

# Textual information retrieval Models

Anne-Laure Ligozat

2018/2019<sup>1</sup>

---

1. freely inspired by Benjamin Piwowarski and Hinrich Schütze's courses

# IR systems

- Document indexing
- Query analysis
- Search model
- Evaluation

# Plan

- 1 Indexation
  - Définition
  - Quels documents ?
  - Du texte aux termes
  - Normalisation
  - Index
  - Pondération des termes
  - Utilisation de l'index
  - Index avancés
- 2 Représentation des documents et de la pertinence
  - Modèle booléen
  - Modèle vectoriel
  - Probabilistic model
- 3 Evaluation

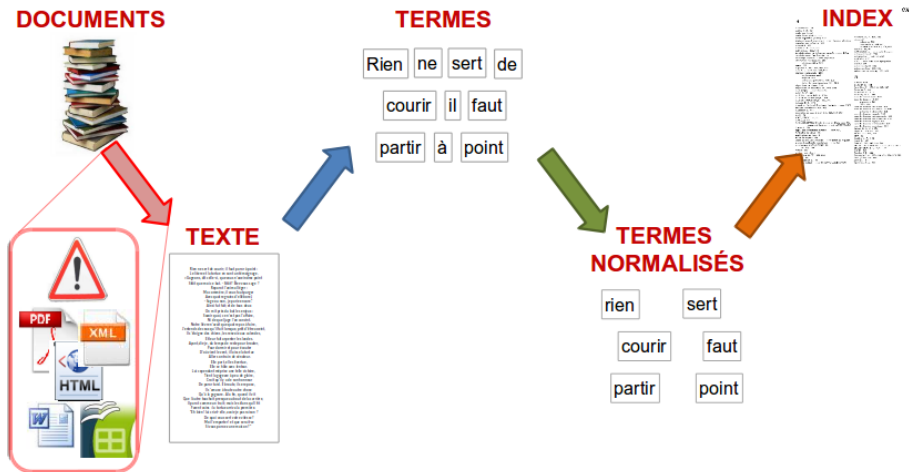
# Why index ?

- Objective : find documents relevant for the user query
  - From the query **words**
  - Impossible to go through the entire collection
    - too many documents → prohibitive response time
    - operations between terms (not, near...) complex
- ⇒ pre-processing = indexing
- goal : “transform documents into substitutes capable to represent the documents content” (Salton et McGill, 1983)

# Free vs controlled indexing

- free indexing : terms from documents
- controlled indexing : predefined terms
  - controlled vocabulary : avoids polysemy, synonymy, granularity problems

# Documents



# Formats

## Formats

- HTML (diff : menus, tables, ads, appearance)
- raw text (structure ?)
- pdf (encoding, appearance)
- word (proprietary format, structure)
- excel (tables)
- openoffice (xml)

## Taking the format into account

- detecting a document type is quite simple
- heuristics specific to each format to extract the text
- search engines rarely use the document structure

# Language and encoding

## Languages

- language identification = difficult problem
- multilingual information retrieval possible

## Encoding

- errors in dealing with encoding → incorrect results

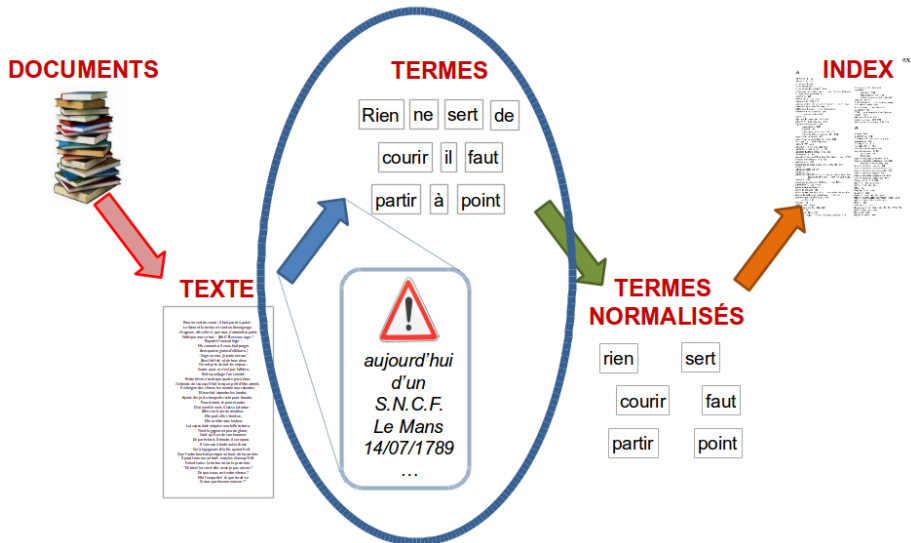


# Document content

unit =

- file ?
- e-mail ?
  - with heading ?
  - with attachments ?
- set of files
  - website
  - document with multiple files
- ...

# From text to terms



# Tokenization

Tokenization = identification of elementary units

- **word** : substring as it appears in the text
- **term** (or type) : normalized word (case, morphology, spelling...)
  - set of terms = dictionary
- **token** : instance of a word or a term in a document

## Difficulties of tokenization

- graphical variants of words with possible separators
  - États-Unis ou États Unis
- composed words in agglutinative languages
  - Lebensversicherungsgesellschaftsangestellter (employee from a life insurance company)
- multiple alphabets in Japanese for example
- bidirectionnality of the reading sense in Arabic (numerals and letters)
- numbers : 555 3424, 24.09.2018
- ...

## DOCUMENTS



## TEXTE

[illegible]

## TERMES

Rien ne sert de  
courir il faut  
partir à point



aujourd'hui  
d'un  
S.N.C.F.  
Le Mans  
14/07/1789

## INDEX

[illegible]

## TERMES NORMALISÉS

rien      sert

courir      faut

partir      point

# Variant normalization (1/2)

- in documents and query
- variants that should be grouped
  - word variants including punctuation
    - U.S.A. and USA
    - morpho-syntaxe and morphosyntaxe
  - diacritic variants
    - in German, Tuebingen, Tübingen and Tübingen
    - in English, resume = résumé
  - proper names variants
    - Gorbachov and Gorbachev
- but
  - accents can be relevant
    - sur and sûr
    - pêche and péché
  - the case can be relevant
    - in German, mit and MIT
    - (interaction between normalization and language detection)
    - in English, fed and Fed

# Variant normalization (2/2)

- possibly asymmetric
  - window → window, windows
  - windows → Windows, windows
  - Windows : non expansion
- + typos or spelling, OCR errors
- important criterion : how will the users write their query most often ?

# Morphological normalization

- use of analyses
  - lemmatization : chanteurs → chanteur, chantions → chanter
  - stemming : automate, automatique, automatiser → automat
    - in particular Porter algorithm, classical algorithm for English
    - stemming useful for certain queries, strongly worsens the results for others
  - POS tagging
- technics considered as mainly solved : low percentage of errors but difficult to reduce

# Stop words

- word that do not add sense to the text
  - determiners : le, la
  - pronouns : je, nous
  - prepositions : sur, contre
- they are the most frequent in a language
  - the 30 most frequent words represent 30% of the word occurrences
  - delete them enables to gain much space in the index
- but
  - useful for multi-term queries : “pomme de terre”, “les Chevaliers du Zodiaque”
  - sometimes bring sense in particular cases : “Let it be”, “The Who”, “être ou ne pas être”
  - compression actually enables to keep stopwords with little space



# Index

[illegible]

Rien ne sert de  
courir il faut  
partir à point

  
aujourd'hui  
d'un  
S.N.C.F.  
Le Mans  
14/07/1789  
...

rien      sert  
courir      faut  
partir      point

**INDEX** <sup>EX</sup>

**A**

1911-1912

1913-1914

1915-1916

1917-1918

1919-1920

1921-1922

1923-1924

1925-1926

1927-1928

1929-1930

1931-1932

1933-1934

1935-1936

1937-1938

1939-1940

1941-1942

1943-1944

1945-1946

1947-1948

1949-1950

1951-1952

1953-1954

1955-1956

1957-1958

1959-1960

1961-1962

1963-1964

1965-1966

1967-1968

1969-1970

1971-1972

1973-1974

1975-1976

1977-1978

1979-1980

1981-1982

1983-1984

1985-1986

1987-1988

1989-1990

1991-1992

1993-1994

1995-1996

1997-1998

1999-2000

2001-2002

2003-2004

2005-2006

2007-2008

2009-2010

2011-2012

2013-2014

2015-2016

2017-2018

2019-2020

2021-2022

2023-2024

2025-2026

2027-2028

2029-2030

2031-2032

2033-2034

2035-2036

2037-2038

2039-2040

2041-2042

2043-2044

2045-2046

2047-2048

2049-2050

2051-2052

2053-2054

2055-2056

2057-2058

2059-2060

2061-2062

2063-2064

2065-2066

2067-2068

2069-2070

2071-2072

2073-2074

2075-2076

2077-2078

2079-2080

2081-2082

2083-2084

2085-2086

2087-2088

2089-2090

2091-2092

2093-2094

2095-2096

2097-2098

2099-2100

2101-2102

2103-2104

2105-2106

2107-2108

2109-2110

2111-2112

2113-2114

2115-2116

2117-2118

2119-2120

2121-2122

2123-2124

2125-2126

2127-2128

2129-2130

2131-2132

2133-2134

2135-2136

2137-2138

2139-2140

2141-2142

2143-2144

2145-2146

2147-2148

2149-2150

2151-2152

2153-2154

2155-2156

2157-2158

2159-2160

2161-2162

2163-2164

2165-2166

2167-2168

2169-2170

2171-2172

2173-2174

2175-2176

2177-2178

2179-2180

2181-2182

2183-2184

2185-2186

2187-2188

2189-2190

2191-2192

2193-2194

2195-2196

2197-2198

2199-2200

2201-2202

2203-2204

2205-2206

2207-2208

2209-2210

2211-2212

2213-2214

2215-2216

2217-2218

2219-2220

2221-2222

2223-2224

2225-2226

2227-2228

2229-2230

2231-2232

2233-2234

2235-2236

2237-2238

2239-2240

2241-2242

2243-2244

2245-2246

2247-2248

2249-2250

2251-2252

2253-2254

2255-2256

2257-2258

2259-2260

2261-2262

2263-2264

2265-2266

2267-2268

2269-2270

2271-2272

2273-2274

2275-2276

2277-2278

2279-2280

2281-2282

2283-2284

2285-2286

2287-2288

2289-2290

2291-2292

2293-2294

2295-2296

2297-2298

2299-2300

2301-2302

2303-2304

2305-2306

2307-2308

2309-2310

2311-2312

2313-2314

2315-2316

2317-2318

2319-2320

2321-2322

2323-2324

2325-2326

2327-2328

2329-2330

2331-2332

2333-2334

2335-2336

2337-2338

2339-2340

2341-2342

2343-2344

2345-2346

2347-2348

2349-2350

2351-2352

2353-2354

2355-2356

2357-2358

2359-2360

2361-2362

2363-2364

2365-2366

2367-2368

2369-2370

2371-2372

2373-2374

2375-2376

2377-2378

2379-2380

2381-2382

2383-2384

2385-2386

2387-2388

2389-2390

2391-2392

2393-2394

2395-2396

2397-2398

2399-2400

2401-2402

2403-2404

2405-2406

2407-2408

2409-2410

2411-2412

2413-2414

2415-2416

2417-2418

2419-2420

2421-2422

2423-2424

2425-2426

2427-2428

2429-2430

2431-2432

2433-2434

2435-2436

2437-2438

2439-2440

# Incidence matrix



	Antoine & Cléopâtre	Jules César	La Tempête	Hamlet	Othello	Macbeth
Antoine	1	1	0	0	0	0
Brutus	1	1	0	1	0	0
César	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cléopâtre	1	0	0	0	0	0
pitié	1	0	1	1	1	1
pire	1	0	1	1	1	0

# Incidence matrix

**Brutus ET Cléopâtre ET PAS Calpurnia**

	Antoine & Cléopâtre	Jules César	La Tempête	Hamlet	Othello	Macbeth
Antoine	1	1	0	0	0	0
Brutus	1	1	0	1	0	0
César	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cléopâtre	1	0	0	0	0	0
pitié	1	0	1	1	1	1
pire	1	0	1	1	1	0

Vecteurs d'incidence

$\neg$ Calpurnia	1	0	1	1	1	1
------------------	---	---	---	---	---	---

**ET "bit à bit"**

1	0	0	0	0	0
---	---	---	---	---	---

# Incidence matrix

- impossible to use in practice
    - collection of a million documents
    - about 1000 words per document
    - total vocabulary of 500,000 distinct words
- how many objects in the matrix? how many 0s? and 1s?

# Inverted index



## Index

### A

Abiteboul S. 799, 948  
 Abern J. E. 424  
 accès aux données distantes – voir RDBA  
 accès direct (index) 846  
 accès séquentiel (index) 846  
 accès séquentiel (physique) – voir séquenceur physique  
 accéder une référence 911  
 Adams S. 620  
 Adiba M. E. 575, 734  
 Adkins L. 550, 519  
 administrateur de la base de données – voir DBA  
 administrateur des données 15  
 adossement dispersé – voir dispersion  
 affectation relationnelle 188  
 relation cible 569  
 appel 707  
 Agrawal R. 861, 841, 845, 848  
 Aho A. V. 382, 382, 824, 946  
 algèbre relationnelle 188  
 implémentation 808  
 objectif 188  
 opérations primitives 188, 203  
 règles de transformation 181, 892  
 algorithme de chasse 380  
 algorithme de réduction de Codd 378  
 ALI2 (SQL) – voir implémentation  
 Allen F. W. 424  
 ALPHA – voir DSL ALPHA  
 ALTER DOMAIN (SQL) 258  
 ALTER TABLE (SQL) 186, 203  
 Alliman E. B. 873  
 American National Standards Institute – voir ANSI  
 analyse membership 456, 465  
 Anderson R. 665  
 annuaire de notes à jour 345, 349, 356, 378  
 ANSI 73  
 ANSI/SPARC 33, 57  
 ANSI/X3/SPARC Study Group on Data Base Management Systems – voir ANSI/SPARC  
 Anton J. 664  
 appel des procédures distantes – voir RPC  
 APPEND (SQL) 456  
 applications en ligne 9  
 arbre de requête 388  
 arbre de trierage abstrait – voir arbre de requête  
 arbres de recherche binaire – voir trie  
 architecture ANSI/SPARC 33  
 voir SQL 90  
 ARJIS 453  
 arité – voir degré  
 Armstrong W. W. 326, 328  
 Arpa M. 563  
 Ashenburt R. L. 82  
 association (SQL) – voir CREATE ASSOCIATION

association 12, 485, 486, 419  
 OO 790  
 réécriture 414  
 association (RDB, 71) 459  
 associativité 170  
 Astrahan M. M. 205, 206, 873  
 Atkinson M. P. 890, 872  
 atomique  
 relations 353  
 transactions 436, 442  
 valeurs scalaires 65, 104, 642  
 attribut 89, 97  
 attribution – voir mot de passe  
 attribution locale 700  
 autorisation – voir sécurité  
 auxiliaire 427  
 AVL – voir fonction d'aggrégation  
 axiome 906  
 axiome déductif 925  
 axiome de base 909, 923  
 axiomes de Armstrong 220, 528

### B

B-trees 850  
 Babel D. E. 534  
 Bancher P. 823, 944, 945, 947  
 Banerjee J. 802  
 Barnes G. M. 872  
 Bentley J. E. 884  
 base de connaissances 880  
 base de données 3, 10  
 ouvrages 18  
 DB2 890  
 base de données déductive 310  
 base de données distribuée 52, 695  
 politique fondamentale 698  
 base de données experte 928  
 base de données calculatoire 322  
 base de données transactionnelle 520  
 base de données logique 940  
 base de données relationnelle 118  
 base de données statistique 368  
 Banerjee D. S. 813, 874  
 Bayer R. 478, 851, 875  
 BCNP 137, 354, 361  
 BIA 16  
 Buckley D. A. 884  
 Bush B. 825  
 Bussell C. 382, 382, 383, 946  
 BEGIN DECLARE SECTION (SQL) 283  
 BEGIN TRANSACTION 429  
 Bels D. 729  
 Bentley J. L. 884  
 Benoit P. A. 343, 419, 479, 504, 576, 739  
 Bird (DB) 894, 905  
 Bitton D. 817  
 Björnerstedt A. 880

# Inverted index

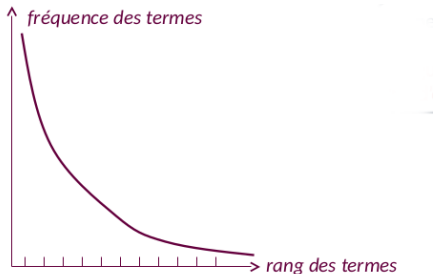
- classical notion of index
- associates indexes to documents that contain them (unique identifier)
  - a  $\rightarrow$  d1, d2, d3, d4, d5...
  - à  $\rightarrow$  d1, d2, d3, d4, d5...
  - abaissa  $\rightarrow$  d3, d4...
  - abaissable  $\rightarrow$  d5
  - abandon  $\rightarrow$  d1, d5
  - abandonna  $\rightarrow$  d2
  - ...

# Vocabulary size

- The vocabulary increases with the collection
- Heaps law :  $M = kT^b$  with
  - M vocabulary size
  - T number of tokens in the collection
  - b and k constants (typically  $b=0,5$  and  $k = 30$  to  $100$ )
  - empirical law
- much worse for the Web !

# Term frequency

- few frequent words and many rare words
- Zipf law : the  $n^{\text{th}}$  most frequent word has a frequency (= number of occurrences) proportional to  $1/n$





# tfidf

- in a query like in a document, all terms do not have the same importance
- intuition 1 : the more occurrences of a word a document contains, the more it is about this term  $\rightarrow tf_{t,d}$  = number of occurrences of term  $t$  in document  $d$
- intuition 2 : very frequent words in all documents are less important (less discriminating)  $\rightarrow df_t$  = number of documents that contain term  $t$
- weight of a term  $tf.idf_{t,d} = tf_{t,d} \times \log_{10} \frac{N}{df_t}$  ( $N$  = number of documents)

# Weighted matrix

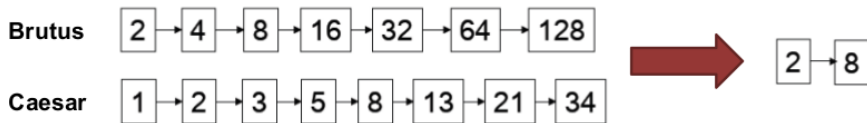
	Antoine & Cléopâtre	Jules César	La Tempête	Hamlet	Othello	Macbeth
Antoine	13,1	11,4	0	0	0	0
Brutus	3,0	8,3	0	1	0	0
César	2,3	2,3	0	0,5	0,3	0,3
Calpurnia	0	11,2	0	0	0	0
Cléopâtre	17,7	0	0	0	0	0
pitié	0,5	0	0,7	0,9	0,9	0,3
pire	1,2	0	0,6	0,6	0,6	0

each document is a vector in  $\mathbb{R}^{|V|}$

# Finding the documents

## Brutus AND Caesar

- On **recherche** « **Brutus** » dans le dictionnaire  
→ On récupère la liste de documents
- On **recherche** « **Caesar** » dans le dictionnaire  
→ On récupère la liste de documents
- On **fusionne** les deux listes



# Queries with phrases

- If a user formulates a query such as "Paris Saclay", a document containing "Le maire de Paris s'est arrêté dans un restaurant de Saclay aujourd'hui" is not likely to be relevant
- concept of phrase easy to understand for users
- a large part of web queries
- simple index insufficient
- two solutions :
  - n-grams index
  - positional index

# N-gram

- n-gram = subsequence of  $n$  elements extracted from a given sequence
- here, word n-grams
  - unigrams : all words
  - bigrams : sequences of 2 words
- différent from a linguistic phrase

# Bigrams index

- index, in addition to simple words, the text bigrams
- 1 bigram = 1 term from dictionary
- actually rarely used
  - query n-grams difficult to determine ([Stanford University Palo Alto](#), [Université Paris-Saclay Orsay](#))
  - index vocabulary very large

# Positional index

- Idea : in the lists of documents in the index, add the position of each term occurrence in the document

terme	fréquence	→	D1	D3	D4
-------	-----------	---	----	----	----



terme	fréquence	→	D1 : pos1, pos2, pos3
			D3 : pos1, pos2
			D4 : pos1, pos2, pos3

# Traversing a positional index

- "Université Paris Saclay"

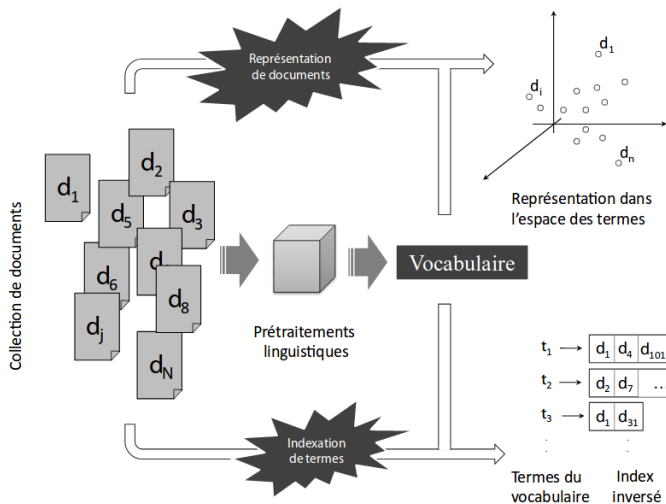
- extraction of all dictionary entries
- recursive use of the fusion algorithm, for documents then for positions
- use of an incremental comparison instead of a strict equality

université	1252	→	D2 : 546
			D6 : 34, 87, 145, 243
			D7 : 44, 87, 34
			...
paris	45	→	D2 : 547
			D6 : 88, 543
			...
saclay	15345	→	D2 : 54, 90
			D6 : 89
			D4 : 43



# Plan

- 1 Indexation
  - Définition
  - Quels documents ?
  - Du texte aux termes
  - Normalisation
  - Index
  - Pondération des termes
  - Utilisation de l'index
  - Index avancés
- 2 Représentation des documents et de la pertinence
  - Modèle booléen
  - Modèle vectoriel
  - Probabilistic model
- 3 Evaluation



# Search models : the three currents

- models based on set theory
  - boolean model
- algebraic models
  - vectorial model
- probabilistic models
  - notion of relevance

from the beginning of IR (60s, 70s)

# Boolean model

- first and simplest model
- based on set theory and Boolean algebra
- query terms either present or absent : binary weighting of terms, 0 or 1
- document soit pertinent soit non pertinent : pertinence binaire (modèle exact)
- query expressed with logical operators : AND, OR, NOT
  - (cyclisme OR natation) AND NOT dopage
  - document relevant iff its content respects the logical formula

# Boolean model : example

Requête **Q** : (cyclisme OR natation) AND NOT dopage

Le document contient					Pertinence du document
cyclisme	natation	cyclisme OR natation	dopage	NOT dopage	
0	0	0	0	1	0
0	0	0	1	0	0
0	1	1	0	1	1
0	1	1	1	0	0
1	0	1	0	1	1
1	0	1	1	0	0
1	1	1	0	1	1
1	1	1	1	0	0

# Boolean model : pros and cons

## Pros

- precise : document contains terms or non
- understandable
- still used in many tools, such as electronic mail
- adapted for specialists when constraint vocabulary preferred (law)

## Cons

- difficult to express long queries with boolean form
- binary criterion not very efficient
  - often too many of too few results
  - term weighting improves results (cf. extended boolean model)
- impossible to rank the results
  - all returned documents are at the same level
  - users prefer a ranking when the list is long

# Ranked lists of results

## Why rank results?

- most users
  - do not know how to write boolean queries
  - do not want to parse too many results (possibly millions of them)
- prefer ranked lists
  - from most useful to the user (relevant) to less useful
  - the number of results is no more a problem
  - the user parses as many as they want
- but requires an effective ranking algorithm
- statistical model
  - quantitative aspect of terms and documents
  - similarity measure between query and document

# Ranked models

## Ranking

- the number of results is not a problem any more : first 10
- if the ranking algorithm works correctly

## Principle

- attribure a score to each query-document pair
  - according to the document relevance to the query
  - generally presence of query terms in the document
- rank documents by descending score



# Vector space model

## Vector space model

- Similarity measure : close representations  $\Rightarrow$  high probability that same information
- Documents and queries represented by vectors in a euclidian space with  $n$  dimensions ( $n$  : number of terms)
  - terms = axes
  - docs = vectors (sparse)
- Document relevance = degree of similarity between vectors of query and document
- Documents ranked from most similar to the query to less similar

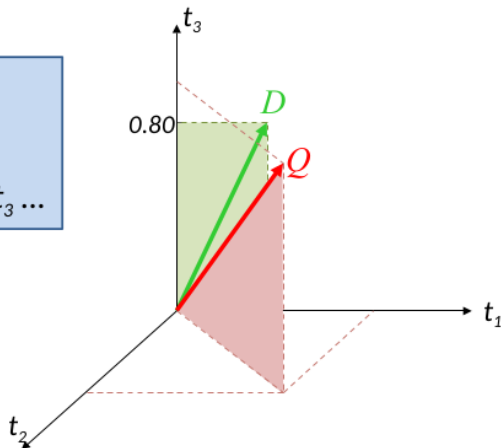
# Vector space model

Requête  $Q$ :  $t_1 t_2 t_3$

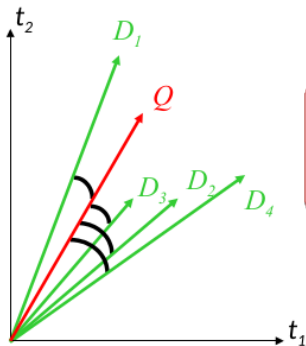
Document  $D$ :  $\dots t_1 \dots t_3 \dots$

Poids  $w_{D,t_1} = 0.45$

Poids  $w_{D,t_3} = 0.80$



# Similarity measure



## Cosinus

$$\text{sim}(\vec{Q}, \vec{D}) = \frac{\vec{Q} \cdot \vec{D}}{|\vec{Q}| \times |\vec{D}|} = \frac{\sum_{i=1}^n w_{i,Q} \times w_{i,D}}{\sqrt{\sum w_{i,Q}^2} \times \sqrt{\sum w_{i,D}^2}}$$

(Le produit scalaire avec  
normalisation de la longueur  
des vecteurs)

# Vector space models : pros and cons

## Pros

- simple query language : keywords list
- performance better than boolean thanks to term weighting
- partial relevance of documents possible
- ranking of documents possible

## Cons

- terms considered as independent
- query language less expressive
- understanding harder

# Probabilistic model

## Problem modelling

- estimation of the relevance probability of a document for a query
- binary notion of relevance :
  - $R_{d,q} = 1$  if d relevant for q
  - $R_{d,q} = 0$  otherwise
- documents ranked by descending relevance probability
- relevance of each document supposed to be independent
- dyssymetry between query and document ( $\neq$  vector space)

# Probabilistic model : conclusion

- leading model :
  - Okapi BM25
    - non binary probabilistic model (term frequency) with document length normalization
    - robust and used a lot
- other probabilistic models
  - bayesian networks
  - language model
    - document = generative model that generates the query
- conclusion
  - problem of initial probabilities
  - independent terms
  - results close to vector space model

# Learning to rank

## Basic principle

- features that influence the relevance
  - query : length, idf means...
  - document : PageRank, spam probability, document date...
  - query + document : cosine similarity between query and document, minimal window in which query terms appear, zones...
- binary classification : relevant 1, not relevant 0
  - but how to learn regression score
- (document) pair learning
  - but errors do not have the same importance
  - sensitive to the number of relevant documents per query
  - first returned documents more important
- apprentissage de listes : optimiser MAP directement par exemple
  - difficult ML problem

# Other models

- generalized vector space model
  - represents dependencies between terms
  - theoretically interesting, but effectiveness not proved
- Latent Semantic Indexing
  - proposes to study "concepts" instead of terms (idead of a text)
  - links documents between them and to query
  - enables to return documents containing no query word
  - less dimensions
  - better recall, lower precision



# Plan

- 1 Indexation
  - Définition
  - Quels documents ?
  - Du texte aux termes
  - Normalisation
  - Index
  - Pondération des termes
  - Utilisation de l'index
  - Index avancés
- 2 Représentation des documents et de la pertinence
  - Modèle booléen
  - Modèle vectoriel
  - Probabilistic model
- 3 Evaluation

# What is a good search engine ?

## Criteria

- main criterion : user satisfaction ?
- fast
  - fast query analysis
  - fast index search
  - fast result ranking
- complete and up-to-date
  - all (or many) documents from the collection are used
  - new documents quickly added
  - ⇒ fast index building
  - ⇒ (on the Web) permanent, effective and fast discovery of the new documents
- most important : relevant

# How to measure relevance ?

- Web search engine
  - the user clicks on some links and not others
  - the user comes back to the search engine
  - the user manages to perform a certain task
- e-commerce website
  - the user buys
  - the user buys quickly
  - a large proportion of visitors buy
- firm website
  - the user gains productivity
  - secured access
  - ...

# What is a good evaluation ?

- evaluating a system aims at knowing :
  - if it performs the expected task
  - if it is better than the competition
  - how to improve it
- that means that the evaluation should be
  - reproducible
    - to evaluate several systems the same way
    - to assess the progress made
  - understandable
    - to identify the possible progress zones
  - fast
    - to evaluate each modification independently
  - objective

# How to make relevance objective ?

- user need transformed into a query → 1st information loss
  - information need : je voudrais savoir si boire du vin rouge réduit le risque de problèmes de coeur
  - query : vin rouge problèmes coeur
  - doc : le coeur de son discours concernait le problème de l'industrie du vin qui peine à reconnaître le rôle de la consommation de vin rouge dans les accidents de voiture
- doc relevant for the query but not for the need
- yet, relevance of the results compared to the information need
- relevance not binary : very relevant, not at all, a little, why not...
- in order to make relevance objective, simplified definition :
  - documents treated independently from each other
  - relevance transformed into a binary notion
- and use of test collections

# Standard methodology

- document collection
  - representative of real documents
- set of information needs/queries
  - also representative
- relevance score for each document and each query
  - human judgements

standard benchmarks : TREC

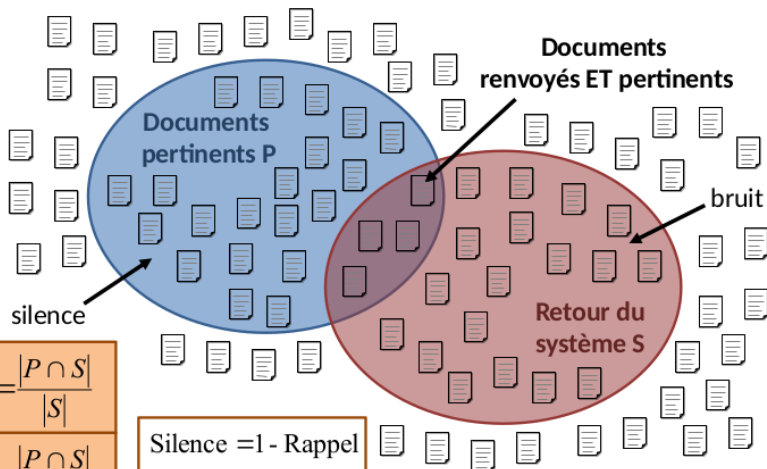
- Ad Hoc in particular (1992 to 1999)
- scores for the first k documents returned by systems

# Test collections

Test collections make the experiments reproducible

- A protocol is defined
- A significative number of examples is manually judged
  - gold standard
  - A part can also serve as a development and/or training set
- An inter-annotator agreement is computed
  - To validate the objectivity
- System results are compared to expected results
- Imperfect but precise metrics are defined

# Evaluation : precision and recall



$$\text{Précision} = \frac{|P \cap S|}{|S|}$$

$$\text{Rappel} = \frac{|P \cap S|}{|P|}$$

Silence = 1 - Rappel

Bruit = 1 - Précision



# Complementarity of precision and recall

## Why not just precision ?

- precision = capacity of a system to return MOSTLY relevant documents
- Returning a single relevant document  $\Rightarrow$  100% precision
- not compatible with user satisfaction !

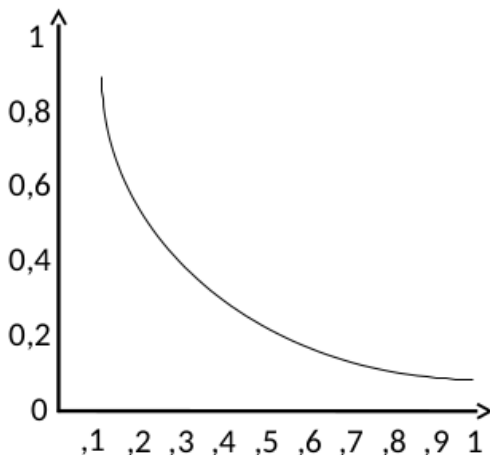
## Why not just recall ?

- recall = capacity of a system to return ALL relevant documents
- Returning all the collection  $\Rightarrow$  100% recall
- not compatible with user satisfaction !

# Recall/precision curve

- recall improves with number of answers
- precision decreases

recall/precision curve used to characterize IR systems



# F-measure

to obtain a unique value, F-measure = harmonic mean

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(\beta^2+1) \times P \times R}{\beta^2 P + R} \text{ with } \alpha = \frac{1}{\beta^2+1}$$

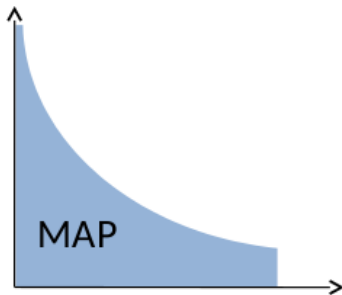
$\beta < 1$  favors precision,  $\beta > 1$  recall

to give as much importance to precision and recall,  $\beta = 1$

$$F = \frac{2PR}{P+R}$$

# Other metrics

- MAP (Mean Average Precision) : area under the R/P curve
- to take into account result ranking :
  - P@5, P@10 : precision for first documents retrieved ; favors high precision
  - P@100
  - R/P curve for k variable
- error rate = (false positives + false negatives) / relevant
- and others...



# Présentation des résultats

- the user must be able to identify the possibly relevant documents with their description → title, url, metadata
- often abstract too
  - static (independent of query)
  - dynamic : “snippets” from document that contains the query terms

# References

- Recherche d'information, Applications, modèles et algorithmes de Massih-Reza Amini et Éric Gaussier (2e édition en 2017)
- Introduction to Information Retrieval, Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze (2008)