

# Introduction to Information Retrieval

Hinrich Schütze and Christina Lioma  
Feature Selection  
(from Ch. 13 of IIR Book)

1

## 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**.

2

## 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.

3

## 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    $A(t, c) \leftarrow \text{COMPUTEFEATUREUTILITY}(\mathbb{D}, t, c)$ 
5    $\text{APPEND}(L, (A(t, c), t))$ 
6 return  $\text{FEATURESWITHLARGESTVALUES}(L, k)$ 
How do we compute  $A$ , the feature utility?
    
```

4

## 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)

5

## 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:

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

**Pointwise mutual information (PMI)**  
just for a fixed  $e_t, e_c$

6

## 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  $c$  ( $e_c = 0$ );  $N = N_{00} + N_{01} + N_{10} + N_{11}$ .

7

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

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

$$e_t = e_{\text{EXPORT}} = 0 \quad N_{11} = 49 \quad N_{10} = 27,652 \quad \text{Plug}$$

$$N_{01} = 141 \quad N_{00} = 774,106$$

these values into formula:

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

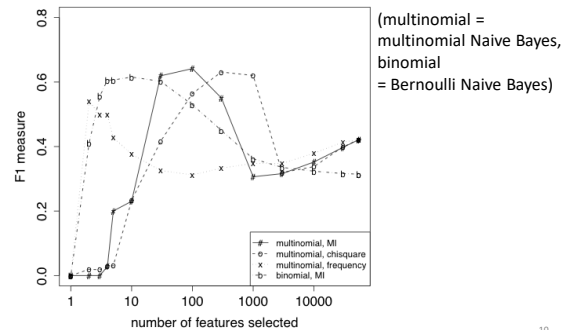
8

## MI feature selection on Reuters

Class: <i>coffee</i>		Class: <i>sports</i>	
term	MI	term	MI
COFFEE	0.0111	SOCCER	0.0681
BAGS	0.0042	CUP	0.0515
GROWERS	0.0025	MATCH	0.0441
KG	0.0019	MATCHES	0.0408
COLOMBIA	0.0018	PLAYED	0.0388
BRAZIL	0.0016	LEAGUE	0.0386
EXPORT	0.0014	BEAT	0.0301
EXPORTERS	0.0013	GAME	0.0299
EXPORTS	0.0013	GAMES	0.0284
CROP	0.0012	TEAM	0.0264

9

## Naive Bayes: Effect of feature selection



10

## Feature selection for Naive Bayes

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

11

## 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:

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

$$e_t = e_{\text{EXPORT}} = 0 \quad N_{11} = 49 \quad N_{10} = 27,652$$

$$N_{01} = 141 \quad N_{00} = 774,106$$

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

12