



A definition of IR [Manning et.al. 2007]

Information retrieval (IR) is finding material (usually documents) of unstructured nature (usually text) that satisfy an information need from within large collections (usually on local computer servers or on the internet).

Dip. di Matematica Pura ed Applicata F. Aiolli - Information Retrieval 2009/2010

- 1. Finding material (documents) in *large* collections
- Unstructured nature as opposed to structured data: data do not have a fixed semantic/structure
- 3. Satisfy an *information need*

Dip. di Matematica Pura ed Applicata F. Aiolli - Information Retrieval 2009/2010

5



•.	
	A desire (possibly specified in an imprecise way) of information <mark>useful</mark> to the solution of the problem, or resources <mark>useful</mark> to a given goal;
	Useful (Relevant), according to the subjective opinion of the user.
Dip. (di Matematica F. Aiolli - Information Retrieval 9
Pura	ed Applicata 2009/2010
Pura	ed Applicata 2009/2010
Pura	ed Applicata 2009/2010 pical tasks covered in IR
Ti	ed Applicata 2009/2010 pical tasks covered in IR Search ('ad hoc' retrieval) Static document collection,
Ti D	ed Applicata 2009/2010 pical tasks covered in IR Search ('ad hoc' retrieval) Static document collection, Dynamic queries Filtering Static query, Dynamic document feeds
	ed Applicata 2009/2010 pical tasks covered in IR Search ('ad hoc' retrieval) Static document collection, Dynamic queries Filtering Static query, Dynamic document feeds Categorization Clustering
	ed Applicata 2009/2010 pical tasks covered in IR Search ('ad hoc' retrieval) Static document collection, Dynamic queries Filtering Static query, Dynamic document feeds Categorization Clustering Collaborative Filtering or Recommendation Browsing Summarization

	Text monolingual o multilingual text Structured text (XML) OCR->Text Spoken text typertext Ausic Graphics Emages /ideo / Animation
 	 monolingual o multilingual text Structured text (XML) OCR->Text Spoken text Hypertext Ausic Graphics Images /ideo / Animation
 	Structured text (XML) OCR->Text Spoken text Hypertext Music Graphics Emages /ideo / Animation
 	Spoken text Aypertext Ausic Graphics Emages Video / Animation
- - - - - - - - -	lypertext Ausic Graphics Cmages /ideo / Animation
- · · / - / - (-] -] -] 	Ausic Graphics Emages /ideo / Animation
- (-) -) -) -) - , - , - , - , - , - , - , - , - , - ,	Graphics Emages /ideo / Animation
]]] \] . 'ip. di ura ec	mages /ideo / Animation
ip. di	/ideo / Animation
ip. di	
ip. di ura ec	
ip. di ura ec	
ura ec	Matematica F. Aiolli - Information Retrieval 11
	a mature of TD
1 7 1	e nuture of IR
] [information Retrieval is difficult because of the
	nuererminancy of relevance:
	The system might interpret differently from the user
	The system might interpret differently from the user the meaning of the documents and/or the query, due to
	The system might interpret differently from the user the meaning of the documents and/or the query, due to inherent ambiguities in natural language;
	The system might interpret differently from the user the meaning of the documents and/or the query, due to inherent ambiguities in natural language; The user might not know exactly what she wants (vague or imprecise information need) and/or finds difficult to formalize it;
Th	e nature of IR

	Data Retrieval , as in DBs Information need cannot directly be expressed as a simple guery and documents
	have not a precise semantic. A translation of them into logical representations is needed.
	 In IR the set of objects to be retrieved are not clearly determined -> Slightly different retrieved sets should not be necessarily considered as a 'fatal' error of the system User satisfaction is the issue of IR
	 Knowledge Retrieval, as in AI In AI a fact α is inferred from a knowledge base Γ of facts expressed in a certain formalism
	 Question Answering In QA a query an answer is returned generated from a semantic analysis of documents. Huge amount of domain knowledge needed
	 Information Browsing, as in Hypermedia Relevant documents are retrieved by an active intervention of the user and not by a search routine
	The goal of a browsing task is less clear in the mind of the user
Dip Pura	. di Matematica F. Aiolli - Information Retrieval 13 a ed Applicata 2009/2010
E	volution of IR
E	volution of IR In the past, IR systems were used only by expert librarians as reference retrieval systems in batch modality.
E	volution of IR In the past, IR systems were used only by expert librarians as reference retrieval systems in batch modality.
E	volution of IR In the past, IR systems were used only by expert librarians as reference retrieval systems in batch modality. Many libraries still use categorization hierarchies to classify their volumes (e.g. Dewey Decimal Classification DDC) The advent of novel computers and the Web have brought to
	volution of IR In the past, IR systems were used only by expert librarians as reference retrieval systems in batch modality. • Many libraries still use categorization hierarchies to classify their volumes (e.g. Dewey Decimal Classification DDC) The advent of novel computers and the Web have brought to • efficient indexes, capable to index and displaying entire documents • processing of user queries with high performance

- methods to deal with multimedia
- interaction with the user
- methods to deal with distributed document collections (e.g. WWW)

	An IR model (can be defined by M=[D,Q,R] where	
	D is a repr	esentation for the documents in the	collection
	Q is a repr (queries)	resentation for the user information	needs
	R(d _i ,q _i) is a	a ranking function which associates	
	a real num represente	ber with a query $q_j \in Q$ and a document of $d_i \in D$.	ent
	N.B. It def a given	nes an ordering among the documents wi query q _j .	th regard to
		E Aiolli - Information Retrieval	15
Dip Pur	. di Matematica a ed Applicata	2009/2010	
Dip Pura	. di Matematica a ed Applicata	2009/2010	
Dip. Pura	Unlike query satisf documents D and in i.e. determined by	action in DBs, the relationship of <i>relevance</i> R be formation needs Q is not formally defined, but i the user. Therefore,	tween s subjective,
Dip.	. di Matematica a ed Applicata Unlike query satisf documents D and in i.e. determined by unlike in DBs, ef a degree of eff	action in DBs, the relationship of <i>relevance</i> R be formation needs Q is not formally defined, but i the user. Therefore, fectiveness is an issue ectiveness (user satisfaction) can be defined	tween S subjective,
Dip. Pura	Unlike query satisf documents D and in i.e. determined by unlike in DBs, ef a degree of effe In an IR system or R: D × Q → ℝ as	action in DBs, the relationship of <i>relevance</i> R be formation needs Q is not formally defined, but i the user. Therefore, fectiveness is an issue ectiveness (user satisfaction) can be defined model, it is necessary to choose whether to tree	tween s subjective, it relevance
Dip. Pura	Unlike query satisf documents D and in i.e. determined by unlike in DBs, ef a degree of effn In an IR system or R: D × Q → R as 1. Boolean-valued 2. Finite-valued R 3. Infinite-valued R	action in DBs, the relationship of <i>relevance</i> R be action in DBs, the relationship of <i>relevance</i> R be aformation needs Q is not formally defined, but i the user. Therefore, fectiveness is an issue activeness (user satisfaction) can be defined model, it is necessary to choose whether to treac $R \in \{0,1\}$ $R \in \{1,,N\}$ $R \in \mathbb{R}$	tween s subjective, it relevance



	Binary Retrieval RS	V(d _i ,q _j)∈ {0,1}	
	Does not allow th fact, for a given	e user to control the magnitude of t query, the system may return	he output. In
	under-dimensioover-dimensio	oned output ned output	
	Ranked Retrieval (c	rdering induced by $RSV(d_{i},q_{j}) \in \mathbb{R}$	2)
	 Allows the user t down until she se of users, those t those that want 	o start from the top of the ranked l es fit. This caters for the need of d hat want just few highly relevant do nany more.	ist and explore ifferent types cuments, and
Dip Pur	. di Matematica a ed Applicata	F. Aiolli - Information Retrieval 2009/2010	19
Dip Pur	di Matematica a ed Applicata	F. Aiolli - Information Retrieval 2009/2010	19 erates
	di Matematica a ed Applicata	F. Aiolli - Information Retrieval 2009/2010	19 erates
	di Matematica a ed Applicata	F. Aiolli - Information Retrieval 2009/2010 TRSYSTEM OP t representations and loading them i ad-hoc retrieval) representations and loading them ir filtering)	19 erates nto an internal nto an internal
	di Matematica a ed Applicata low a mode (off-line) building documen index structure (Or building query index structure ((on-line) [ad-hoc co	F. Aiolli - Information Retrieval 2009/2010 TR System op t representations and loading them i ad-hoc retrieval) representations and loading them ir filtering)	19 erates nto an internal nto an internal





Incidence Matrix

The result of the indexing process: the incidence matrix.

	d ₁	 d _i		d _m	
† ₁	w ₁₁	 w _{1i}		W _{1m}	
•••		 			
t _k	w _{k1}	 w _{ki}	•••	w _{km}	N. be
		 	•••	•••	
t _n	w _{n1}	 w _{ni}	•••	w _{nm}	

N.B. W_{ki} can either be binary or real.

 $T = \{t_{1,...}, t_{n}\}$ is the **dictionary** of the document base

Dip. di Matematica Pura ed Applicata F. Aiolli - Information Retrieval 2009/2010

24

	Weights can be assigned	
	 Manually (typically binary weights are used) by trained human indexers or intermediaries who are familiar with The discipline the documents deal with The indexing technique (e.g. the optimum number of terms for an IREP, the controlled vocabulary,) The contents of the collection (e.g. topic distribution) Automatically (either binary or real weights are used): by indexing process in the documents, in the docum	ne Sses
	query and in the collection.	IE
	Approach 2. is nowadays the only one left in <i>text retrieval</i> (cheaper of more effective).	and
	Approach 1. and 2. has been used in conjunction until recently becaus of the difficulty in producing effective automatic indexing technique for non-textual media.	e :S
D:		25
Pura	. di Matematica F. Aiolli - Information Retrieval a ed Applicata 2009/2010	25
Pura	ndexina – considerations	23
Dip. Pura	ndexing - considerations	25
	A di Matematica a ed Applicata F. Aiolli - Information Retrieval 2009/2010 Adexing - considerations The use of the same indexing technique for documents and query alike tends to guarantee correct matching process	25
	The use of the same indexing technique for documents and query alike tends to guarantee correct matching process There is indeterminacy in the indexing proces different indexers (human or automatic) do n produce in general the same IREP for the san document!	e a ss: ne



Recall: the "degree of completeness" of the system $ho = Pr(Ret|Rel) = rac{|\hat{Rel} \cap \hat{Ret}|}{|\hat{Rel}|}$

Dip. di Matematica Pura ed Applicata F. Aiolli - Information Retrieval 2009/2010



	In a ranked re relative to a re	trieval system, precision and recall are ank position r.	e values
	These systems function of re What is the recall has a	s can be evaluated by computing precises call, i.e. $\pi(\rho)$ precision $\pi(r(\rho))$ at the first rank position $r(\rho)$	sion as a b) for which
	We compute t relevant docur are interpolate	his function at each rank position in wl nent has been retrieved, and the resul ed yielding a <i>precision/recall plot</i>	hich a ting values
	A unique nume computing e.g.	rical value of the effectiveness can be the integral of precision as a function	e obtained by of recall
Dip. Pur:	. di Matematica a ed Applicata	F. Aiolli - Information Retrieval 2009/2010	31

Collection D, |D| = 100, a query q with |Rel|=20

Rel?	ρ	π(ρ)
У	1/20=0.05	1/1=1.00
У	2/20=0.10	2/2=1.00
Ν		
У	3/20=0.15	3/4=0.75
Ν		
Ν		
У	4/20=0.20	4/7=0.57
	•••	
	Rel? y y N y N y y	Rel?ρY1/20=0.05Y2/20=0.10N3/20=0.15N3/20=0.15N4/20=0.20Υ4/20=0.20

Dip. di Matematica Pura ed Applicata

F. Aiolli - Information Retrieval 2009/2010







	TREC web sit	te <pre>http://trec.nist.gov</pre>	
	The corpus: ' competitions	'Ad hoc' track in the first 8 T between '92 and '99.	REC
	Several millio	ons documents and 450 inform	nation needs
	In TREC, as i documents ai approximatin below	in many other big collections r re identified by a <i>data-pooling</i> ng the set of relevant documer	elevant 7 method thus 1ts from
Dip. Pura	di Matematica ed Applicata	F. Aiolli - Information Retrieval 2009/2010	39
		Relevant Known Re	levant

