

---

---

# Introduction to Information Retrieval

— Text Classification & Naive Bayes —

---

---

# Outline

- **Text classification**
- Naive Bayes
- Evaluation of TC
- Feature Selection

# A text classification task: Email spam filtering

From: "" <takworld@hotmail.com>

Subject: real estate is the only way... gem oalvgkay

Anyone can buy real estate with no money down

Stop paying rent TODAY !

There is no need to spend hundreds or even thousands for similar courses

I am 22 years old and I have already purchased 6 properties using the methods outlined in this truly INCREDIBLE ebook.

Change your life NOW !

=====

Click Below to order:

<http://www.wholesaledaily.com/sales/nmd.htm>

=====

How would you write a program that would automatically detect and delete this type of message?

# Formal definition of TC: Training

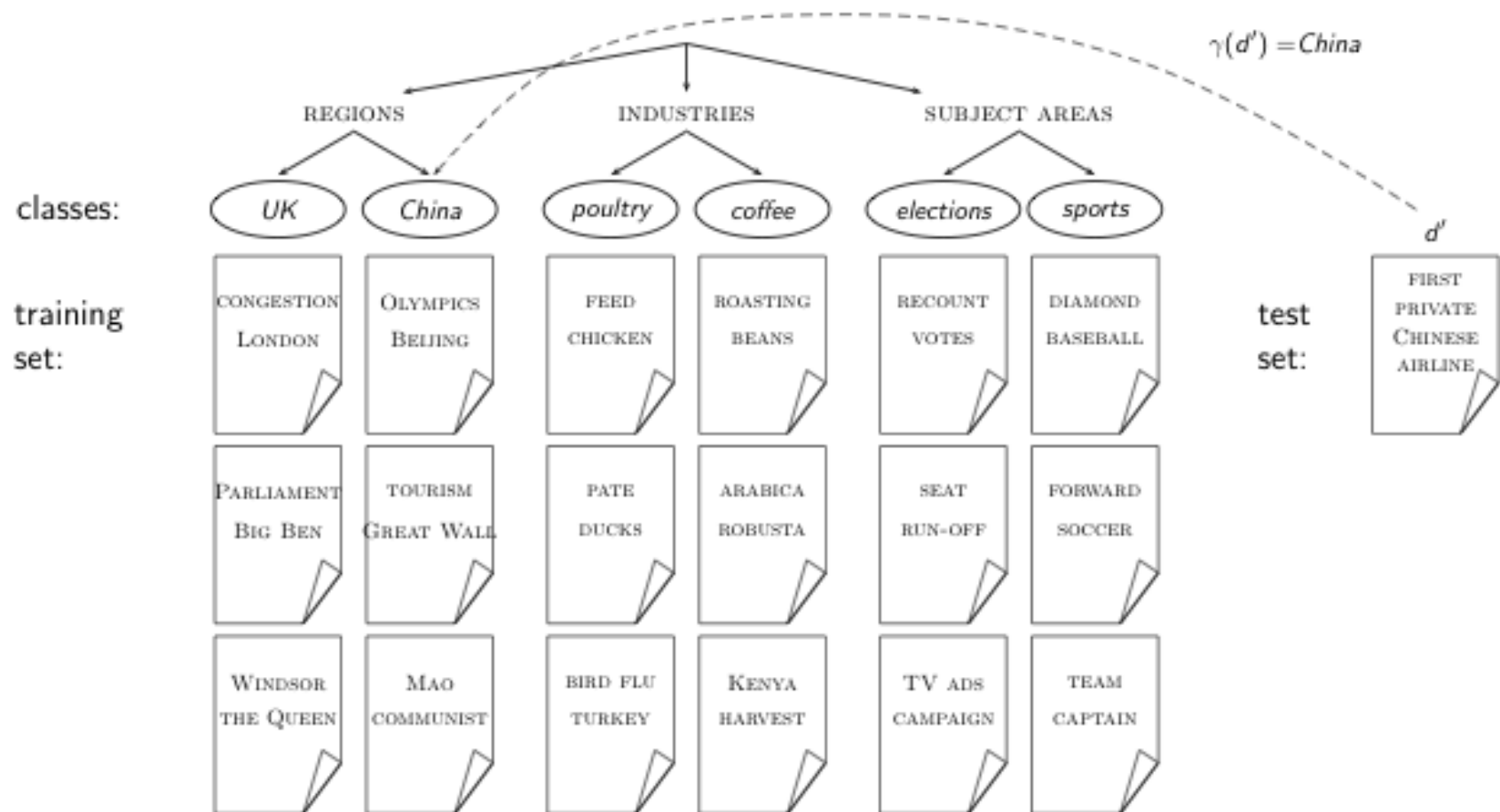
- Given
  - A document space  $X$ ,
    - Documents are represented in this space – typically some type of high-dimensional space.
  - A fixed set of classes  $C = \{c_1, c_2, \dots, c_J\}$ 
    - The classes are human-defined for the needs of an application (e.g., relevant vs. nonrelevant).
  - A training set  $D$  of labeled documents with each labeled document  $\langle d, c \rangle \in X \times C$
- Using a learning method or learning algorithm, we then wish to learn a classifier  $\Upsilon$  that maps documents to classes:

$$\Upsilon : X \rightarrow C$$

# Formal definition of TC: Application/Testing

- Given: a description  $d \in X$  of a document, determine:  $\Upsilon(d) \in C$ , that is, the class that is most appropriate for  $d$

# Topic classification



# Examples of how search engines use classification

- Language identification (classes: English vs. French etc.)
- The automatic detection of spam pages (spam vs. nonspam)
- The automatic detection of sexually explicit content (sexually explicit vs. not)
- Topic-specific or *vertical* search – restrict search to a “vertical” like “related to health” (relevant to vertical vs. not)
- Standing queries (e.g., Google Alerts)
- Sentiment detection: is a movie or product review positive or negative (positive vs. negative)

# Classification methods: 1. Manual

- Manual classification was used by Yahoo in the beginning of the web. Also: ODP, PubMed
- Very accurate if job is done by experts
- Consistent when the problem size and team is small
- Scaling manual classification is difficult and expensive.

→ We need automatic methods for classification.



## Classification methods: 2. Rule-based

- Our Google Alerts example was rule-based classification.
- There are IDE-type development environments for writing very complex rules efficiently. (e.g., Verity)
- Often: Boolean combinations (as in Google Alerts)
- Accuracy is very high if a rule has been carefully refined over time by a subject expert.
- Building and maintaining rule-based classification systems is cumbersome and expensive.

# A Verity topic (a complex classification rule)

comment line	# Beginning of art topic definition		
top-level topic	art ACCRUE		
topic definition modifiers	/author = "fsmith" /date = "30-Dec-01" /annotation = "Topic created by fsmith"		
subtopic		subtopic	* 0.70 film ACCRUE
subtopic	* 0.70 performing-arts ACCRUE		** 0.50 STEM
evidencetopic	** 0.50 WORD		/wordtext = film
topic definition modifier	/wordtext = ballet	subtopic	** 0.50 motion-picture PHRASE
evidencetopic	** 0.50 STEM		*** 1.00 WORD
topic definition modifier	/wordtext = dance		/wordtext = motion
evidencetopic	** 0.50 WORD		*** 1.00 WORD
topic definition modifier	/wordtext = opera		/wordtext = picture
evidencetopic	** 0.30 WORD		** 0.50 STEM
topic definition modifier	/wordtext = symphony		/wordtext = movie
subtopic	* 0.70 visual-arts ACCRUE	subtopic	* 0.50 video ACCRUE
	** 0.50 WORD		** 0.50 STEM
	/wordtext = painting		/wordtext = video
	** 0.50 WORD		** 0.50 STEM
	/wordtext = sculpture		/wordtext = vcr
			# End of art topic

# Classification methods: 3. Statistical/Probabilistic

- This was our definition of the classification problem – text classification as a learning problem
- Supervised learning of the classification function  $Y$  and its application to classifying new documents
- Naive Bayes, Rocchio, kNN, SVMs
- No free lunch: requires hand-classified training data
- But this manual classification can be done by non-experts.

# Outline

- Text classification
- **Naive Bayes**
- Evaluation of TC
- Feature Selection

# The Naive Bayes classifier

- The Naive Bayes classifier is a probabilistic classifier.
- We compute the probability of a document  $d$  being in a class  $c$  as follows:

$$P(c|d) \propto P(c) \prod_{1 \leq k \leq n_d} P(t_k|c)$$

- $n_d$  is the length of the document. (number of tokens)
- $P(t_k | c)$  is the conditional probability of term  $t_k$  occurring in a document of class  $c$
- $P(t_k | c)$  as a measure of how much evidence  $t_k$  contributes that  $c$  is the correct class.
- $P(c)$  is the prior probability of  $c$ .
- If a document's terms do not provide clear evidence for one class vs. another, we choose the  $c$  with highest  $P(c)$ .

# Maximum a posteriori class

- Our goal in Naive Bayes classification is to find the “best” class.
- The best class is the most likely or maximum a posteriori (MAP) class

$c_{\text{map}}$ :

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} \hat{P}(c|d) = \arg \max_{c \in \mathbb{C}} \hat{P}(c) \prod_{1 \leq k \leq n_d} \hat{P}(t_k|c)$$

# Taking the log

- Multiplying lots of small probabilities can result in floating point underflow.
- Since  $\log(xy) = \log(x) + \log(y)$ , we can sum log probabilities instead of multiplying probabilities.
- Since log is a monotonic function, the class with the highest score does not change.
- So what we usually compute in practice is:

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} [\log \hat{P}(c) + \sum_{1 \leq k \leq n_d} \log \hat{P}(t_k | c)]$$

# Naive Bayes classifier

- Classification rule:

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} [\log \hat{P}(c) + \sum_{1 \leq k \leq n_d} \log \hat{P}(t_k | c)]$$

- Simple interpretation:

- Each conditional parameter  $\log \hat{P}(t_k | c)$  is a weight that indicates how good an indicator  $t_k$  is for  $c$ .
- The prior  $\log \hat{P}(c)$  is a weight that indicates the relative frequency of  $c$ .
- The sum of log prior and term weights is then a measure of how much evidence there is for the document being in the class.
- We select the class with the most evidence.



# Parameter estimation take 1: Maximum likelihood

- Estimate parameters  $\hat{P}(c)$  and  $\hat{P}(t_k|c)$  from train data: How?

- Prior:

$$\hat{P}(c) = \frac{N_c}{N}$$

- $N_c$  : number of docs in class  $c$ ;  $N$ : total number of docs

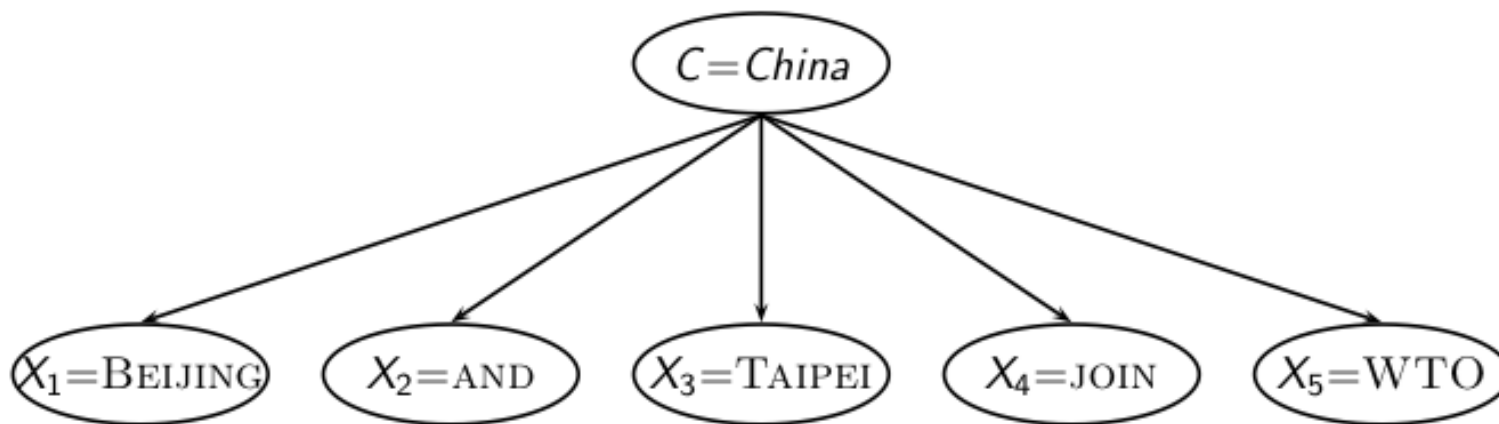
- Conditional probabilities:

$$\hat{P}(t|c) = \frac{T_{ct}}{\sum_{t' \in V} T_{ct'}}$$

- $T_{ct}$  is the number of tokens of  $t$  in training documents from class  $c$   
(includes multiple occurrences)

- We've made a Naive Bayes independence assumption  $\hat{P}(t_{k_1}|c) = \hat{P}(t_{k_2}|c)$  occurring in position  $k_1$  and  $k_2$ .

# The problem with maximum likelihood estimates: Zeros



- $P(\text{China} | d) \propto P(\text{China}) \cdot P(\text{BEIJING} | \text{China}) \cdot P(\text{AND} | \text{China})$   
•  $P(\text{TAIPEI} | \text{China}) \cdot P(\text{JOIN} | \text{China}) \cdot P(\text{WTO} | \text{China})$
- If WTO never occurs in class China in the train set:

$$\hat{P}(\text{WTO} | \text{China}) = \frac{T_{\text{China}, \text{WTO}}}{\sum_{t' \in V} T_{\text{China}, t'}} = \frac{0}{\sum_{t' \in V} T_{\text{China}, t'}} = 0$$

# The problem with maximum likelihood estimates: Zeros (cont)

- If there were no occurrences of WTO in documents in class China, we'd get a zero estimate:

$$\hat{P}(\text{WTO} | \text{China}) = \frac{T_{\text{China}, \text{WTO}}}{\sum_{t' \in V} T_{\text{China}, t'}} = 0$$

- We will get  $P(\text{China} | d) = 0$  for any document that contains WTO!
- Zero probabilities cannot be conditioned away.

# To avoid zeros: Add-one smoothing

- Before:

$$\hat{P}(t|c) = \frac{T_{ct}}{\sum_{t' \in V} T_{ct'}}$$

- Now: Add one to each count to avoid zeros:

$$\hat{P}(t|c) = \frac{T_{ct} + 1}{\sum_{t' \in V} (T_{ct'} + 1)} = \frac{T_{ct} + 1}{(\sum_{t' \in V} T_{ct'}) + B}$$

- B is the number of different words (in this case the size of the vocabulary:

$$|V| = M)$$

# To avoid zeros: Add-one smoothing

- Estimate parameters from the training corpus using add-one smoothing
- For a new document, for each class, compute sum of (i) log of prior and (ii) logs of conditional probabilities of the terms
- Assign the document to the class with the largest score

# Naive Bayes: Training

TRAINMULTINOMIALNB( $\mathbb{C}, \mathbb{D}$ )

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

# Naive Bayes: Testing

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

# Exercise

	docID	words in document	in $c = \textit{China}$ ?
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Chinese Chinese Chinese Tokyo Japan	?

- Estimate parameters of Naive Bayes classifier
- Classify test document



	docID	words in document	in c
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Chinese Chinese Chinese Tokyo Japan	?

## Example: Parameter estimates

Priors:  $\hat{P}(c) = 3/4$  and  $\hat{P}(\bar{c}) = 1/4$  Conditional probabilities:

$$\hat{P}(\text{CHINESE}|c) = (5 + 1)/(8 + 6) = 6/14 = 3/7$$

$$\hat{P}(\text{TOKYO}|c) = \hat{P}(\text{JAPAN}|c) = (0 + 1)/(8 + 6) = 1/14$$

$$\hat{P}(\text{CHINESE}|\bar{c}) = (1 + 1)/(3 + 6) = 2/9$$

$$\hat{P}(\text{TOKYO}|\bar{c}) = \hat{P}(\text{JAPAN}|\bar{c}) = (1 + 1)/(3 + 6) = 2/9$$

The denominators are  $(8 + 6)$  and  $(3 + 6)$  because the lengths of  $\text{text}_c$  and  $\text{text}_{\bar{c}}$  are 8 and 3, respectively, and because the constant  $B$  is 6 as the vocabulary consists of six terms.

# Example: Classification

	docID	words in document	in c :
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Chinese Chinese Chinese Tokyo Japan	?

$$\hat{P}(c|d_5) \propto 3/4 \cdot (3/7)^3 \cdot 1/14 \cdot 1/14 \approx 0.0003$$

$$\hat{P}(\bar{c}|d_5) \propto 1/4 \cdot (2/9)^3 \cdot 2/9 \cdot 2/9 \approx 0.0001$$

Thus, the classifier assigns the test document to  $c = \textit{China}$ . The reason for this classification decision is that the three occurrences of the positive indicator CHINESE in  $d_5$  outweigh the occurrences of the two negative indicators JAPAN and TOKYO.

# Time complexity of Naive Bayes

mode	time complexity
training	$\Theta( \mathbb{D} L_{\text{ave}} +  \mathbb{C}  V )$
testing	$\Theta(L_a +  \mathbb{C} M_a) = \Theta( \mathbb{C} M_a)$

- $L_{\text{ave}}$ : average length of a training doc,  $L_a$ : length of the test doc,  $M_a$ : number of distinct terms in the test doc,  $\mathbb{D}$ : training set,  $V$ : vocabulary, set of classes  $\mathbb{C}$ :
- $\Theta(|\mathbb{D}|L_{\text{ave}})$  is the time it takes to compute all counts.
- $\Theta(|\mathbb{C}||V|)$  is the time it takes to compute the parameters from the counts.
- Generally test time is also linear (in the length of the test document).
- Thus, Naive Bayes is linear in the size of the training set (training) and the test document  $|\mathbb{C}||V| < |\mathbb{D}|L_{\text{ave}}$  (testing). This is optimal.

# Violation of Naive Bayes independence assumption

- The independence assumptions do not really hold of documents written in natural language.
- Conditional independence:

$$P(\langle t_1, \dots, t_{n_d} \rangle | c) = \prod_{1 \leq k \leq n_d} P(X_k = t_k | c)$$

- Positional independence:  $\hat{P}(t_{k_1} | c) = \hat{P}(t_{k_2} | c)$
- Exercise
  - Examples for why conditional independence assumption is not really true?
  - Examples for why positional independence assumption is not really true?
- How can Naive Bayes work if it makes such inappropriate assumptions?

# Why does Naive Bayes work?

- Naive Bayes can work well even though conditional independence assumptions are badly violated.
- Example:

	$c_1$	$c_2$	class selected
true probability $P(c d)$	0.6	0.4	$c_1$
$\hat{P}(c) \prod_{1 \leq k \leq n_d} \hat{P}(t_k c)$	0.00099	0.00001	
NB estimate $\hat{P}(c d)$	0.99	0.01	$c_1$

- Double counting of evidence causes underestimation (0.01) and overestimation (0.99).
- Classification is about predicting the correct class and not about accurately estimating probabilities.
- Correct estimation  $\Rightarrow$  accurate prediction.
- But not vice versa!

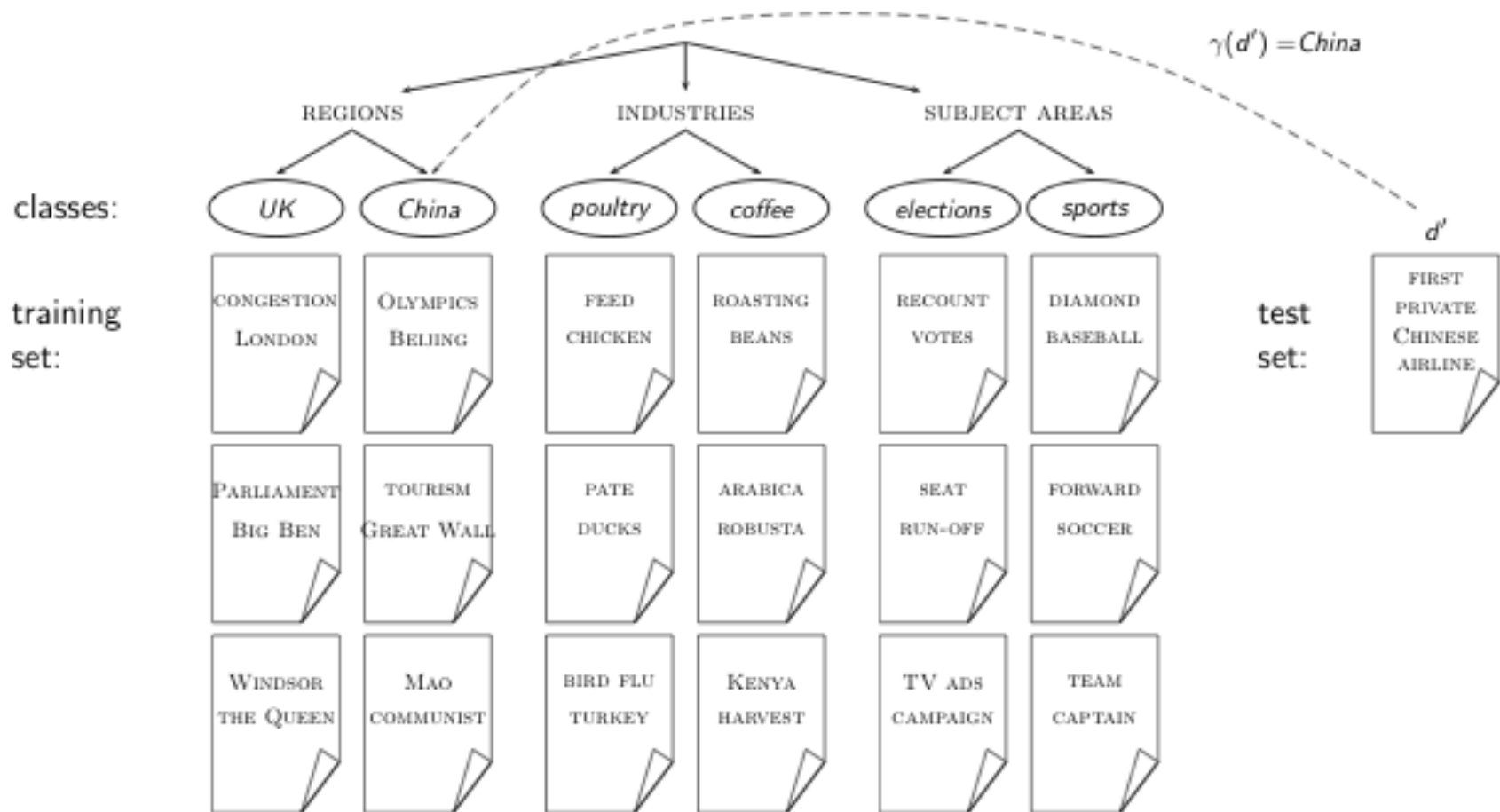
# Naive Bayes is not so naive

- Naive Naive Bayes has won some bakeoffs (e.g., KDD-CUP 97)
- More robust to non-relevant features than some more complex learning methods
- More robust to concept drift (changing of definition of class over time) than some more complex learning methods
- Better than methods like decision trees when we have many equally important features
- A good dependable baseline for text classification (but not the best)
- Optimal if independence assumptions hold (never true for text, but true for some domains)
- Very fast
- Low storage requirements

# Outline

- Text classification
- Naive Bayes
- **Evaluation of TC**
- Feature Selection

# Evaluation on Reuters





## Example: The Reuters collection

symbol	statistic	value
<i>N</i>	documents	800,000
<i>L</i>	avg. # word tokens per document	200
<i>M</i>	word types	400,000
	avg. # bytes per word token (incl. spaces/punct.)	6
	avg. # bytes per word token (without spaces/punct.)	4.5
	avg. # bytes per word type	7.5
	non-positional postings	100,000,000
type of class	number	examples
region	366	UK, China
industry	870	poultry, coffee
subject area	126	elections, sports

# A Reuters document



You are here: [Home](#) > [News](#) > [Science](#) > [Article](#)

Go to a Section: [U.S.](#) [International](#) [Business](#) [Markets](#) [Politics](#) [Entertainment](#) [Technology](#) [Sports](#) [Oddly Enough](#)

## Extreme conditions create rare Antarctic clouds

Tue Aug 1, 2006 3:20am ET

[Email This Article](#) [Print This Article](#) [Reprints](#)

[\[-\]](#) Text [\[+\]](#)



SYDNEY (Reuters) - Rare, mother-of-pearl colored clouds caused by extreme weather conditions above Antarctica are a possible indication of global warming, Australian scientists said on Tuesday.

Known as nacreous clouds, the spectacular formations showing delicate wisps of colors were photographed in the sky over an Australian

# Evaluating classification

- Evaluation must be done on test data that are independent of the training data (usually a disjoint set of instances).
- It's easy to get good performance on a test set that was available to the learner during training (e.g., just memorize the test set).
- Measures: Precision, recall,  $F_1$ , classification accuracy

# Precision $P$ and recall $R$

	in the class	not in the class
predicted to be in the class	true positives (TP)	false positives (FP)
predicted to not be in the class	false negatives (FN)	true negatives (TN)

$$P = TP / (TP + FP)$$

$$R = TP / (TP + FN)$$

## A combined measure: $F$

- $F_1$  allows us to trade off precision against recall.

$$F_1 = \frac{1}{\frac{1}{2} \frac{1}{P} + \frac{1}{2} \frac{1}{R}} = \frac{2PR}{P + R}$$

- This is the harmonic mean of  $P$  and  $R$ :  $\frac{1}{F} = \frac{1}{2} \left( \frac{1}{P} + \frac{1}{R} \right)$

# Averaging: Micro vs. Macro

- We now have an evaluation measure ( $F_1$ ) for one class.
- But we also want a single number that measures the aggregate performance over all classes in the collection.
- Macro-averaging
  - Compute  $F_1$  for each of the  $C$  classes
  - Average these  $C$  numbers
- Micro-averaging
  - Compute TP, FP, FN for each of the  $C$  classes
  - Sum these  $C$  numbers (e.g., all TP to get aggregate TP)
  - Compute  $F_1$  for aggregate TP, FP, FN

# Naive Bayes vs. other methods (Reuters- 21578)

(a)		NB	Rocchio	kNN		SVM
	micro-avg-L (90 classes)	80	85	86		89
	macro-avg (90 classes)	47	59	60		60

(b)		NB	Rocchio	kNN	trees	SVM
	earn	96	93	97	98	98
	acq	88	65	92	90	94
	money-fx	57	47	78	66	75
	grain	79	68	82	85	95
	crude	80	70	86	85	89
	trade	64	65	77	73	76
	interest	65	63	74	67	78
	ship	85	49	79	74	86
	wheat	70	69	77	93	92
	corn	65	48	78	92	90
	micro-avg (top 10)	82	65	82	88	92
	micro-avg-D (118 classes)	75	62	n/a	n/a	87

Evaluation measure:  $F_1$  Naive Bayes does pretty well, but some methods beat it consistently (e.g., SVM).

# Outline

- Text classification
- Naive Bayes
- Evaluation of TC
- **Feature Selection**



# Feature selection

- In text classification, we usually represent documents in a high-dimensional space, with each dimension corresponding to a term.
- Axis = dimension = word = term = feature
- Many dimensions correspond to rare words.
- Rare words can mislead the classifier.
- Rare misleading features are called **noise features**.
- Eliminating noise features from the representation increases efficiency and effectiveness of text classification.
- Eliminating features is called **feature selection**.

# Example for a noise feature

- Let's say we're doing text classification for the class *China*.
- Suppose a rare term, say ARACHNOCENTRIC, has no information about *China* . . .  
... but all instances of ARACHNOCENTRIC happen to occur in *China* documents in our training set.
- Then we may learn a classifier that incorrectly interprets ARACHNOCENTRIC as evidence for the class *China*.
- Such an incorrect generalization from an accidental property of the training set is called **overfitting**.
- **Feature selection reduces overfitting and improves the accuracy of the classifier.**

# Basic feature selection algorithm

```
SELECTFEATURES( $\mathbb{D}$ ,  $c$ ,  $k$ )
1   $V \leftarrow \text{EXTRACTVOCABULARY}(\mathbb{D})$ 
2   $L \leftarrow []$ 
3  for each  $t \in V$ 
4  do  $A(t, c) \leftarrow \text{COMPUTEFEATUREUTILITY}(\mathbb{D}, t, c)$ 
5      $\text{APPEND}(L, \langle A(t, c), t \rangle)$ 
6  return  $\text{FEATURESWITHLARGESTVALUES}(L, k)$ 
How do we compute  $A$ , the feature utility?
```

# Different feature selection methods

- A feature selection method is mainly defined by the feature utility measure it employs
- Feature utility measures:
  - Frequency – select the most frequent terms
  - Mutual information – select the terms with the highest mutual information
  - Mutual information is also called information gain in this context.
  - Chi-square

# Mutual information

- Compute the feature utility  $A(t, c)$  as the expected mutual information (MI) of term  $t$  and class  $c$ .
- MI tells us “how much information” the term contains about the class and vice versa.
  - For example, if a term’s occurrence is independent of the class (same proportion of docs within/without class contain the term), then MI is 0.
- Definition:

$$I(U; C) = \sum_{e_t \in \{1,0\}} \sum_{e_c \in \{1,0\}} P(U=e_t, C=e_c) \log_2 \frac{P(U=e_t, C=e_c)}{P(U=e_t)P(C=e_c)}$$

- Definition:

$$I(U; C) = \sum_{e_t \in \{1,0\}} \sum_{e_c \in \{1,0\}} P(U=e_t, C=e_c) \log_2 \frac{P(U=e_t, C=e_c)}{P(U=e_t)P(C=e_c)}$$

## How to compute MI values

- Based on maximum likelihood estimates, the formula we use is:

$$I(U; C) = \frac{N_{11}}{N} \log_2 \frac{NN_{11}}{N_{1.}N_{.1}} + \frac{N_{01}}{N} \log_2 \frac{NN_{01}}{N_{0.}N_{.1}} \\ + \frac{N_{10}}{N} \log_2 \frac{NN_{10}}{N_{1.}N_{.0}} + \frac{N_{00}}{N} \log_2 \frac{NN_{00}}{N_{0.}N_{.0}}$$

- $N_{10}$ : number of documents that contain  $t$  ( $et = 1$ ) and are not in  $c$  ( $ec = 0$ )
- $N_{11}$ : number of documents that contain  $t$  ( $et = 1$ ) and are in  $c$  ( $ec = 1$ )
- $N_{01}$ : number of documents that do not contain  $t$  ( $et = 0$ ) and are in  $c$  ( $ec = 1$ )
- $N_{00}$ : number of documents that do not contain  $t$  ( $et = 0$ ) and are not in  $c$  ( $ec = 0$ );  $N = N_{00} + N_{01} + N_{10} + N_{11}$ .

$$I(U; C) = \frac{N_{11}}{N} \log_2 \frac{NN_{11}}{N_{1.}N_{.1}} + \frac{N_{01}}{N} \log_2 \frac{NN_{01}}{N_{0.}N_{.1}} \\ + \frac{N_{10}}{N} \log_2 \frac{NN_{10}}{N_{1.}N_{.0}} + \frac{N_{00}}{N} \log_2 \frac{NN_{00}}{N_{0.}N_{.0}}$$

## MI example for *poultry*/EXPORT in Reuters

$$e_t = e_{\text{EXPORT}} = 1 \quad e_c = e_{\text{poultry}} = 1 \quad e_t = e_{\text{EXPORT}} = 0 \quad e_c = e_{\text{poultry}} = 0$$

$N_{11} = 49$	$N_{10} = 27,652$
$N_{01} = 141$	$N_{00} = 774,106$

Plug these values into formula:

$$I(U; C) = \frac{49}{801,948} \log_2 \frac{801,948 \cdot 49}{(49 + 27,652)(49 + 141)} \\ + \frac{141}{801,948} \log_2 \frac{801,948 \cdot 141}{(141 + 774,106)(49 + 141)} \\ + \frac{27,652}{801,948} \log_2 \frac{801,948 \cdot 27,652}{(49 + 27,652)(27,652 + 774,106)} \\ + \frac{774,106}{801,948} \log_2 \frac{801,948 \cdot 774,106}{(141 + 774,106)(27,652 + 774,106)} \\ \approx 0.000105$$

## MI feature selection on Reuters

Class: *coffee*

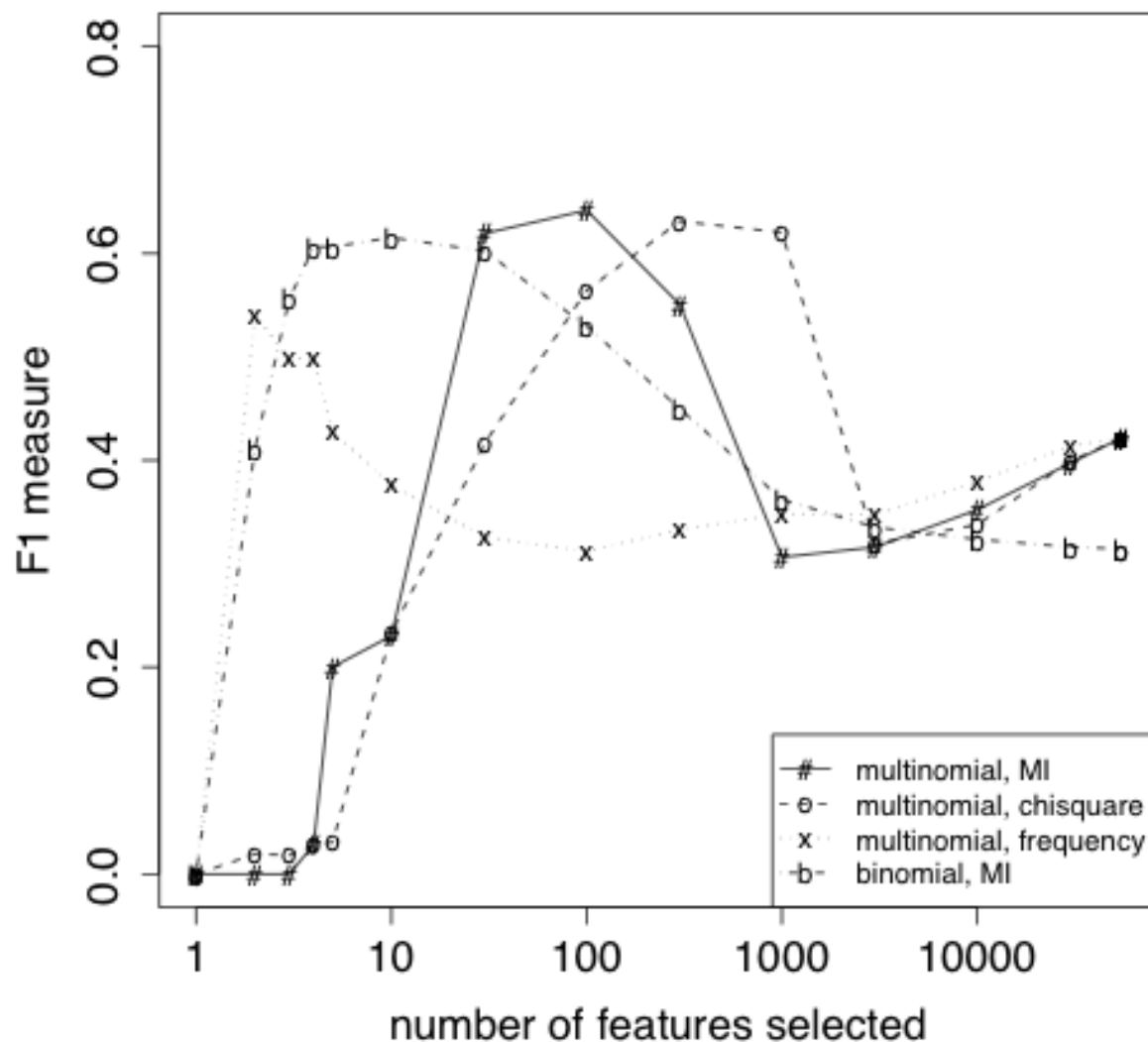
term	MI
COFFEE	0.0111
BAGS	0.0042
GROWERS	0.0025
KG	0.0019
COLOMBIA	0.0018
BRAZIL	0.0016
EXPORT	0.0014
EXPORTERS	0.0013
EXPORTS	0.0013
CROP	0.0012

Class: *sports*

term	MI
SOCCER	0.0681
CUP	0.0515
MATCH	0.0441
MATCHES	0.0408
PLAYED	0.0388
LEAGUE	0.0386
BEAT	0.0301
GAME	0.0299
GAMES	0.0284
TEAM	0.0264



# Naive Bayes: Effect of feature selection



- multinomial = multinomial Naive Bayes
- binomial = Bernoulli Naive Bayes

# Take-away today

- Text classification: definition & relevance to information retrieval
- Naive Bayes: simple baseline text classifier
- Evaluation of text classification: how do we know it worked / didn't work?
- Feature Selection improves performance.

---

---

# Thank you

*Questions?*

---

---