognizing Contextual Polarity in Phrase-Level Sentiment Analysis. HLT-EMNLP 2005
  - ▶ MPQA Subjectivity Lexicon: [http://mpqa.cs.pitt.edu/lexicons/subj\\_lexicon/](http://mpqa.cs.pitt.edu/lexicons/subj_lexicon/)

# More Context Issues...

- ▶ **Feelings**: "it is too bad", "I am very sorry for you", ...
- ▶ **Agreement**: "I must agree to", "it is the case that", ...
- ▶ **Hedging**: "a little bit", "kind/sort of", ...
- ▶ ....

*(not part of this course)*

- Requires the use of **dependency parsing** to detect the relevant **context**

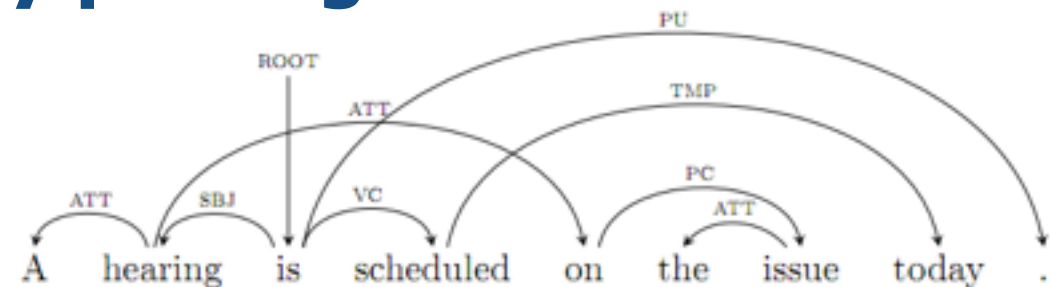


Image Source: WikiBooks, Daniele Pighin

*The Take-Home Message:*

- Inter-rater **human agreement** for sentiment tasks
  - ▶ often only at around 80% (Cohen's  $\kappa \sim 0.7$ )
  - ▶ sentiment expressions have a very high "uncertainty" (ambiguity)



# Evaluation Metrics

Evaluation is all about answering questions like:

How to measure a change to an approach?

Did adding that feature improve or decrease performance?

Is the approach good at locating the relevant pieces or good at excluding the irrelevant bits?

How to compare entirely different methods?

*this one is really hard...*

# Basic Evaluation Metrics: Accuracy, F-Measure, MCC Score

- **True/False  
Positive/Negative**

- counts; TP, TN, FP, FN

- **Precision** (P)

- correct hits [TP] ÷  
all hits [TP + FP]

- **Recall** (R; **Sensitivity**, TPR)

- correct hits [TP] ÷  
true cases [TP + FN]

- **Specificity** (True Negative Rate)

- correct misses [TN] ÷  
negative cases [FP + TN]

*NB: no result order:  
lesson 5 (tomorrow)*

- **Accuracy**

- correct classifications [TP + TN] ÷  
all cases [TP + TN + FN + FP]
- highly **sensitive to** class **imbalance**

- **F-Measure** (F-Score)

- the harmonic mean between P & R  
 $= 2 TP \div (2 TP + FP + FN)$   
 $= (2 P R) \div (P + R)$

- does **not** require a **TN** count

- **MCC Score** (Mathew's  
Correlation Coefficient)

- $\chi^2$ -**based**:  $(TP TN - FP FN) \div$   
 $\text{sqrt}[(TP+FP)(TP+FN)(TN+FP)(TN+FN)]$
- **robust against** class **imbalance**

# Practical: Twitter Sentiment Detection

- Implement a MaxEnt classifier to detect the sentiment of Twitter “tweets” for Apple/iPhone and Google/Android products.
- Try to improve the result of 70 % accuracy by choosing better features and implementing a better tokenization strategy.
- Experiment with making use of the sentiment clues and polarity intensifiers from the Subjectivity Lexicon.



“Romance should never begin with sentiment.  
It should begin with science and end with a settlement.”

Oscar Wilde, 1895