

# Text classification & Naive Bayes

Giorgio Gambosi

Course of Information Retrieval  
CdLM in Computer Science  
University of Rome Tor Vergata

Derived from slides produced by C. Manning and by H. Schütze

# A text classification task: Email spam filtering

From: "" <takworlld@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**  $\mathbb{X}$ 
  - Documents are represented in this space – typically some type of high-dimensional space.
- A fixed set of **classes**  $\mathbb{C} = \{c_1, c_2, \dots, c_J\}$ 
  - The classes are human-defined for the needs of an application (e.g., spam vs. nonspam).
- A **training set**  $\mathbb{D}$  of labeled documents. Each labeled document  $\langle d, c \rangle \in \mathbb{X} \times \mathbb{C}$

Using a learning method or **learning algorithm**, we then wish to learn a **classifier**  $\gamma$  that maps documents to classes:

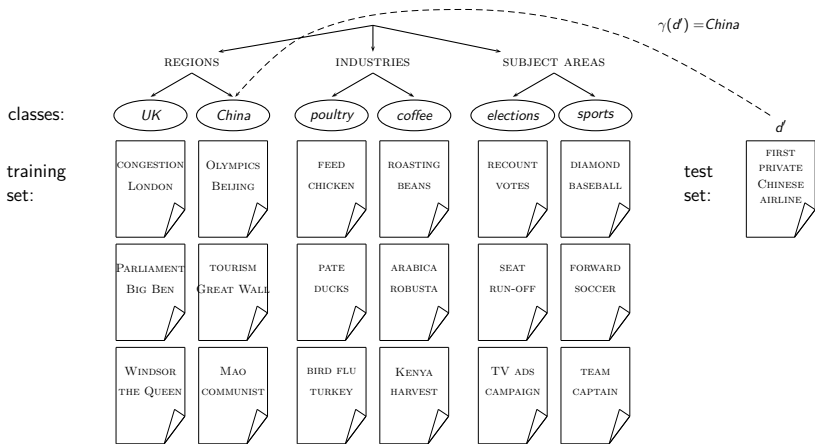
$$\gamma : \mathbb{X} \rightarrow \mathbb{C}$$

# Formal definition of TC: Application/Testing

Given: a description  $d \in \mathbb{X}$  of a document

Determine:  $\gamma(d) \in \mathbb{C}$ , that is, the class that is most appropriate for  $d$

# Topic classification



# Exercise

- Find examples of uses of text classification in information retrieval

# Examples of how search engines use classification

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

# 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

- E.g., Google Alerts is 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          * 0.70 film ACCRUE
                  ** 0.50 STEM
                  /wordtext = film
subtopic          ** 0.50 motion-picture PHRAS
                  *** 1.00 WORD
                  /wordtext = motion
                  *** 1.00 WORD
                  /wordtext = picture
                  ** 0.50 STEM
                  /wordtext = movie
subtopic          * 0.50 video ACCRUE
                  ** 0.50 STEM
                  /wordtext = video
                  ** 0.50 STEM
                  /wordtext = vcr
                  # End of art topic
subtopic          * 0.70 performing-arts ACCRUE
                  ** 0.50 WORD
                  /wordtext = ballet
subtopic          ** 0.50 WORD
                  /wordtext = dance
                  ** 0.50 WORD
                  /wordtext = opera
                  ** 0.30 WORD
                  /wordtext = symphony
subtopic          * 0.70 visual-arts ACCRUE
                  ** 0.50 WORD
                  /wordtext = painting
                  ** 0.50 WORD
                  /wordtext = sculpture

```

## Classification methods: 3. Statistical/Probabilistic

- This was our definition of the classification problem – text classification as a learning problem
- (i) Supervised learning of a the classification function  $\gamma$  and (ii) application of  $\gamma$  to classifying new documents
- We will look at two methods for doing this: Naive Bayes and SVMs
- No free lunch: requires hand-classified training data
- But this manual classification can be done by non-experts.

# 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:

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

- 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:

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

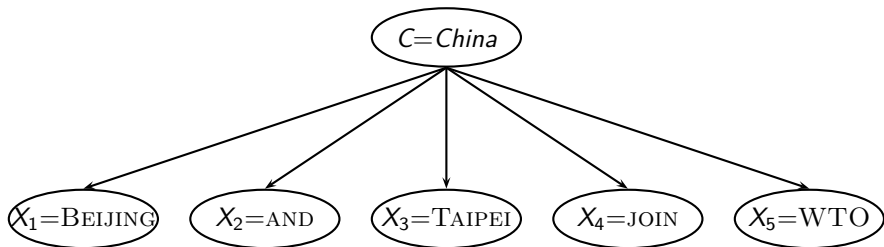
- $N_c$ : number of docs in class  $c$ ;  $N$ : total number of docs
- Conditional probabilities:

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

- $T_{ct}$  is the number of tokens of  $t$  in training documents from class  $c$  (includes multiple occurrences)
- We've made a **Naive Bayes independence assumption** here:  
 $\hat{P}(t_k|c) = \hat{P}(t_k|c)$ , independent of position



# The problem with maximum likelihood estimates: Zeros



$$P(\text{China}|d) \propto P(\text{China}) \cdot P(\text{BEIJING}|\text{China}) \cdot P(\text{AND}|\text{China}) \\ \cdot P(\text{TAIPEI}|\text{China}) \cdot P(\text{JOIN}|\text{China}) \cdot P(\text{WTO}|\text{China})$$

- 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 are no occurrences of *WTO* in documents in class *China*, we get a zero estimate:

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

- $\rightarrow$  We will get  $P(\text{China}|d) = 0$  for any document that contains *WTO*!

# 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 bins – in this case the number of different words or the size of the vocabulary  $|V| = M$

# Naive Bayes: Summary

- 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: Estimate parameters, classify test set

|              | 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 | ?                         |

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

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

( $B$  is the number of bins – in this case the number of different words or the size of the vocabulary  $|V| = M$ )

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





## 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

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

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

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 | $\Theta( \mathbb{D} L_{\text{ave}} +  \mathbb{C}  \mathcal{V} )$ |
| 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,  $\mathcal{V}$ : vocabulary,  $\mathbb{C}$ : set of classes
- $\Theta(|\mathbb{D}|L_{\text{ave}})$  is the time it takes to compute all counts.
- $\Theta(|\mathbb{C}||\mathcal{V}|)$  is the time it takes to compute the parameters from the counts.
- Generally:  $|\mathbb{C}||\mathcal{V}| < |\mathbb{D}|L_{\text{ave}}$
- 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 (testing). This is **optimal**.

# Naive Bayes: Analysis

- Now we want to gain a better understanding of the properties of Naive Bayes.
- We will formally derive the classification rule ...
- ...and make our assumptions explicit.

# Derivation of Naive Bayes rule

We want to find the class that is most likely given the document:

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} P(c|d)$$

Apply Bayes rule  $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$ :

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} \frac{P(d|c)P(c)}{P(d)}$$

Drop denominator since  $P(d)$  is the same for all classes:

$$c_{\text{map}} = \arg \max_{c \in \mathbb{C}} P(d|c)P(c)$$

# Too many parameters / sparseness

$$\begin{aligned}c_{\text{map}} &= \arg \max_{c \in \mathbb{C}} P(d|c)P(c) \\ &= \arg \max_{c \in \mathbb{C}} P(\langle t_1, \dots, t_k, \dots, t_{n_d} \rangle | c)P(c)\end{aligned}$$

- There are too many parameters  $P(\langle t_1, \dots, t_k, \dots, t_{n_d} \rangle | c)$ , one for each unique combination of a class and a sequence of words.
- We would need a very, very large number of training examples to estimate that many parameters.
- This is the problem of **data sparseness**.

# Naive Bayes conditional independence assumption

To reduce the number of parameters to a manageable size, we make the **Naive Bayes conditional independence assumption**:

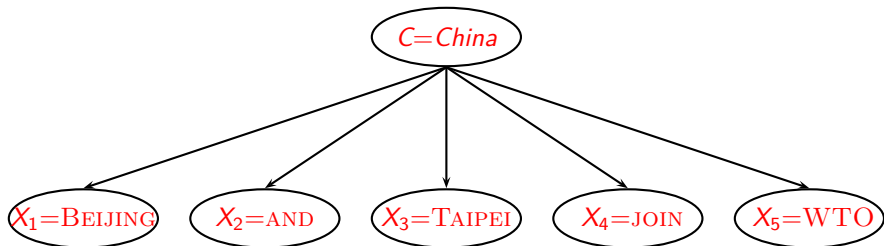
$$P(d|c) = P(\langle t_1, \dots, t_{n_d} \rangle | c) = \prod_{1 \leq k \leq n_d} P(X_k = t_k | c)$$

We assume that the probability of observing the conjunction of attributes is equal to the product of the individual probabilities  $P(X_k = t_k | c)$ .

Recall from earlier the estimates for these conditional probabilities:

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

# Generative model



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

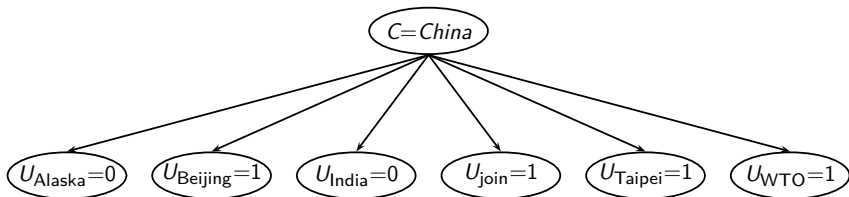
- Generate a class with probability  $P(c)$
- Generate each of the words (in their respective positions), conditional on the class, but independent of each other, with probability  $P(t_k|c)$
- To classify docs, we “reengineer” this process and find the class that is most likely to have generated the doc.



## Second independence assumption

- $\hat{P}(X_{k_1} = t|c) = \hat{P}(X_{k_2} = t|c)$
- For example, for a document in the class *UK*, the probability of generating *QUEEN* in the first position of the document is the same as generating it in the last position.
- The two independence assumptions amount to the **bag of words** model.

# A different Naive Bayes model: Bernoulli model



# Violation of Naive Bayes independence assumptions

- 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:
- $\hat{P}(X_{k_1} = t | c) = \hat{P}(X_{k_2} = t | c)$
- The independence assumptions do not really hold of documents written in natural language.
- 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:

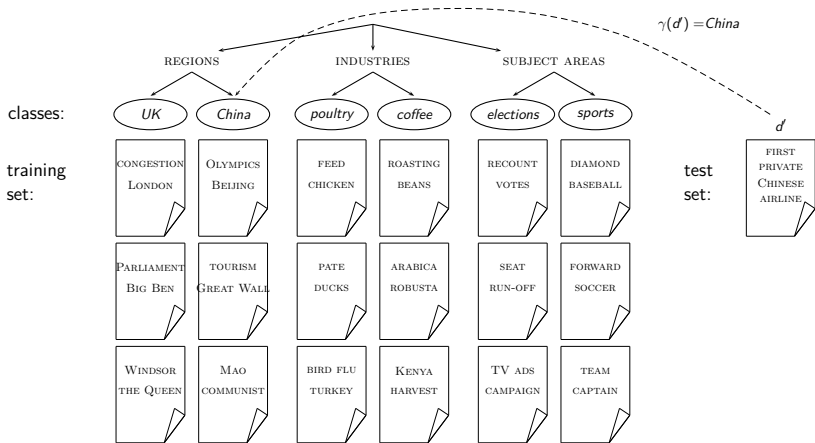
|   | $c_1$   | $c_2$   | class selected |
|---|---------|---------|----------------|
| true probability $P(c d)$                             | 0.6     | 0.4     | $c_1$          |
| $\hat{P}(c) \prod_{1 \leq k \leq n_d} \hat{P}(t_k c)$ | 0.00099 | 0.00001 |                |
| NB estimate $\hat{P}(c d)$                            | 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.
- Naive Bayes is terrible for correct estimation ...
- ...but it often performs well at accurate prediction (choosing the correct class).

# Naive Bayes is not so naive

- Naive Bayes has won some bakeoffs (e.g., KDD-CUP 97)
- More robust to nonrelevant 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

# Evaluation on Reuters



# Example: The Reuters collection

| symbol | statistic                       | value   |
|--------|---------------------------------|---------|
| $N$    | documents                       | 800,000 |
| $L$    | avg. # word tokens per document | 200     |
| $M$    | word types                      | 400,000 |

| type of class | number | examples          |
|---------------|--------|-------------------|
| region        | 366    | UK, China         |
| industry      | 870    | poultry, coffee   |
| subject area  | 126    | elections, sports |

# A Reuters document

**REUTERS** 

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

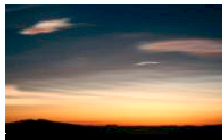
Go to a Section: [U.S.](#) [International](#) [Business](#) [Markets](#) [Politics](#) [Entertainment](#) [Technology](#) [Sports](#) [Oddly Enough](#)

## 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, i.e., training and test sets are disjoint.
- 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)  |

TP, FP, FN, TN are counts of documents. The sum of these four counts is the total number of documents.

$$\text{precision: } P = TP / (TP + FP)$$

$$\text{recall: } R = TP / (TP + FN)$$

## 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.
- **Macroaveraging**
  - Compute  $F_1$  for each of the  $C$  classes
  - Average these  $C$  numbers
- **Microaveraging**
  - 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

# $F_1$ scores for Naive Bayes vs. other methods

| (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  |

Naive Bayes does pretty well, but some methods beat it consistently (e.g., SVM).

# Feature selection

- In text classification, we usually represent documents in a **high-dimensional** space, with each dimension corresponding to a term.
- In this lecture: 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   APPEND( $L$ ,  $\langle A(t, c), t \rangle$ )  
6 return 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 (see book)

# Mutual information

- Compute the feature utility  $A(t, c)$  as the **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:

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

# How to compute MI values

- Based on maximum likelihood estimates, the formula we actually 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  $t$  ( $e_t = 1$ ) and are not in  $c$  ( $e_c = 0$ );  $N_{11}$ : number of documents that contain  $t$  ( $e_t = 1$ ) and are in  $c$  ( $e_c = 1$ );  $N_{01}$ : number of documents that do not contain  $t$  ( $e_t = 0$ ) and are in  $c$  ( $e_c = 1$ );  $N_{00}$ : number of documents that do not contain  $t$  ( $e_t = 0$ ) and are not in  $c$  ( $e_c = 0$ );  $N = N_{00} + N_{01} + N_{10} + N_{11}$ .

## How to compute MI values (2)

- Alternative way of computing MI:

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

- $N(U=e_t, C=e_c)$  is the count of documents with values  $e_t$  and  $e_c$ .
- $E(U=e_t, C=e_c)$  is the expected count of documents with values  $e_t$  and  $e_c$  if we assume that the two random variables are independent.

# MI example for *poultry*/EXPORT in Reuters

$$e_c = e_{poultry} = 1 \quad e_c = e_{poultry} = 0$$

|                        |                |                    |
|------------------------|----------------|--------------------|
| $e_t = e_{EXPORT} = 1$ | $N_{11} = 49$  | $N_{10} = 27,652$  |
| $e_t = e_{EXPORT} = 0$ | $N_{01} = 141$ | $N_{00} = 774,106$ |

Plug these values into formula:

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

# MI feature selection on Reuters

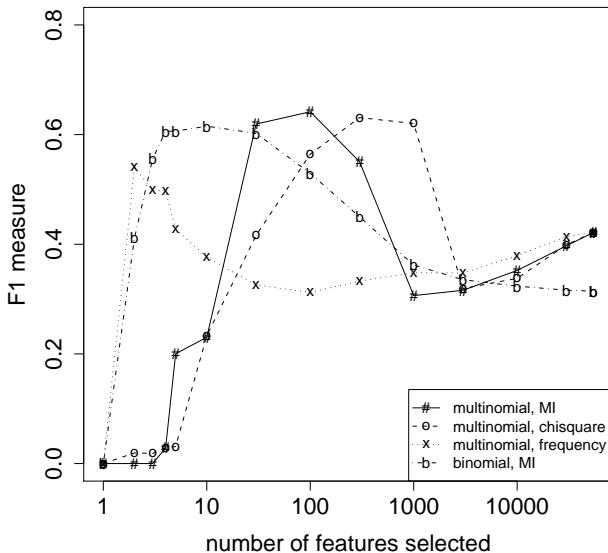
Class: *coffee*

| 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)

# Feature selection for Naive Bayes

- In general, feature selection is necessary for Naive Bayes to get decent performance.
- Also true for many other learning methods in text classification: **you need feature selection for optimal performance.**



# Exercise

(i) Compute the “export”/POULTRY contingency table for the “Kyoto”/JAPAN in the collection given below. (ii) Make up a contingency table for which MI is 0 – that is, term and class are independent of each other.

“export”/POULTRY table:

|                               |                                |                                |
|-------------------------------|--------------------------------|--------------------------------|
| $e_t = e_{\text{EXPORT}} = 1$ | $e_c = e_{\text{poultry}} = 1$ | $e_c = e_{\text{poultry}} = 0$ |
| $e_t = e_{\text{EXPORT}} = 0$ | $N_{11} = 49$                  | $N_{10} = 27,652$              |
|                               | $N_{01} = 141$                 | $N_{00} = 774,106$             |

Collection:

|              | docID | words in document     | in $c = \text{Japan?}$ |
|--------------|-------|-----------------------|------------------------|
| training set | 1     | Kyoto Osaka Taiwan    | yes                    |
|              | 2     | Japan Kyoto           | yes                    |
|              | 3     | Taipei Taiwan         | no                     |
|              | 4     | Macao Taiwan Shanghai | no                     |
|              | 5     | London                | no                     |

## Feature selection: MI for *poultry*/EXPORT

Goal of feature selection: eliminate noise and useless features for better effectiveness and efficiency

$$e_c = e_{poultry} = 1 \quad e_c = e_{poultry} = 0$$

|                        |                |                    |
|------------------------|----------------|--------------------|
| $e_t = e_{EXPORT} = 1$ | $N_{11} = 49$  | $N_{10} = 27,652$  |
| $e_t = e_{EXPORT} = 0$ | $N_{01} = 141$ | $N_{00} = 774,106$ |

Plug these values into formula:

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

# Feature selection for Reuters classes coffee and sports

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 |

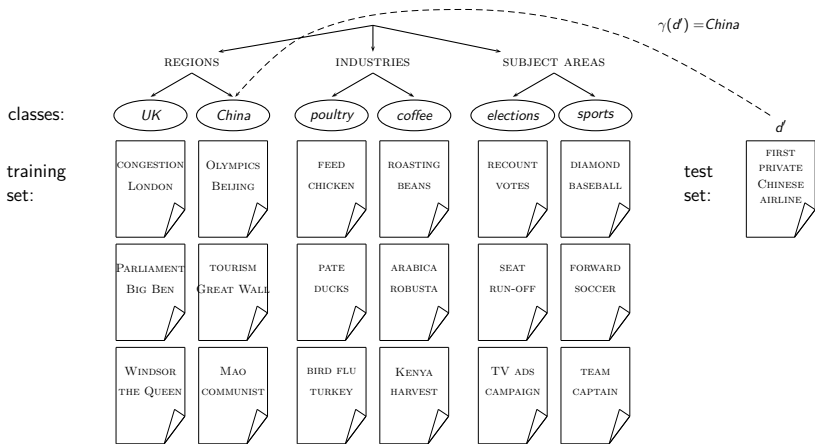
Class: *sports*

| 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 |

# Recall vector space representation

- Each document is a vector, one component for each term.
- Terms are axes.
- High dimensionality: 100,000s of dimensions
- Normalize vectors (documents) to unit length
- How can we do classification in this space?

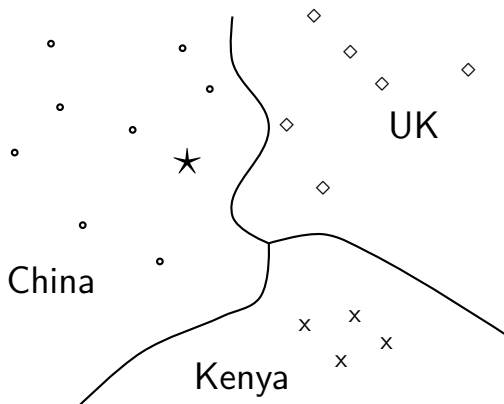
# Basic text classification setup



# Vector space classification

- As before, the training set is a set of documents, each labeled with its class.
- In vector space classification, this set corresponds to a labeled set of points or vectors in the vector space.
- Premise 1: Documents in the same class form a **contiguous region**.
- Premise 2: Documents from different classes **don't overlap**.
- We define lines, surfaces, hypersurfaces to divide regions.

# Classes in the vector space



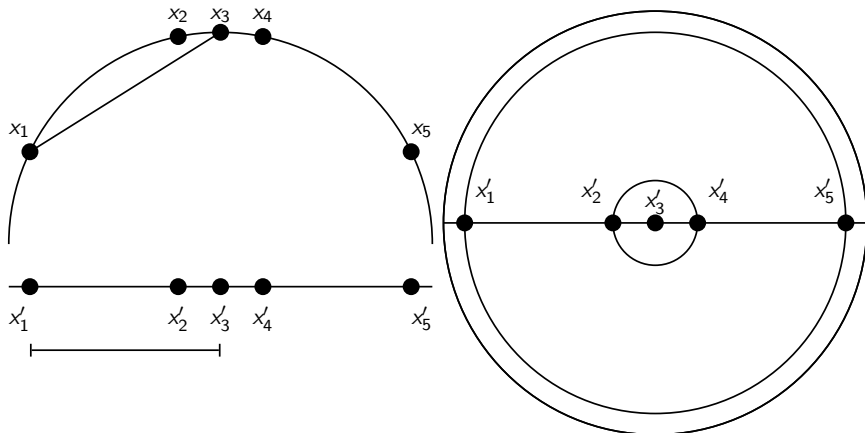
Should the document  $\star$  be assigned to *China*, *UK* or *Kenya*?

Find separators between the classes

Based on these separators:  $\star$  should be assigned to *China*

How do we find separators that do a good job at classifying new documents like  $\star$ ? – Main topic of today

## Aside: 2D/3D graphs can be misleading



*Left:* A projection of the 2D semicircle to 1D. For the points  $x_1, x_2, x_3, x_4, x_5$  at  $x$  coordinates  $-0.9, -0.2, 0, 0.2, 0.9$  the distance  $|x_2x_3| \approx 0.201$  only differs by 0.5% from  $|x'_2x'_3| = 0.2$ ; but  $|x_1x_3|/|x'_1x'_3| = d_{\text{true}}/d_{\text{projected}} \approx 1.06/0.9 \approx 1.18$  is an example of a large distortion (18%) when projecting a large area. *Right:* The corresponding projection of the 3D hemisphere to 2D.



# Relevance feedback

- In relevance feedback, the user marks documents as relevant/nonrelevant.
- Relevant/nonrelevant can be viewed as [classes](#) or [categories](#).
- For each document, the user decides which of these two classes is correct.
- The IR system then uses these class assignments to build a better query (“model”) of the information need ...
- ...and returns better documents.
- Relevance feedback is a form of [text classification](#).

# Using Rocchio for vector space classification

- The principal difference between relevance feedback and text classification:
  - The training set is given as part of the input in text classification.
  - It is interactively created in relevance feedback.

# Rocchio classification: Basic idea

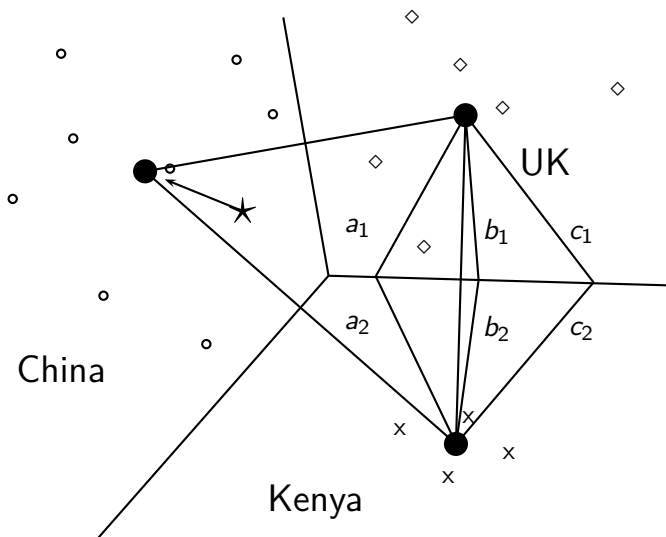
- Compute a centroid for each class
  - The centroid is the average of all documents in the class.
- Assign each test document to the class of its closest centroid.

# Recall definition of centroid

$$\vec{\mu}(c) = \frac{1}{|D_c|} \sum_{d \in D_c} \vec{v}(d)$$

where  $D_c$  is the set of all documents that belong to class  $c$  and  $\vec{v}(d)$  is the vector space representation of  $d$ .

# Rocchio illustrated: $a_1 = a_2, b_1 = b_2, c_1 = c_2$



# Rocchio algorithm

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

- 1 **for each**  $c_j \in \mathbb{C}$
- 2 **do**  $D_j \leftarrow \{d : \langle d, c_j \rangle \in \mathbb{D}\}$
- 3  $\vec{\mu}_j \leftarrow \frac{1}{|D_j|} \sum_{d \in D_j} \vec{v}(d)$
- 4 **return**  $\{\vec{\mu}_1, \dots, \vec{\mu}_J\}$

APPLYROCCHIO( $\{\vec{\mu}_1, \dots, \vec{\mu}_J\}, d$ )

- 1 **return**  $\arg \min_j |\vec{\mu}_j - \vec{v}(d)|$

# Rocchio properties

- Rocchio forms a simple representation for each class: the **centroid**
  - We can interpret the centroid as the **prototype** of the class.
- Classification is based on similarity to / distance from centroid/prototype.
- Does not guarantee that classifications are consistent with the training data!

# Time complexity of Rocchio

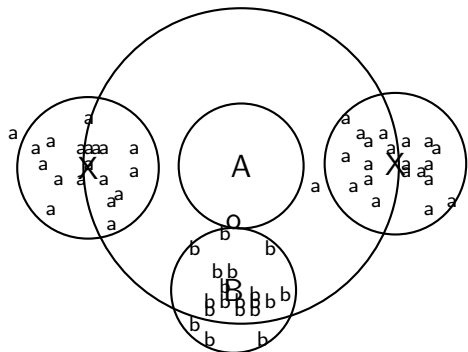
| mode     | time complexity   |
|----------|---|
| training | $\Theta( \mathbb{D} L_{\text{ave}} +  \mathbb{C}  V ) \approx \Theta( \mathbb{D} L_{\text{ave}})$ |
| testing  | $\Theta(L_a +  \mathbb{C} M_a) \approx \Theta( \mathbb{C} M_a)$                                   |



# Rocchio vs. Naive Bayes

- In many cases, Rocchio performs worse than Naive Bayes.
- One reason: Rocchio does not handle nonconvex, multimodal classes correctly.

# Rocchio cannot handle nonconvex, multimodal classes



Exercise: Why is Rocchio not expected to do well for the classification task a vs. b here?

- A is centroid of the a's, B is centroid of the b's.
- The point o is closer to A than to B.
- But o is a better fit for the b class.
- A is a multimodal class with two prototypes.
- But in Rocchio we only have one prototype

# kNN classification

- kNN classification is another vector space classification method.
- It also is very simple and easy to implement.
- kNN is more accurate (in most cases) than Naive Bayes and Rocchio.
- If you need to get a pretty accurate classifier up and running in a short time ...
- ...and you don't care about efficiency that much ...
- ...use kNN.

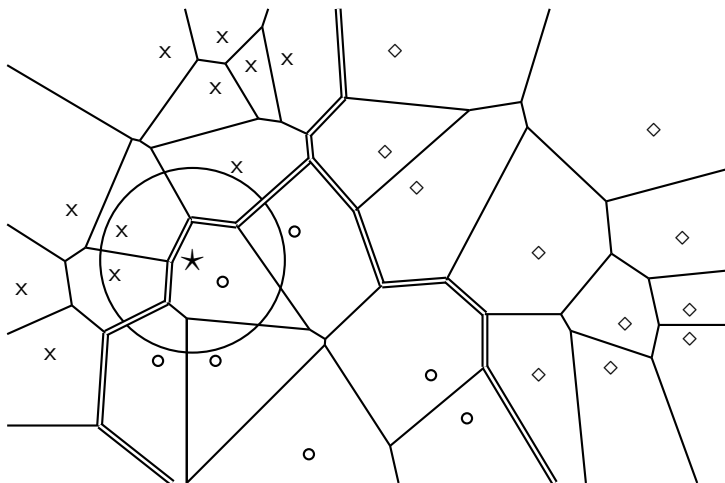
# kNN classification

- kNN =  $k$  nearest neighbors
- **kNN classification rule for  $k = 1$  (1NN)**: Assign each test document to the class of its **nearest neighbor** in the training set.
- 1NN is not very robust – one document can be mislabeled or atypical.
- **kNN classification rule for  $k > 1$  (kNN)**: Assign each test document to the **majority class of its  $k$  nearest neighbors** in the training set.
- Rationale of kNN: contiguity hypothesis
  - We expect a test document  $d$  to have the same label as the training documents located in the local region surrounding  $d$ .

# Probabilistic kNN

- Probabilistic version of kNN:  $P(c|d)$  = fraction of  $k$  neighbors of  $d$  that are in  $c$
- **kNN classification rule for probabilistic kNN:** Assign  $d$  to class  $c$  with highest  $P(c|d)$

# kNN is based on Voronoi tessellation



# kNN algorithm

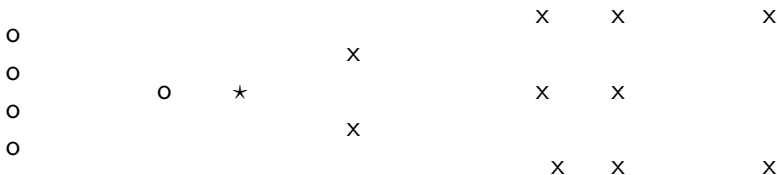
TRAIN-KNN( $\mathbb{C}, \mathbb{D}$ )

- 1  $\mathbb{D}' \leftarrow \text{PREPROCESS}(\mathbb{D})$
- 2  $k \leftarrow \text{SELECT-K}(\mathbb{C}, \mathbb{D}')$
- 3 **return**  $\mathbb{D}', k$

APPLY-KNN( $\mathbb{D}', k, d$ )

- 1  $S_k \leftarrow \text{COMPUTENEARESTNEIGHBORS}(\mathbb{D}', k, d)$
- 2 **for each**  $c_j \in \mathbb{C}(\mathbb{D}')$
- 3 **do**  $p_j \leftarrow |S_k \cap c_j|/k$
- 4 **return**  $\arg \max_j p_j$

# Exercise



How is star classified by:

(i) 1-NN (ii) 3-NN (iii) 9-NN (iv) 15-NN (v) Rocchio?



# Time complexity of kNN

## kNN with preprocessing of training set

training  $\Theta(|\mathbb{D}|L_{ave})$

testing  $\Theta(L_a + |\mathbb{D}|M_{ave}M_a) = \Theta(|\mathbb{D}|M_{ave}M_a)$

- kNN test time proportional to the size of the training set!
- The larger the training set, the longer it takes to classify a test document.
- kNN is inefficient for very large training sets.
- Question: Can we divide up the training set into regions, so that we only have to search in one region to do kNN classification for a given test document? (which perhaps would give us better than linear time complexity)

# Curse of dimensionality

- Our intuitions about space are based on the 3D world we live in.
- Intuition 1: some things are close by, some things are distant.
- Intuition 2: we can carve up space into areas such that: within an area things are close, distances between areas are large.
- These two intuitions don't necessarily hold for high dimensions.
- In particular: for a set of  $k$  uniformly distributed points, let  $d_{\min}$  be the smallest distance between any two points and  $d_{\max}$  be the largest distance between any two points.
- Then

$$\lim_{d \rightarrow \infty} \frac{d_{\max} - d_{\min}}{d_{\min}} = 0$$

# Curse of dimensionality: Simulation

- Simulate

$$\lim_{d \rightarrow \infty} \frac{d_{\max} - d_{\min}}{d_{\min}} = 0$$

- Pick a dimensionality  $d$
- Generate 10 random points in the  $d$ -dimensional hypercube (uniform distribution)
- Compute all 45 distances
- Compute  $\frac{d_{\max} - d_{\min}}{d_{\min}}$
- We see that intuition 1 (some things are close, others are distant) is not true for high dimensions.

## Intuition 2: Space can be carved up

- Intuition 2: we can carve up space into areas such that: within an area things are close, distances between areas are large.
- If this is true, then we have a simple and efficient algorithm for kNN.
- To find the  $k$  closest neighbors of data point  $\langle x_1, x_2, \dots, x_d \rangle$  do the following.
- Using binary search find all data points whose first dimension is in  $[x_1 - \epsilon, x_1 + \epsilon]$ . This is  $O(\log n)$  where  $n$  is the number of data points.
- Do this for each dimension, then intersect the  $d$  subsets.

## Intuition 2: Space can be carved up

- Size of data set  $n = 100$
- Again, assume uniform distribution in hypercube
- Set  $\epsilon = 0.05$ : we will look in an interval of length 0.1 for neighbors on each dimension.
- What is the probability that the nearest neighbor of a new data point  $\vec{x}$  is in this neighborhood in  $d = 1$  dimension?
- for  $d = 1$ :  $1 - (1 - 0.1)^{100} \approx 0.99997$
- In  $d = 2$  dimensions?
- for  $d = 2$ :  $1 - (1 - 0.1^2)^{100} \approx 0.63$
- In  $d = 3$  dimensions?
- for  $d = 3$ :  $1 - (1 - 0.1^3)^{100} \approx 0.095$
- In  $d = 4$  dimensions?
- for  $d = 4$ :  $1 - (1 - 0.1^4)^{100} \approx 0.0095$
- In  $d = 5$  dimensions?
- for  $d = 5$ :  $1 - (1 - 0.1^5)^{100} \approx 0.0009995$

## Intuition 2: Space can be carved up

- In  $d = 5$  dimensions?
- for  $d = 5$ :  $1 - (1 - 0.1^5)^{100} \approx 0.0009995$
- In other words: with enough dimensions, there is only one “local” region that will contain the nearest neighbor with high certainty: the entire search space.
- We cannot carve up high-dimensional space into neat neighborhoods ...
- ...unless the “true” dimensionality is much lower than  $d$ .

# kNN: Discussion

- No training necessary
  - But linear preprocessing of documents is as expensive as training Naive Bayes.
  - We always preprocess the training set, so in reality training time of kNN is linear.
- kNN is very accurate if training set is large.
- Optimality result: asymptotically zero error if Bayes rate is zero.
- But kNN can be very inaccurate if training set is small.

# Linear classifiers

- Definition:
  - A linear classifier computes a linear combination or weighted sum  $\sum_i w_i x_i$  of the feature values.
  - Classification decision:  $\sum_i w_i x_i > \theta$ ?
  - ...where  $\theta$  (the threshold) is a parameter.
- (First, we only consider binary classifiers.)
- Geometrically, this corresponds to a line (2D), a plane (3D) or a hyperplane (higher dimensionalities), the **separator**.
- We find this separator based on training set.
- Methods for finding separator: Perceptron, Rocchio, Naive Bayes – as we will explain on the next slides
- Assumption: The classes are **linearly separable**.

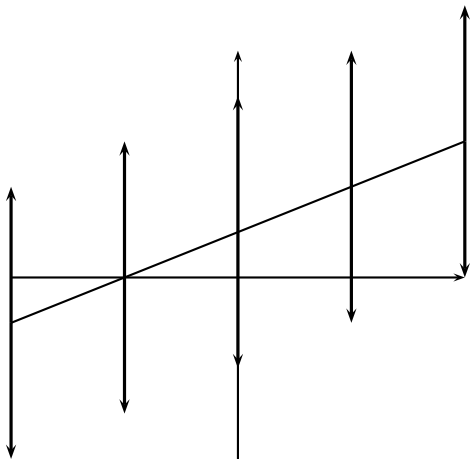


# A linear classifier in 1D



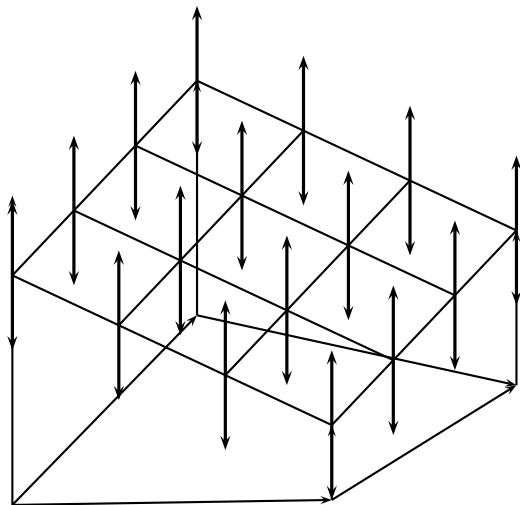
- A linear classifier in 1D is a point described by the equation  $w_1 d_1 = \theta$
- The point at  $\theta/w_1$
- Points  $(d_1)$  with  $w_1 d_1 \geq \theta$  are in the class  $c$ .
- Points  $(d_1)$  with  $w_1 d_1 < \theta$  are in the complement class  $\bar{c}$ .

# A linear classifier in 2D



- A linear classifier in 2D is a line described by the equation  $w_1 d_1 + w_2 d_2 = \theta$
- Example for a 2D linear classifier
- Points  $(d_1 \ d_2)$  with  $w_1 d_1 + w_2 d_2 \geq \theta$  are in the class  $c$ .
- Points  $(d_1 \ d_2)$  with  $w_1 d_1 + w_2 d_2 < \theta$  are in the complement class  $\bar{c}$ .

# A linear classifier in 3D



- A linear classifier in 3D is a plane described by the equation
$$w_1 d_1 + w_2 d_2 + w_3 d_3 = \theta$$
- Example for a 3D linear classifier
- Points  $(d_1 \ d_2 \ d_3)$  with  $w_1 d_1 + w_2 d_2 + w_3 d_3 \geq \theta$  are in the class  $c$ .
- Points  $(d_1 \ d_2 \ d_3)$  with  $w_1 d_1 + w_2 d_2 + w_3 d_3 < \theta$  are in the complement class  $\bar{c}$ .

# Rocchio as a linear classifier

- Rocchio is a linear classifier defined by:

$$\sum_{i=1}^M w_i d_i = \vec{w} \vec{d} = \theta$$

where  $\vec{w}$  is the **normal vector**  $\vec{\mu}(c_1) - \vec{\mu}(c_2)$  and  $\theta = 0.5 * (|\vec{\mu}(c_1)|^2 - |\vec{\mu}(c_2)|^2)$ .

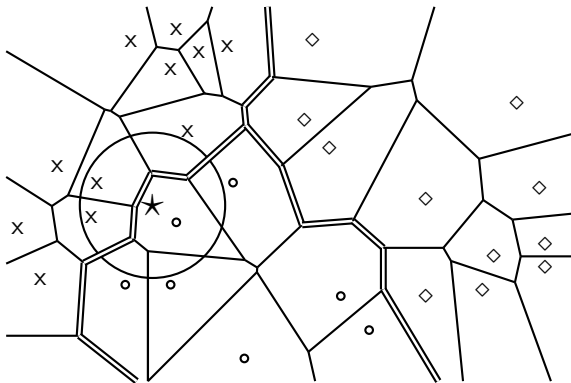
# Naive Bayes as a linear classifier

Multinomial Naive Bayes is a linear classifier (in log space) defined by:

$$\sum_{i=1}^M w_i d_i = \theta$$

where  $w_i = \log[\hat{P}(t_i|c)/\hat{P}(t_i|\bar{c})]$ ,  $d_i =$  number of occurrences of  $t_i$  in  $d$ , and  $\theta = -\log[\hat{P}(c)/\hat{P}(\bar{c})]$ . Here, the index  $i$ ,  $1 \leq i \leq M$ , refers to terms of the vocabulary (not to positions in  $d$  as  $k$  did in our original definition of Naive Bayes)

# kNN is not a linear classifier



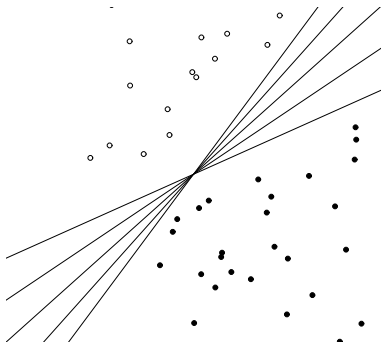
- Classification decision based on majority of  $k$  nearest neighbors.
- The decision boundaries between classes are piecewise linear ...
- ...but they are in general not linear classifiers that can be described as
$$\sum_{i=1}^M w_i d_i = \theta.$$

# Example of a linear two-class classifier

| $t_i$      | $w_i$ | $d_{1i}$ | $d_{2i}$ | $t_i$ | $w_i$ | $d_{1i}$ | $d_{2i}$ |
|------------|-------|----------|----------|-------|-------|----------|----------|
| prime      | 0.70  | 0        | 1        | dlrs  | -0.71 | 1        | 1        |
| rate       | 0.67  | 1        | 0        | world | -0.35 | 1        | 0        |
| interest   | 0.63  | 0        | 0        | sees  | -0.33 | 0        | 0        |
| rates      | 0.60  | 0        | 0        | year  | -0.25 | 0        | 0        |
| discount   | 0.46  | 1        | 0        | group | -0.24 | 0        | 0        |
| bundesbank | 0.43  | 0        | 0        | dlr   | -0.24 | 0        | 0        |

- This is for the class *interest* in Reuters-21578.
- For simplicity: assume a simple 0/1 vector representation
- $d_1$ : “rate discount dlrs world”
- $d_2$ : “prime dlrs”
- $\theta = 0$
- Exercise: Which class is  $d_1$  assigned to? Which class is  $d_2$  assigned to?
- We assign document  $\vec{d}_1$  “rate discount dlrs world” to *interest* since  $\vec{w}^T \vec{d}_1 = 0.67 \cdot 1 + 0.46 \cdot 1 + (-0.71) \cdot 1 + (-0.35) \cdot 1 = 0.07 > 0 = \theta$ .
- We assign  $\vec{d}_2$  “prime dlrs” to the complement class (not in *interest*) since  $\vec{w}^T \vec{d}_2 = -0.01 \leq \theta$ .

# Which hyperplane?





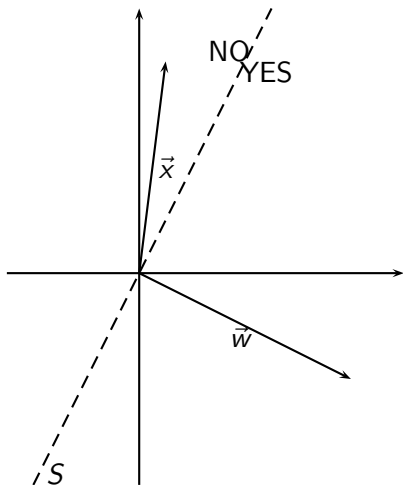
# Learning algorithms for vector space classification

- In terms of actual computation, there are two types of learning algorithms.
- (i) **Simple** learning algorithms that estimate the parameters of the classifier directly from the training data, often **in one linear pass**.
  - Naive Bayes, Rocchio, kNN are all examples of this.
- (ii) **Iterative** algorithms
  - Support vector machines
  - Perceptron (example available as PDF on website: <http://cislmu.org>)
- **The best performing learning algorithms usually require iterative learning.**

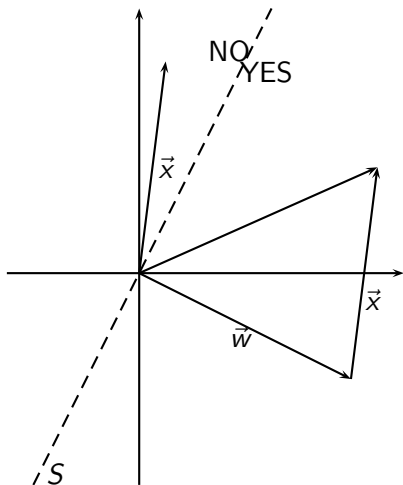
# Perceptron update rule

- Randomly initialize linear separator  $\vec{w}$
- Do until convergence:
  - Pick data point  $\vec{x}$
  - If  $\text{sign}(\vec{w}^T \vec{x})$  is correct class (1 or -1): do nothing
  - Otherwise:  $\vec{w} = \vec{w} - \text{sign}(\vec{w}^T \vec{x}) \vec{x}$

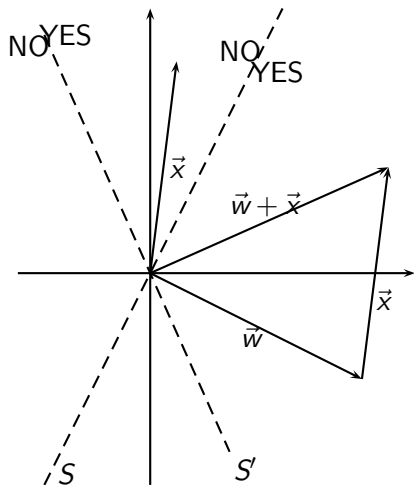
# Perceptron (class of $\vec{x}$ is YES)



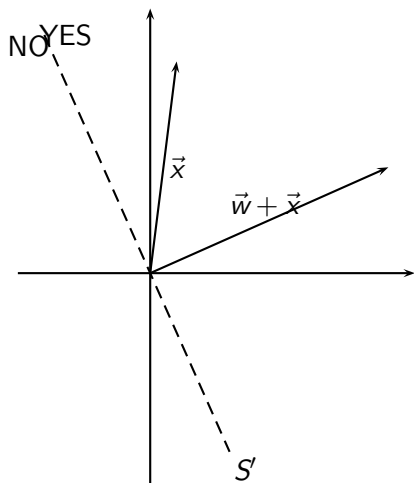
# Perceptron (class of $\vec{x}$ is YES)



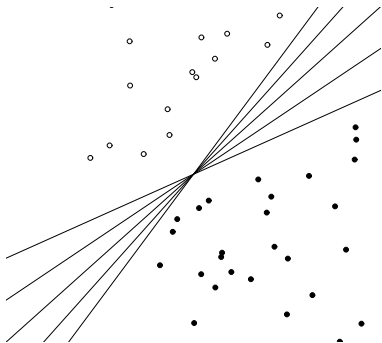
# Perceptron (class of $\vec{x}$ is YES)



# Perceptron (class of $\vec{x}$ is YES)



# Which hyperplane?



# Which hyperplane?

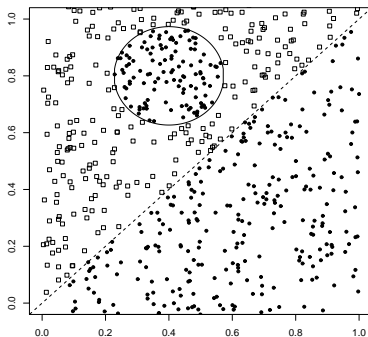
- For linearly separable training sets: there are **infinitely** many separating hyperplanes.
- They all separate the training set perfectly ...
- ...but they behave differently on test data.
- Error rates on new data are low for some, high for others.
- How do we find a low-error separator?
- Perceptron: generally bad; Naive Bayes, Rocchio: ok; linear SVM: good



# Linear classifiers: Discussion

- Many common text classifiers are linear classifiers: Naive Bayes, Rocchio, logistic regression, linear support vector machines etc.
- Each method has a different way of selecting the separating hyperplane
  - Huge differences in performance on test documents
- Can we get better performance with more powerful nonlinear classifiers?
- Not in general: A given amount of training data may suffice for estimating a linear boundary, but not for estimating a more complex nonlinear boundary.

# A nonlinear problem

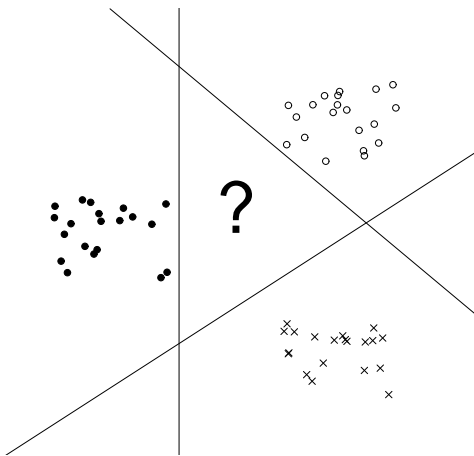


- Linear classifier like Rocchio does badly on this task.
- kNN will do well (assuming enough training data)

# Which classifier do I use for a given TC problem?

- Is there a learning method that is optimal for all text classification problems?
- No, because there is a tradeoff between bias and variance.
- Factors to take into account:
  - How much training data is available?
  - How simple/complex is the problem? (linear vs. nonlinear decision boundary)
  - How noisy is the problem?
  - How stable is the problem over time?
    - For an unstable problem, it's better to use a simple and robust classifier.

# How to combine hyperplanes for $> 2$ classes?



# One-of problems

- One-of or multiclass classification
  - Classes are mutually exclusive.
  - Each document belongs to exactly one class.
  - Example: language of a document (assumption: no document contains multiple languages)

# One-of classification with linear classifiers

- Combine two-class linear classifiers as follows for one-of classification:
  - Run each classifier separately
  - Rank classifiers (e.g., according to score)
  - Pick the class with the highest score

# Any-of problems

- Any-of or multilabel classification
  - A document can be a member of 0, 1, or many classes.
  - A decision on one class leaves decisions open on all other classes.
  - A type of “independence” (but not statistical independence)
  - Example: topic classification
  - Usually: make decisions on the region, on the subject area, on the industry and so on “independently”

# Any-of classification with linear classifiers

- Combine two-class linear classifiers as follows for any-of classification:
  - Simply run each two-class classifier separately on the test document and assign document accordingly