Information Retrieval (Text Processing)

Fabio Aiolli

http://www.math.unipd.it/~aiolli

Dipartimento di Matematica Pura ed Applicata Università di Padova

Anno Accademico 2008/2009

Dip. di Matematica Pura ed Applicata F. Aiolli - Information Retrieval 2008/2009

Techniques for Text Retrieval

Text Indexing:

- using the "laws" of statistical linguistics
- Distributional characteristics of terms within docs and within the collection for estimating relevance
- Text pre-processing
 - Removing noise (including stop words, prefixes, and suffixes) for a better (compact) representation of th 'meaning' of queries and docs
- Lexical resources for polysemy, omonymy, and synonymy resolution
 - "Normalising the (natural) language used in docs and queries in order to alleviate the vocabulary mismatch problem

1



$\nu \mu$.	ui Matematica	
Pura	ed Applicata	

TF and IDF

A Popular way to implement the previous considerations

(1.a) is implemented by making w_{ki} grow with the term frequency of t_k in d_i

$$tf(t_k, d_i) = \begin{cases} 1 + \log \#(t_k, d_i) & \text{if } \#(t_k, d_i) > 0\\ 0 & \text{otherwise} \end{cases}$$

(2.b) is implemented by removing from consideration all the terms which occurs in less than α docs (1 $\leq \alpha \leq$ 5)

(2.a) and (2.c) are implemented by making w_{ki} grow with the inverse document frequency of t_k

$$idf(t_k) = \log \frac{|C|}{\#c(t_k)} \tag{1}$$

Dip. di Matematica Pura ed Applicata F. Aiolli - Information Retrieval 2008/2009 5

6

TFIDF

The final weights are obtained by normalizing by cosine normalization, i.e.

$$w_{ki} = \frac{tfidf(t_k, d_i)}{\sqrt{\sum_{s=1}^{n} tfidf(t_s, d_i)^2}} = \frac{tf(t_k, d_i) \cdot idf(t_k)}{\sqrt{\sum_{s=1}^{n} (tf(t_s, d_i) \cdot idf(t_s))^2}}$$

Note that tf and tfidf equals 0 for terms ${\sf t}_k$ that not occur in ${\sf d}_i.$ This shows why the IREPs of documents and queries are sparse vectors

Several variants of the tfidf but they share (i) a tf-like component (ii) idf-like components and (iii) length normalization

TFIDF class of functions is the most popular class of weighting functions!

Dip. di Matematica	F. Aiolli - Information Retrieval
Pura ed Applicata	2008/2009

Text Pre-processing

	ext pre-proce	ssing consists in general of the following	j steps
1.	Reduction in characters,	to ASCII format (e.g. removal of formatt) etc.)	ng/style
2. 3.	Conversion of Identification delimited by	of uppercase into lowercase on of the 'words' (strings of contiguous cho ' blanks) within the text	aracters
4.	Removal of p	punctuation from 'words'	
5.	Removal of r	numbers	
6. 7.	Removal of s Grouping wo morphologic	stop words rds (conflation), typically those that share al root (stemming)	the same
] No in	ote however t formation to	hat what is 'noise' to a given task may another!	be
Dip. di Ma Pura ed A	itematica pplicata	F. Aiolli - Information Retrieval 2008/2009	7

Stop word removal

In order to improve efficiency, words that occur too frequently in the collection (stop words) are removed, since they have almost null discrimination value:

- E.g. articles, pronouns, prepositions, very frequent verbs (e.g. have) and adverbs (e.g. well)
- This substantially reduces (up to 30%!) the size of data structures in secondary storage

Pre-compiled lists of stop words (stop lists) are available in the public domain for every major language. Alternatively, they can be looked up in grammar books, or automatically extracted from text corpora by exploiting the fact that

- They have very high DF
- They have almost constant (i.e. document-independent) relative TF, i.e. for any stop word s_k and documents d_1 and d_2 it is the case that

$$\frac{\#(s_k, d_1)}{|d_1|} \approx \frac{\#(s_k, d_2)}{|d_2|}$$

Dip. di Matematica Pura ed Applicata

F. Aiolli - Information Retrieval 2008/2009

Stemming



- computation* comput*
- Apt for languages with simple morphological structure.
 - Non null error-rate

Tabular: table of n pairs <word,stem>

- Huge!
- Apt for languages with complex morphological structure.
- Better precision, smaller efficiency



	This would reduce the loss in recall due to				
	Spelling checking for an out-of-vocabulary (OOV) word w_1 can be done by picking w_2 such to minimize the edit distance. E.g. minimum				
	deletions s	switches. N	lon null err	ror rate.	
	 Useful for OCR'ed documents Lower error rate for queries . We can use feedback from the user 				
]					
Dip.	di Matematica	F. Aiolli - Inforr	nation Retrieval		13
	Conder (1.		PaceBank ABA		
	Coogle britnev spears spe Britney Spears spelling co	llina 🛟	PageRank №	Check* KAutoLink \ AutoFill	
	Coogle • Dritney spears spe Britney Spears spelling co The data below shows some o that way. Each of these varia system (data for the correctly Detume to Concilia (data	llina ; rrrection of the misspellings detected b tions was entered by at least s spelled query is shown for c	y our spelling correction system two different unique users omparison).	Check* 《AutoLink 》 AutoFill stem for the query [britney spe within a three month period, a	













Associative thesauri

	Clustered thesaur which no distinction semantic relation; special case of as Associative thesaure represent words of relationship of set	ri can be buit automatically in the case on is made between different types of in this case, clustered thesauri are a sociative thesauri uri are graphs of words, where nodes and edges represent a (generic) mantic similarity between two words	in
	 Edges may be o Edges have an a strength of the The advantage is automatic way, sta semantic relation characteristics of occurrence (or complete the semanteristics) 	riented or not associated numerical weight w _{ij} (the e semantic association) that they may be built in a completely arting from a collection of docs. The between t _i and t _j mirrors the the collection, usually based on co- -absence) between t _i and t _j	
Dip. d	li Matematica F	- Aiolli - Information Retrieval	24





Abstract Indexing Theory

The weighting function $ffiif(f_r, i_s)$ can now be defined and weights can be further normalized by $w_{rs} = ffiif(f_r, i_s)/||ffiif(f_r, i_s)||^2$

Moreover SIM $(i_s, i_t) = w_r w_t$

If F is the set of terms and I is the set of docs we obtain document-document similarity for document search

If F is the set of docs and I is the set of terms we obtain term-term similarity for associative thesauri

Dip. di Matematica Pura ed Applicata F. Aiolli - Information Retrieval 2008/2009

29