





Why Information Retrieval:

Information Retrieval (IR) mainly studies unstructured data:

Text in Web pages or emails; image; audio; video; protein sequences..

Merrill Lynch estimates that more than 85 percent of all business information exists as unstructured data - commonly appearing in emails, memos, notes from call centers and support operations, news, user groups, chats, reports, ... and Web pages.

Unstructured data:

No structure: no primary key as in RDBMS

Semantic meaning unknown: natural language processing systems try to find the meaning in the unstructured text











Sulling S	ome core	concepts of IR
Google	Web Images Groups Ne information retrieval system	wes Froogle Maps Scholar more »
Web	/	Results 1 - 10 of about 75,600,000
Intelligent Information bit.csc.lsu.edu/~kraft/ro Past Performance I Welcome to the Past P	Retrieval · Callan CMU IR Group . etrieval.html - 38k - <u>Cached</u> - <u>Simi</u> nformation Retrieval Syste erformance Information Retrieval	n lar pages M System (PPIRS). PPIRS is a web-
enabled, government-wi www.ppirs.gow/ - 10k - (de application that provides timely ≳ached - <u>Similar pages</u>	and pertinent Text Summarizations
Electronic Digital In EDIS is the Electronic I information on topics edis.ifas.ufl.edu/ - 26k -	formation Source (EDIS) Data Information Source of UF/IF/ relevant to you: profitable and susta <u>Cached</u> - <u>Similar pages</u>	AS Extension, a collection of ainable
Information retriev Automated information explosion in scientific li en.wikipedia.org/wiki/Inf	val - Wikipedia, the free encyon retrieval (IR) systems were origin terature in the last few decades ormation retrieval - 44k - Jun 21.	clopedia nally used to manage information 2006 - Cached - Similar pages









































Text Representation: What computer sees

 Reviewer:

"dage456" (Carmichael, CA USA) -

See all my reviewsSee all my reviewsa deflated wallet (wonder where the money went). I have had my ipod now for 4months and cannot imagine how I used to get by with my old rio 600 with is 64 megs ofram and.. usb connection. Because of its size this little machine goes with myeverywhere and its ten hour battery life means I can listen to stuff all day long.Pros:size, both physical and capacity.
size, bots beautiful
br><connection:</td>FIREWIRE!!Cons: needs the ability to bookmark. Iuse my ipod mostly for audiobooks. the ipod needs to include a bookmark feature forthose like me.
> />

From Amazon Customer Review of IPod

















Text Representation: Indexing Statistical Properties of Text

Word	Frequency	Word	Frequency	
the	1130021	market	52110	
of	547311	bank	47940	
to	516636	stock	47401	
а	464736	trade	47310	
in	390819			
and	387703			

Statistics collected from Wall Street Journal (WSJ), 1987





C. C	Preprocessing				
4 the	1 at	1 different	1 may	1 step	
3 and	1 basal	1 exchange	1 nontarget	1 substance	
3 by	1 be	1 exogenous	1 not	1 suggests	
3 steroids	1 been	1 fluorescent	1 may	1 target	
2 centrioles	1 bodies	1 from	1 of	1 technique	
2 in	1 can	1 growth	1 precise	1 two	
1 affect	1 at	1 has	1 receptor	1 unexpected	
1 already	1 cell	1 identity	1 regularly	1 vitally	
1 Although	1 cells	1 level	1 reveal	1 way	
1 antibodies	1 cilia-bearing	1 localization	1 Specific	1 with	











Text Representation: Inverted Lists					<u> </u>	
Doc ID	Text					
1	kids question noting in 1960s					
2	young man question everything in	estion everything in 1970s				
3	kids question questions in 1980s					
4	young man question nothing in	uestion nothing in 2000s				
		Term 1 2 3 4 5 6 7 8 9 10	ID	Term kids question nothing in 19060s young man everything 1970s questions	Documents 1,3 1,2,3,4 1,4 1,2,3,4 1 2,4 2 2 3	
	Inverted Lists	11		1980s 2000s	3 4	



































of the stand	Retrieval Models: Vector Space Model						
Vector	Vector representation						
Γ		Java	Sun	Starbucks			
	D1	1	1	0			
	D2	1	0	1			
	D3	1	0	0			
	Query	1	0.2	1]		
	Similarity Score	D1	D2	D3			
	Query	0.59	0.99	0.70			



