

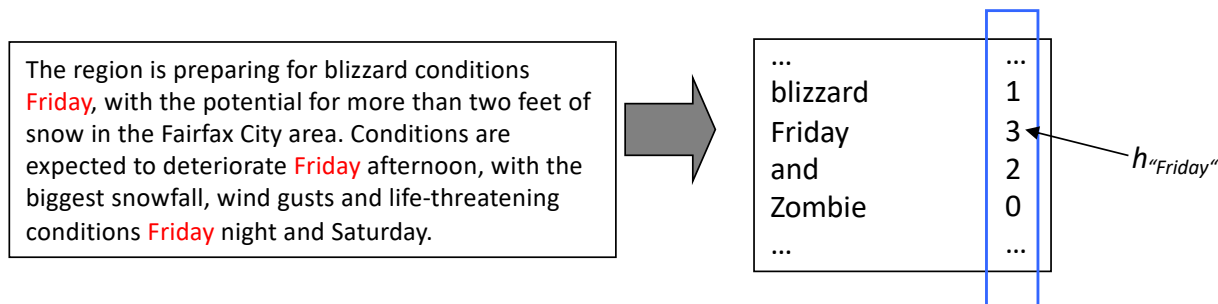
## Outline

---

- Data preprocessing
- Decomposing a dataset: instances and features
- Basic data descriptors
- Feature spaces and proximity (similarity, distance) measures
- Feature transformation for text data

## Feature transformations for text data 1/6

- Text represented as a set of terms (“Bag-Of-Words“ model)
  - Terms:
    - Single words (“cluster“, “analysis“..)  
or
    - bigrams, trigrams, ...n-grams (“cluster analysis“..)
  - Transformation of a document  $d$  in a vector  $r(d) = (h_1, \dots, h_d)$ ,  $h_i \geq 0$ : the frequency of term  $t_i$  in  $d$



## Feature transformations for text data 2/6

---

- Challenges/Problems in Text Mining:
  1. Common words (“e.g.”, “the”, “and”, “for”, “me”)
  2. Words with the same root (“fish”, “fisher”, “fishing”,...)
  3. Very high-dimensional space (dimensionality  $d > 10.000$ )
  4. Not all terms are equally important
  5. Most term frequencies  $h_i = 0$  (“sparse feature space“)
- More challenges due to language:
  - Different words have same meaning (synonyms)
    - “freedom” – “liberty”
  - Words have more than one meanings
    - e.g. “java”, “mouse”

## Feature transformations for text data 3/6

- Problem 1: Common words (“e.g.”, “the”, “and”, “for”, “me”)
  - Solution: ignore these terms (stop-words) or remove stop-words

There are stop-words lists for all languages in WWW.

The following is a list of stop words that are frequently used in English language, but do not carry the thematic component. As a rule in SEO, this set of words trying to exclude in the analysis. We always welcome, if you have any suggestions to change or supplement the list.

Stop Words per Language
Arabic
Brazil
Bulgarian
Czech
Chinese
Dutch
<b>English</b>
German
Greek
Finish
French
Hindi

- 1 a
- 2 able
- 3 about
- 4 above
- 5 abroad
- 6 abst
- 7 accordance
- 8 according
- 9 accordingly
- 10 across
- 11 act

Source: <https://countwordsfree.com/stopwords>

The following is a list of stop words that are frequently used in German language, but do not carry the thematic component. As a rule in SEO, this set of words trying to exclude in the analysis. We always welcome, if you have any suggestions to change or supplement the list.

Stop Words per Language
Arabic
Brazil
Bulgarian
Czech
Chinese
Dutch
English
<b>German</b>
Greek
Finish
French
Hindi

- 1 a
- 2 ab
- 3 aber
- 4 ach
- 5 acht
- 6 achte
- 7 achten
- 8 achter
- 9 achttes
- 10 ag
- 11 alle

The following is a list of stop words that are frequently used in Chinese language, but do not carry the thematic component. As a rule in SEO, this set of words trying to exclude in the analysis. We always welcome, if you have any suggestions to change or supplement the list.

Stop Words per Language
Arabic
Brazil
Bulgarian
Czech
<b>Chinese</b>
Dutch
English
German
Greek
Finish
French
Hindi

- 1 一
- 2 上
- 3 下
- 4 不
- 5 与
- 6 且
- 7 个
- 8 为
- 9 乃
- 10 么
- 11 之

## Feature transformations for text data 3/6

---

- Problem 2: Words with the same root (“fish”, “fisher”, “fishing”,...)

- Solution: Stemming

Map the words to their root

example: "fishing", "fished", "fish", and "fisher" to the root word, "fish".

The root of the words is the output of stemming.

Challenge: Avoid overstemming and/or understemming

Example for overstemming:

university, universal, universities, and universe → “univers”      meaning of the word/term is lost

Example for understemming:

stemming “data” → “dat” and “datum” → “datu” is too weak (understemming),

Better solution would be data, datum → “dat”, but then we would have the problem to treat “date” that would be reduced to “dat” as well.

For English, the Porter stemmer is widely used. (Stemming solutions exist for other languages also)

(Porter's Stemming Algorithms: <http://tartarus.org/~martin/PorterStemmer/index.html>)

More advanced stemmer: snowball stemmer (Porter2), lancaster stemmer

## Feature transformations for text data 4/6

- Problem 3: Too many features/ terms
  - Solution: Select the most important features (“Feature Selection”)
  - Example: average document frequency for a term
    - Very frequent terms appear in almost all documents
    - Very rare terms appear in only a few documents

Ranking procedure:

1. Compute document frequency for all terms  $t_i$  :
2. Sort terms w.r.t.  $DF(t_i)$  and get  $rank(t_i)$
3. Sort terms by  $score(t_i) = DF(t_i) \cdot rank(t_i)$   
e.g.  $score(t_{23}) = 0.82 \cdot 1 = 0.82$   
 $score(t_{17}) = 0.65 \cdot 2 = 1.3$
4. Select the  $k$  terms with the largest  $score(t_i)$

$$DF(t_i) = \frac{\#Docs\ containing\ t_i}{\#All\ documents}$$

Rank	Term	DF
1.	$t_{23}$	0.82
2.	$t_{17}$	0.65
3.	$t_{14}$	0.52
4.	...	...

## Feature transformations for text data 5/6

- Problem 4: Not all terms that are frequent are equally important, what value should be assigned to a term (feature)?
  - Idea: Very frequent terms that are frequent in all documents are less informative than less frequent words. Define such a term weighting schema.
  - Solution: TF-IDF (Term Frequency · Inverse Document Frequency)

Consider both the importance of the term in the document and in the whole collection of documents.

$$TF(t, d) = \frac{n(t, d)}{\sum_{w \in d} n(w, d)} \quad \text{The relative frequency of term } t \text{ in } d \quad [n(t, d) = \# t \text{ in } d]$$

$$IDF(t) = \log\left(\frac{|DB|}{|\{d \mid d \in DB \wedge t \in d\}|}\right) \quad \text{Inverse frequency of term } t \text{ in all DB}$$

$$TF \times IDF = TF(t, d) \cdot IDF(t)$$

Feature vector with TF IDF :  $r(d) = (TF(t_1, d) \cdot IDF(t_1), \dots, TF(t_n, d) \cdot IDF(t_n))$

## Feature transformations for text data 6/6

- Problem 5: in each document, most of the terms have a frequency of  $h_i = 0$   
What distance measure should we use to compare documents?
  - Euclidean distance is not a good idea: it is influenced by vectors lengths (many 0-0 matchings)
  - Idea: use more appropriate distance measures

**Jaccard Coefficient:** Ignore terms absent in both documents (without prior feature transformation)

$$d_{Jaccard}(d_1, d_2) = 1 - \frac{|d_1 \cap d_2|}{|d_1 \cup d_2|} = \frac{|\{t | t \in d_1 \wedge t \in d_2\}|}{|\{t | t \in d_1 \vee t \in d_2\}|}$$

**Cosine Coefficient:** Consider feature-transformed documents (e.g. TF-IDF vectors)

$$d_{\cosinus}(d_1, d_2) = 1 - \frac{\langle d_1, d_2 \rangle}{\|d_1\| \cdot \|d_2\|} = 1 - \frac{\sum_{i=0}^n (d_{1,i} \cdot d_{2,i})}{\sqrt{\sum_{i=0}^n d_{1,i}^2} \cdot \sqrt{\sum_{i=0}^n d_{2,i}^2}}$$