

# Text Classification

Sargur Srihari

# Information Retrieval – Standing Query

- Ad-hoc retrieval
  - Transient information need
    - Pose one or more queries to a search engine
  - Example
    - Need to track developments in *Multicore Computer Chips*
    - Issue the query against index of recent newswire articles
      - multicore AND computer AND chip
  - Standing Query
    - Same as every other query but executed periodically on collection which is incremented periodically
    - Will miss new articles such as *multicore processors*
    - Queries get complex
      - (multicore OR multi-core) AND (chip OR processor OR microprocessor)
- Classification
  - Captures generality and scope of problem space of standing queries

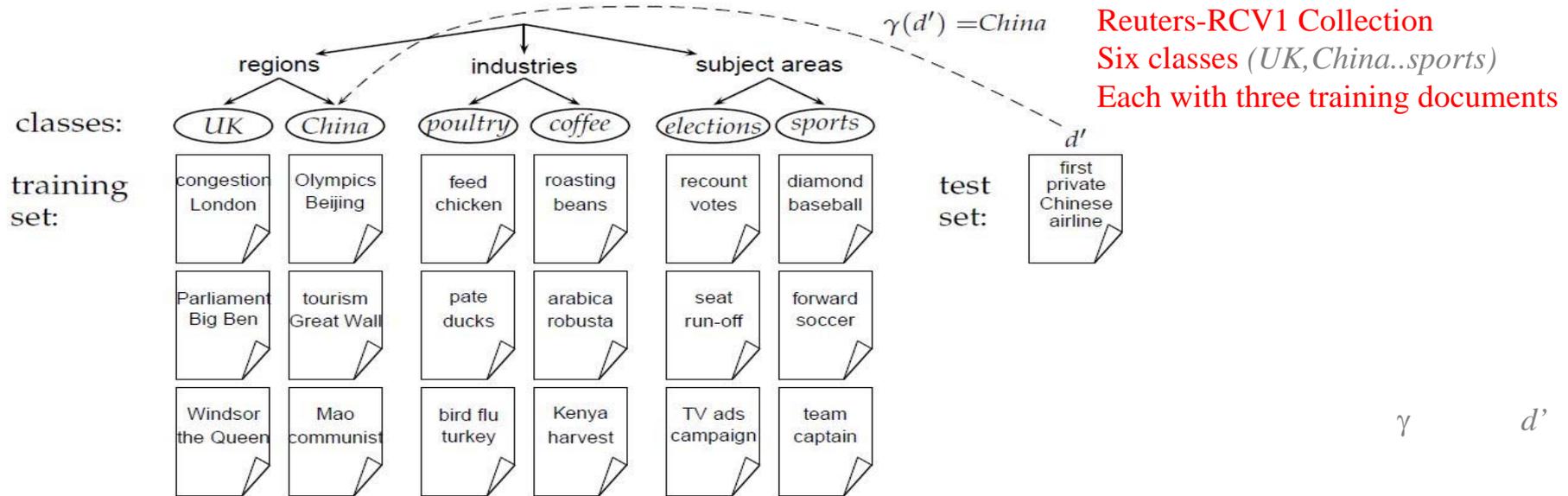
# Classification

- Given a set of classes
- Determine which class a given object belongs to
- Standing query divides new newswire articles into into two classes (two-class classification)
  - Documents about **multicore computer chips**
  - Documents not about **multicore computer chips**
- Classification with standing queries is also called **Routing** or **Filtering**

# Text Classification Synonyms

- Class need not be as focused as **multicore computer chips**
- Can be a more generic subject area
  - Like **coffee** or **China**
- Referred to as Topics
- Classification task is then referred to as
  - Text Classification
  - Text Categorization
  - Topic Classification
  - Topic Spotting

# Example of "China" as a class:



► Figure 13.1 Classes, training set and test set in text classification.

# Classification Examples in IR

- Pre-processing for indexing
  - Detecting a document's encoding
    - ASCII, Unicode, UTF-8
  - Word segmentation
    - Is gap between two letters a word boundary or not?
  - Truecasing
  - Identifying the language
- Automatic detection of
  - spam pages
  - Automatic detection of sexually explicit content
    - Include only if SafeSearch is off

# More Classification Examples in IR

- **Sentiment Detection**
  - Automatic detection of movie or product review as positive or negative
    - User checks for negative reviews before buying a camera
- **Email Sorting**
  - User has folders such as
    - Talk announcements, bills, family and friends, spam
- **Vertical Search Engine**
  - Restrict searches to a particular topic
    - Query computer science for China will return computer science departments in China with higher precision and recall than the query computer science China on a general database
    - Because vertical engine will not have china in a different sense (hard white ceramic)

# Text Classification Methods

- Manual

- Computer is not essential
- Traditionally solved manually
  - Books in a library are assigned Library of Congress categories by a librarian

- Rules

- Standing queries can be thought of as rules
  - (multicore OR multi-core) AND (chip OR processor OR multiprocessor)
- Have good scaling properties but labor intensive to create/maintain
- A domain expert can do better than automatically generated classifiers, but difficult to find someone with specialized skills

# Statistical Text Classification

- Third Approach is machine learning
- Decision criterion is automatically learned from training data
- Require a number of good example documents for each class
  - Called training documents
- Need for manual labeling is not eliminated
  - Labeling is easier than writing rules
  - Anyone can look at a document and say whether or not it is related to China
  - Sometimes part of work-flow
    - News articles returned by a standing query are classified each day as relevant or not

# Text Classification Example

- Three predefined categories: **religion, politics** and **computers**
- Input document ( from 20\_newsgroups ):

- Newsgroups: comp.os.ms-windows.misc

Subject: WIn 3.0 ICON HELP PLEASE!

Message-ID: <1993Apr27.193102.12730@iscsvax.uni.edu>

Date: 27 Apr 93 19:31:02 -0600

Organization: University of Northern Iowa

I have win 3.0 and downloaded several icons and BMP's but I cant figure out how to change the "wallpaper" or use the icons. Any help would be appreciated.

# Text Categorization as Statistical Learning

- Classify new documents based on likelihood suggested by training set

- **Input:**

- a new text document

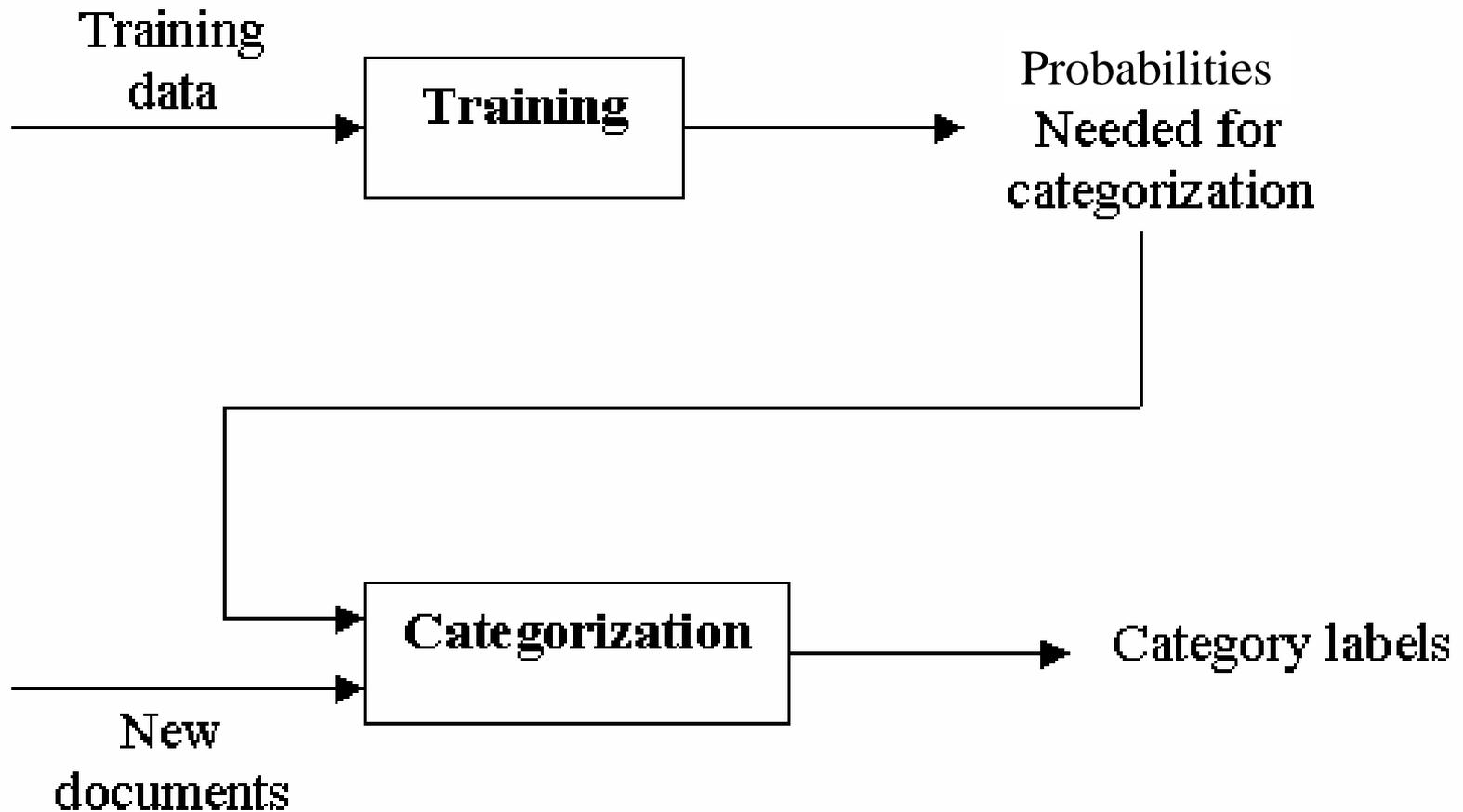
- **Output:**

- score(confidence) corresponding to each category. Category with highest score(confidence) is chosen

Categories:                      Scores:

- |            |    |
|------------|----|
| • Religion | 30 |
| • Politics | 27 |
| • Computer | 98 |

# Machine Learning Approach to Text Categorization



# Two main procedures

- *Training ( Learning )*: Calculate probabilities based on the training data
  - Input: training documents
  - Output: probabilities needed for categorization
- *Categorization ( Classification )*: Assign category labels to new documents based on likelihood
  - Input: a new document
  - Output: the class labels assigned to the document by the system.

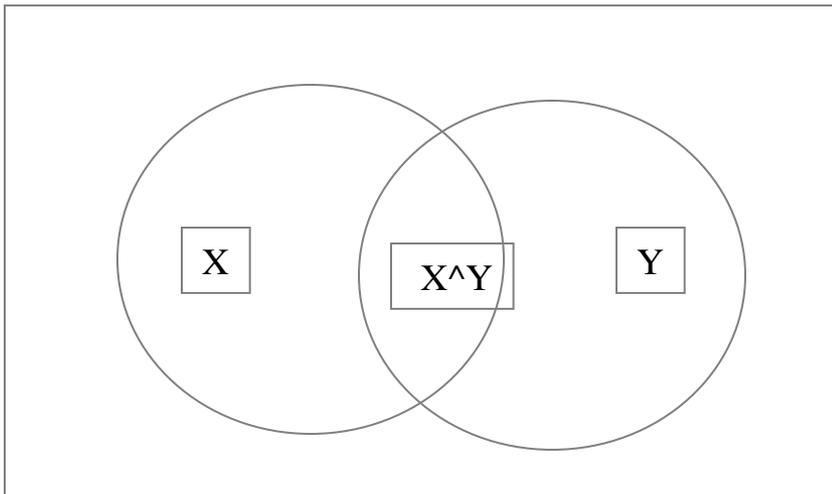
# Naïve Bayes Classifier

*Naïve Bayes classifier* is a highly practical Bayesian learning method.

Naïve Bayes classifier =  
Bayes theorem +  
Independence (Naïve )assumption

# Bayes Theorem

$$P(X | Y) = \frac{P(Y | X)P(X)}{P(Y)}$$



**Proof:**

From Venn Diagram

$$P(X/Y) = P(X^Y) / P(Y)$$

Symmetrically,

$$P(Y/X) = P(X^Y) / P(X)$$

Or

$$P(X^Y) = P(Y/X)P(X)$$

Called the chain rule

# Bayes Classification with $n$ classes

- Suppose we have:

- $n$  predefined categories:  $C_1, C_2, \dots, C_n$

- New document to be categorized:  $d$

In order to categorize new document  $d$ , we calculate:

$$\arg \underset{i}{\text{Max}} P(C_i | d)$$

$$P(C_i | d) = \frac{P(d | C_i) P(C_i)}{P(d)}$$

$$\arg \max_i P(C_i | d) = \arg \max_i \frac{P(d | C_i) P(C_i)}{P(d)} = \arg \max_i P(d | C_i) P(C_i)$$

Where:

$P(C_i)$  is the prior probability of  $C_i$

$P(C_i | d)$  is called a maximum a posteriori (MAP) hypothesis.

# Document representation

- The new document can be represented as a word vector:

$$\langle w_1, w_2, \dots, w_{docLength} \rangle$$

Where:  $w_i$  is the  $i$ th word in the document.

So, now we have:

$$\arg \max_i P(C_i | d) = \arg \max_i P(d | C_i)P(C_i) = \arg \max_i P(\langle w_1, w_2, \dots, w_{|d|} \rangle | C_i)P(C_i)$$

# Naïve Bayes Assumptions

In order to calculate  $P(\langle w_1, w_2, \dots, w_{|d|} \rangle | C_i)$ , we make the following two reasonable Assumptions.

- The words  $w_1, w_2, \dots, w_{|d|}$  are conditionally independent with each other.

$$e.g. \quad P(w_1, w_2, \dots, w_{|d|}) = P(w_1)P(w_2) \cdots P(w_{|d|})$$

- Word positions are independent and identically distributed.

$$e.g. \quad P(\langle w_1, w_2, \dots, w_{|d|} \rangle) = P(\langle w_2, w_1, \dots, w_{|d|} \rangle) = \cdots = P(w_1, w_2, \dots, w_{|d|})$$

# Naïve Bayes Classifier

Naïve Bayes classifier is based on the above assumptions:

$$\begin{aligned} \arg \max_i P(C_i | d) &= \arg \max_i P(d | C_i)P(C_i) \\ &= \arg \max_i P(\langle w_1, w_2, \dots, w_{|d|} \rangle | C_i)P(C_i) \\ &= \arg \max_i P(w_1 | C_i)P(w_2 | C_i) \cdots P(w_{|d|} | C_i)P(C_i) \end{aligned}$$

# The overall implementation

- Training:
  - Construct a vocabulary that is composed by all the distinct words in the training documents
  - Calculate

– Calculate

$$P(C_i) = \frac{\text{number of training docs in } C_i}{\text{total number of training docs}}$$

$$P(v_j | C_i) = \frac{\text{times of } v_j \text{ occurring in } C_i + 1}{\text{total distinct positions in } C_i + |\text{Vocabulary}|}$$

- Categorization

- Suppose the new document  $d$  can be represented as:

$$\langle w_1, w_2, \dots, w_{|d|} \rangle$$

- During categorization, we need to calculate:

$$\begin{aligned} \arg \max_i P(C_i | d) &= \arg \max_i P(d | C_i)P(C_i) \\ &= \arg \max_i P(\langle w_1, w_2, \dots, w_{|d|} \rangle | C_i)P(C_i) \\ &= \arg \max_i P(w_1 | C_i)P(w_2 | C_i) \cdots P(w_{|d|} | C_i)P(C_i) \end{aligned}$$

- If we can find  $w_j$  in the vocabulary, just use the probability value obtained during training
- Otherwise, just disregard  $w_j$

# Electronic Newsgroups Considered in Text Classification Experiment

Given 1000 training documents from each group  
Learn to classify new documents according to  
which newsgroup it came from

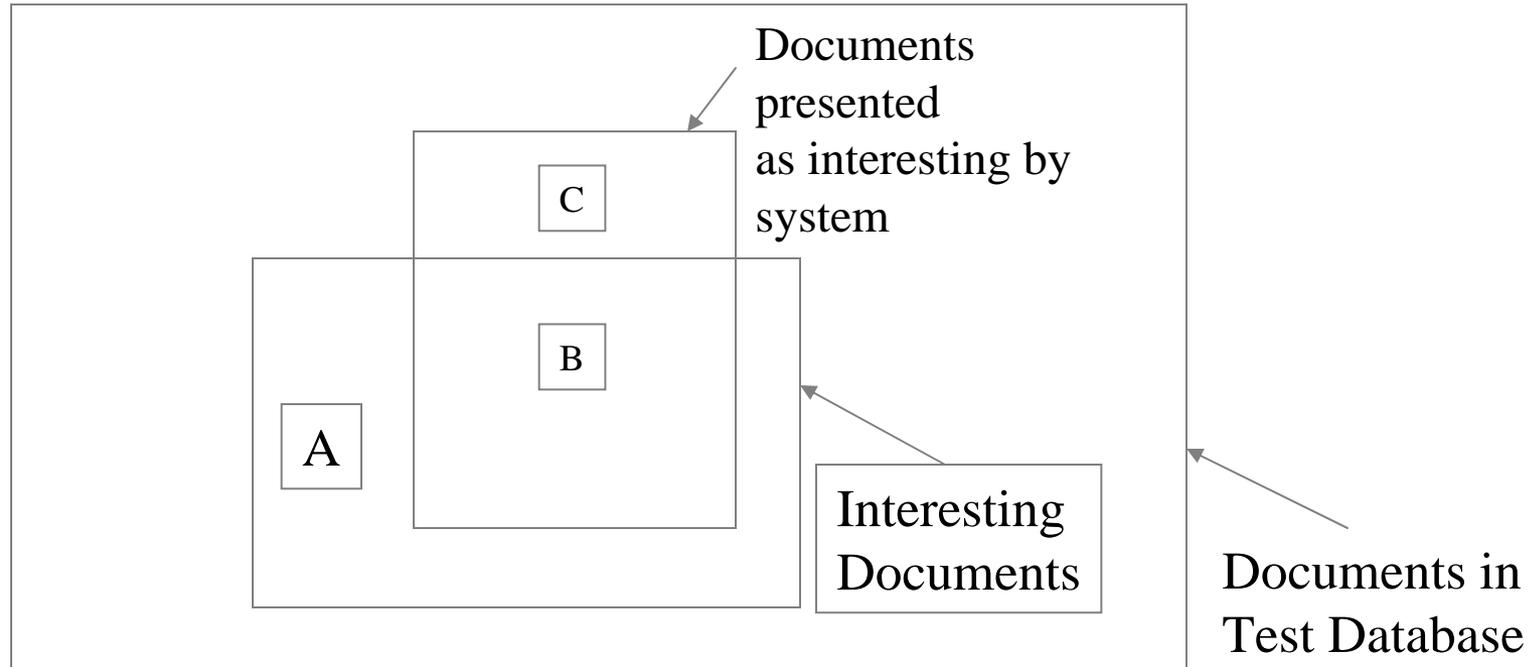
comp.graphics	misc.forsale
comp.os.ms-windows.misc	rec.autos
comp.sys.ibm.pc.hardware	rec.motorcycles
comp.sys.mac.hardware	rec.sport.baseball
comp.windows.x	rec.sport.hockey

alt.atheism	sci.space
soc.religion.christian	sci.crypt
talk.religion.misc	sci.electronics
talk.politics.mideast	sci.med
talk.politics.misc	
talk.politics.guns	

Naive Bayes: 89% classification accuracy

# Precision and Recall Measures

Text Categorizer is run on every document on a correctly labeled database



$$\text{Precision} = B / B + C$$

$$\text{Recall} = B / A$$

# Precision and Recall Evaluation

- Suppose we have 2 predefined categories:
- **C1** (interesting) and **C2** (not interesting)

<b>Test documents</b> \ <b>category</b>	<i>The category the document belongs to</i>	<i>The category assigned to the document by the categorizer</i>
D1	C1	C1
D2	C1	C2
D3	C2	C2
D4	C2	C1
D5	C2	C1

A = 2 = no of interesting documents

B = 1 = no correctly labeled as interesting

C = 2 = no wrongly labeled as interesting

- The precision and recall for C1:

$$\textit{recall} = \frac{1}{2} = 50\% \quad \textit{precision} = \frac{1}{3} = 33.33\%$$

- The precision and recall for C2:

$$\textit{recall} = \frac{1}{3} = 33.33\% \quad \textit{precision} = \frac{1}{2} = 50\%$$

- The overall precision and recall:

$$\textit{recall} = \textit{precision} = \frac{2}{5} = 40\%$$