

**Fig. 1.1** Steps in the KDD process  
(**k**nowledge **d**iscovery in **d**ata)

[3]

Step 2: – text-preperation

**syntax:**

tokenization/normalization (98%)\*  
 simplest thing/important thing  
 identifying the units in your text  
 to read the punctuation, e.g.:  
 - dr.  
 - This is a sentence.

**lemmatization:**

reduce wordforms to their dictionary item  
 is/been/was/be  
 --> belongs to 'to be'  
 + plurals --> singulars

**syntactical:**

part-of-speech tagging  
 important elements for object text-mining  
 --> nouns  
 for subjective text-mining  
 --> adjectives  
 word sense disambiguation  
 bank / bank  
 --> river bank / money bank  
 semantic role labeling

**pragmatics: (?)**

named entity recognition  
 co-reference resolution (50%)\*  
 <-- meaning output

\*(% refers to accuracy)

[2]

Step 4: – interpretation

# P A T T E R N

## # data mining (step 3)

- a Google, Twitter and Wikipedia API
- a web crawler
- a HTML DOM parser

## # natural language processing (step 2)

- part-of-speech taggers
- n-gram search
- sentiment analysis
- WordNet

## # machine learning (step 4)

- vector space model
- clustering
- SVM

## # network analysis

## # <canvas> visualization

[1]

PATTERN modules:

### pattern.web

- url downloads
- interval requests
- search engine requests
- use google translate
- crawl
  - wikipedia articles
  - fb comments + reactions
  - dbpedia
  - twitter
  - rss
- parse HTML elements, PDFs
- retrieve emails via imap
- retrieve local information (eg. tweets)

### pattern.db

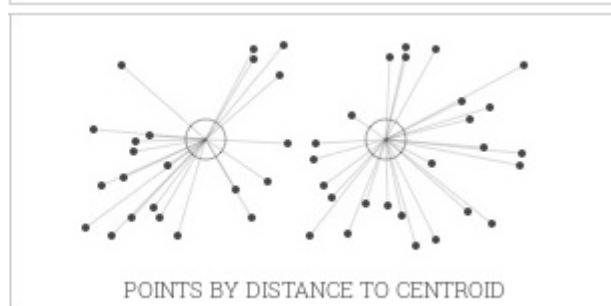
- built database
- work with time/date

### pattern.en |es|de|fr|it|nl

- text preparation
- sentiment analysis tool
- WordNet interface
- wordlists interface

### pattern.search

- a pattern matching system similar to regular expressions, that can be used to search a string by *syntax* (word function) or by *semantics* (word meaning).



[1]



*semantics* (word meaning).  
 - eg.: ('{NP} be \* than {NP}')

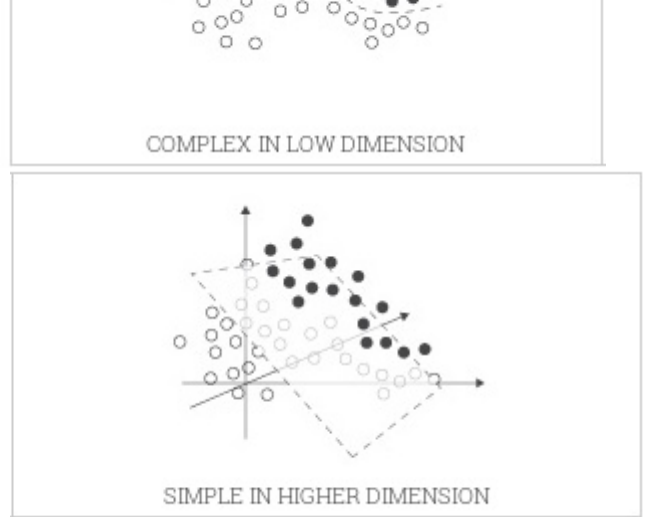
ng

**pattern.vector**

- machine learning tools:
  - word count functions
  - bag-of-words documents
  - a vector space model
  - latent semantic analysis (context analysis)
  - algorithms for *clustering*
    - k-means (similar clusters)
    - hierarchical (nested clusters)
  - and *classification*
    - NB (Naive Bayes)
    - KNN (k-nearest neighbor)
    - SLP (Single-layer perceptron)
    - SVM (Support vector machine)

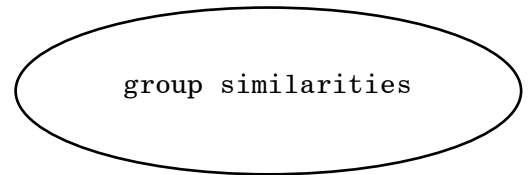
**pattern.graph**

[1]

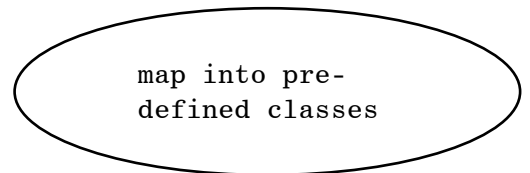


[1]

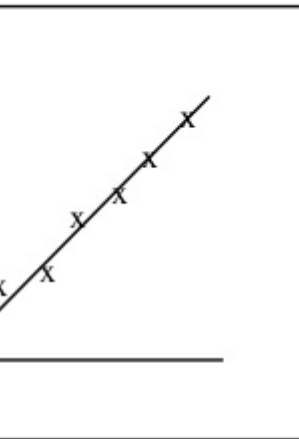
*clustering* unsupervised learning



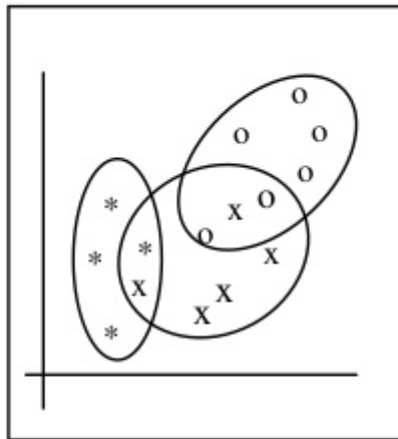
*classification* supervised learning



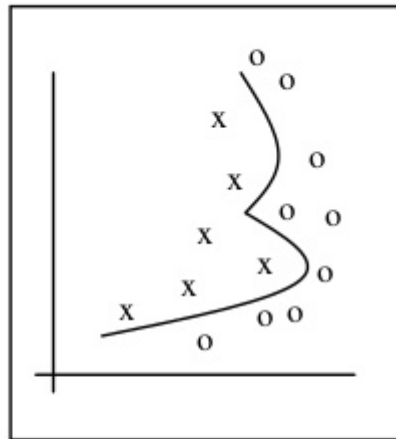
[3]



(A)



(B)



(C)

2 Examples of different types of discovery algorithms: Pattern mining with a linear decision function (A), clustering (B), and classification (C)

**sources**

[1]: <http://www.clips.ua.ac.be/pages/pattern>

[2]: CLiPS - Guy de Pauw, Pattern workshop - Cqrrrelations, January 2015

[3]: Data Mining and Profiling in Large Databases, Bart Custers, Toon Calders, Bart Schermer, and Tal Zarsky (Eds.) (2013)