	Binary Retrieval RSV(d _i ,q _j)∈ {0,1}	
	Does not allow the user to control the magnitude of the output fact, for a given query, the system may return	t. In
	 under-dimensioned output over-dimensioned output 	
	Ranked Retrieval (ordering induced by RSV(d_i,q_j) \in \mathbb{R})	
	Allows the user to start from the top of the ranked list and exdown until she sees fit. This caters for the need of different of users, those that want just few highly relevant documents, those that want many more.	kplore types and
Dip Pur	o. di Matematica F. Aiolli - Sistemi Informativi ra ed Applicata 2007/2008	17
Dip Pur	o. di Matematica ra ed Applicata F. Aiolli - Sistemi Informativi 2007/2008	17 P CS
	e. di Matematica ra ed Applicata F. Aiolli - Sistemi Informativi 2007/2008 How a modern IR system operat (off-line)	17 P es
	 Aiolli - Sistemi Informativi 2007/2008 How a modern IR system operat (off-line) building document representations and loading them into an int index structure; post-processing the retrieved documents 	17 ernal
	 di Matematica ra ed Applicata F. Aiolli - Sistemi Informativi 2007/2008 dow a modern IR system operat (off-line) building document representations and loading them into an int index structure; post-processing the retrieved documents (on-line) 	17 ernal





Incidence Matrix

The result of the indexing process: the incidence matrix.

	d ₁		di		d _m	
† ₁	w ₁₁	•••	w _{1i}	•••	W _{1m}	
•••	•••	•••	•••	•••	•••	
t _k	w _{k1}	•••	w _{ki}	•••	w _{km}	N.E 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 - Sistemi Informativi 2007/2008 22

 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, th controlled vocabulary) The controlled vocabulary) The controlled vocabulary) Automatically (either binary or real weights are used): by indexing process based on a statistical analysis of word occurrence in the documents, in the query and in the collection. Approach 2. is nowadays the only one left in <i>text retrieval</i> (cheaper or more effective). 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. Dip. di Matematica Pura ed Applicata F. Aiolli - Sistemi Informativi 2007/2008 The use of the same indexing technique for documents and query alike tends to guarantee correct matching process There is indeterminacy in the indexing process different indexers (human or automatic) do n produce in general the same IREP for the sam document! Unlike what happened in reference retrieval 	 Manually (typically binary weights are used) by trained human indexers on termediaries who are familiar with The discipline the documents deal with The indexing technique (e.g. the optimum number of terms for an IREP, th controlled vocabulary,) The contents of the collection (e.g. topic distribution) Automatically (either binary or real weights are used): by indexing proce based on a statistical analysis of word occurrence in the documents, in the user and in the collection. oach 2, is nowadays the only one left in text retrieval (cheaper state) 	r he sses
 Approach 2. is nowadays the only one left in <i>text retrieval</i> (cheaper of more effective). Approach 1. and 2. has been used in conjunction until recently becauss of the difficulty in producing effective automatic indexing technique for non-textual media. Dip. di Matematica Pura ed Applicata F. Aiolli - Sistemi Informativi 2007/2008 The use of the same indexing technique for documents and query alike tends to guarantee correct matching process There is indeterminacy in the indexing process different indexers (human or automatic) do n produce in general the same IREP for the sam document! Unlike what happened in reference retrieval 	oach 2. is nowadays the only one left in <i>text retrieval</i> (cheaper	he
 Approach 1. and 2. has been used in conjunction until recently because of the difficulty in producing effective automatic indexing technique for non-textual media. Dip. di Matematica Pura ed Applicata F. Aiolli - Sistemi Informativi 2007/2008 Indexing - considerations The use of the same indexing technique for documents and query alike tends to guarantee correct matching process There is indeterminacy in the indexing process different indexers (human or automatic) do n produce in general the same IREP for the same index ing recent what happened in reference retrieval 	ettective).	and
Dip. di Matematica F. Aiolli - Sistemi Informativi 2007/2008 Indexing - considerations The use of the same indexing technique for documents and query alike tends to guarantee correct matching process There is indeterminacy in the indexing process different indexers (human or automatic) do n produce in general the same IREP for the sam document! Unlike what happened in reference retrieval	oach 1. and 2. has been used in conjunction until recently becaus le difficulty in producing effective automatic indexing technique lon-textual media.	3e 25
 Pura ed Applicata 2007/2008 Indexing - considerations The use of the same indexing technique for documents and query alike tends to guarantee correct matching process There is indeterminacy in the indexing process different indexers (human or automatic) do n produce in general the same IREP for the sam document! Unlike what happened in reference retrieval 	matica F. Aiolli - Sistemi Informativi	23
 The use of the same indexing technique for documents and query alike tends to guarantee correct matching process There is indeterminacy in the indexing process different indexers (human or automatic) do n produce in general the same IREP for the sam document! Unlike what happened in reference retrieval 	exina - considerations	
 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 sam document! Unlike what happened in reference retrieval 		
 There is indeterminacy in the indexing proces different indexers (human or automatic) do n produce in general the same IREP for the sam document! Unlike what happened in reference retrieval 	e use of the same indexing technique for suments and query alike tends to guaranted rect matching process	ea
Unlike what happened in reference retrieval	ere is indeterminacy in the indexing proces ferent indexers (human or automatic) do r duce in general the same IREP for the sar sument!	ss: 10t ne
systems, the on-line availability of the entire document allows the use of the entire document also for indexing	ike what happened in reference retrieval	



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 - Sistemi Informativi 2007/2008





Collection D, C	= 100, a query q	with Rel =20
------------------	------------------	---------------

					_
	Rank Pos. q	Rel?	ρ	π(ρ)	
	1	У	1/20=0.05	1/1=1.00	
	2	У	2/20=0.10	2/2=1.00	
	3	N			
	4	У	3/20=0.15	3/4=0.75	
	5	N			
	6	N			
	7	У	4/20=0.20	4/7=0.57	
				•••	
. d	Matematica	F. Aiolli	- Sistemi Informativi		

Dip Pura ed Applicata 2007/2008

Note that

- □ The effectiveness of a system is typically evaluated by *averaging* over different queries (*macroaveraging*)
 - Different searchers are equally important
 - Partial view of a problem: different methods may work best for different types of queries
- A typical precision/recall plot
 - Is monotonically decreasing
 - For $\rho=1$ it takes the value $\pi=q$, where q=|Rel|/D is the generality (frequency) of the query
- □ When the document base is big, it is very important to have high precisions for small recall values.
 - Measures such as precision at 10 (P@10) are often used in place of $\pi(\rho)$
- 'Typical' values are not beyond .4 precision at .4 recall





	C web site	<pre>http://trec.nist.gov</pre>	
□ The comp	corpus: 'A Detitions	Ad hoc' track in the first 8 TF between '92 and '99.	REC
] Seve	eral millio	ns documents and 450 inform	ation needs
□ In T docu appr belov	REC, as ir ments ar oximating N	n many other big collections ro e identified by a <i>data-pooling</i> g the set of relevant documen	elevant method thus its from
Dip. di Maten Pura ed Appli	natica cata	F. Aiolli - Sistemi Informativi 2007/2008	37
Dip. di Maten Pura ed Appli	natica cata	F. Aiolli - Sistemi Informativi 2007/2008	37
Dip. di Maten Pura ed Appli	natica cata	F. Aiolli - Sistemi Informativi 2007/2008 Relevant Known	37

