































TEXT TRANSFORMATION

Stopping

- Remove common words (stop words)
 - e.g., "and", "or", "the", "in"
- Some impact on efficiency and effectiveness
- Can be a problem for some queries
- Stemming
 - Group words derived from a common stem
 - e.g., "computer", "computers", "computing", "compute"
 - Usually effective, but not for all queries
 - Benefits vary for different languages



- Link Analysis
 - Makes use of *links* and *anchor text* in web pages
 - Link analysis identifies *popularity* and *community* information
 - e.g., PageRank
 - Anchor text can significantly enhance the representation of pages pointed to by links
 - Significant impact on web search
 - Less importance in other applications

TEXT TRANSFORMATION

- Information Extraction
 - Identify classes of index terms that are important for some applications
 - e.g., *named entity recognizers* identify classes such as *people*, *locations*, *companies*, *dates*, etc.
- Classifier
 - Identifies class-related metadata for documents
 - i.e., assigns labels to documents
 - e.g., topics, reading levels, sentiment, genre
 - Use depends on application



<section-header><section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item>









USER INTERACTION

- Results output
 - Constructs the display of ranked documents for a query
 - Generates *snippets* to show how queries match documents
 - Highlights important words and passages
 - Retrieves appropriate *advertising* in many applications
 - May provide *clustering* and other visualization tools







Logging

- Logging user queries and interaction is crucial for improving search effectiveness and efficiency
- Query logs and clickthrough data used for query suggestion, spell checking, query caching, ranking, advertising search, and other components

(29)

- Ranking analysis
 - Measuring and tuning ranking effectiveness
- Performance analysis
 - Measuring and tuning system efficiency









	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0 🔪	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0
Duriture			1	if pla	v cont	ains
ыни	AND Caesar E	SUTINOT		uard (otho	nuico



<section-header><section-header><section-header><text><text><text><text><text>

































SCHEMA FOR INDEX CONSTRUCTION IN MAPREDUCE

Schema of map and reduce functions

• map: input \rightarrow list(k, v) reduce: (k,list(v)) \rightarrow output

Instantiation of the schema for index construction

- map: web collection \rightarrow list(termID, docID)
- reduce: (<termID1, list(docID)>, <termID2, list(docID)>, ...) → (postings list1, postings list2, ...)

Example for index construction

- map: d2 : C died. d1 : C came, C c'ed. → (<C, d2>, <died,d2>, <C,d1>, <came,d1>, <C,d1>, <c'ed, d1>
- reduce: (<C,(d2,d1,d1)>, <died,(d2)>, <came,(d1)>, <c'ed,(d1)>)
 → (<C,(d1:2,d2:1)>, <died,(d2:1)>, <came,(d1:1)>,
 <c'ed,(d1:1)>)



















POSITIONAL INDEX SIZE

- Need an entry for each occurrence, not just once per document

- Index size depends on average document size
 - Average web page has <1000 terms
 - SEC filings, PDF files, ... easily 100,000 terms
- Consider a term with frequency 0.1% in a doc

Document size	Postings	Positional postings
1000	1	1
100,000	1	100

63

















ELIAS-Г C • To encode a number	CO <i>k</i> , co	DE	• $k_d = \lfloor \log_2 k \rfloor$ • $k_r = k - 2^{\lfloor \log_2 k \rfloor}$
 <i>k_d</i> is number of bina 	ry digi	ts, enco	oded in unary
Number (k)	k_d	k_r	Code
1	0	0	0
2	1	0	10 0
3	1	1	10 1
6	2	2	110 10
15	3	7	1110 111
16	4	0	11110 0000
255	7	127	11111110 1111111
1023	9	511	1111111110 111111111
			(12)



ELIA2-7 COI	Æ					
Split k _d into:		• k_{de}	$d = \lfloor$	\log_2	(k_d+1)	
		• k_{di}	r = k	$c_d - 2$	$2\lfloor \log_2(k_d+1) \rfloor$	
• encode k _{dd} in una	ary, <i>k</i>	_{dr} in bi	nary, a	and k _r	in binary	
Number (k)	k_d	k_r	k_{dd}	k_{dr}	Code	
1	0	0	0	0	0	
2	1	0	1	0	10 0 0	
3	1	1	1	0	10 0 1	
6	2	2	1	1	$10\ 1\ 10$	
15	3	7	2	0	$110 \ 00 \ 111$	
16	4	0	2	1	$110 \ 01 \ 0000$	
055	7	127	3	0	$1110\ 000\ 1111111$	
255		F11	2	2	1110 010 111111111	



	k	Number of	bytes	
	$k < 2^7$	1		
	$2^{\prime} \leq k < 2^{14}$	$\begin{vmatrix} 2\\ 2 \end{vmatrix}$		
	$2^{-1} \le \kappa < 2^{-1}$ $2^{21} \le k < 2^{28}$	3 4		
k	E	Binary Code	Hexadecimal	_
$\frac{k}{1}$	F	Binary Code 1 0000001	Hexadecimal 81	_
$\frac{k}{1}$	E	Binary Code 1 0000001 1 0000110 1 1111111	Hexadecimal 81 86 FF	_
$\frac{k}{1}$ 6 127 128	E	Binary Code 1 0000001 1 0000110 1 111111 1 1 0000000	Hexadecimal 81 86 FF 01 80	_
$ \frac{k}{1} \\ 6 \\ 127 \\ 128 \\ 130 $	E 0 000000 0 000000	Binary Code 1 0000001 1 0000110 1 111111 1 1 0000000 1 1 0000000 1 1 0000010	Hexadecimal 81 86 FF 01 80 01 82	_



















EXAMPLE: WESTLAW HTTP://WWW.WESTLAW.COM/

- Largest commercial (paying subscribers) legal search service (started 1975; ranking added 1992)
- Tens of terabytes of data; 700,000 users
- Majority of users still use boolean queries
- Example query:
 - What is the statute of limitations in cases involving the federal tort claims act?
 - LIMIT! /3 STATUTE ACTION /S FEDERAL /2 TORT /3 CLAIM
 - /3 = within 3 words, /S = in same sentence

<section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item>











TE	RM-DOCUME	NT COU	NT MAI	RICE	S	
• Con • Ea	sider the number ach document is a <mark>co</mark>	of occurre unt vector in l	nces of a te ℕ ^v : a column	erm in a below	docum	ent:
	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	157	73	0	0	0	0
Brutus	4	157	0	1	0	0
Caesar	232	227	0	2	1	1
Calpurnia	0	10	0	0	0	0
Cleopatra	57	0	0	0	0	0
mercy	2	0	3	5	5	1
worser	2	0	1	1	1	0
			1			94)









DOCUMENT FREQUENCY, CONTINUED

- Frequent terms are less informative than rare terms
- Consider a query term that is frequent in the collection (e.g., high, increase, line)
- A document containing such a term is more likely to be relevant than a document that doesn't
- But it's not a sure indicator of relevance.
- → For frequent terms, we want high positive weights for words like high, increase, and line
- But lower weights than for rare terms.
- We will use document frequency (df) to capture this.









• The tf-idf weight of a term is the product of its tf weight and its idf weight.

$$\mathbf{W}_{t,d} = (1 + \log \mathrm{tf}_{t,d}) \times \log_{10}(N/\mathrm{df}_t)$$

- Best known weighting scheme in information retrieval
 Note: the "-" in tf-idf is a hyphen, not a minus sign!
 - Alternative names: tf.idf, tf x idf
- Increases with the number of occurrences within a document
- Increases with the rarity of the term in the collection



Antony5.25Brutus1.21Caesar8.59Calpurnia0	3.18 6.1	0	0	0	0 35
Brutus1.21Caesar8.59Calpurnia0	6.1	0			0.00
Caesar 8.59 Calpurnia 0		Ū	1	0	0
Calpurnia 0	2.54	0	1.51	0.25	0
	1.54	0	0	0	0
Cleopatra 2.85	0	0	0	0	0
mercy 1.51	0	1.9	0.12	5.25	0.88
worser 1.37	0	0.11	4.15	0.25	1.95













COSINE SIMILA	RITY AN	IONGST	3 DOCU	MENTS
the novels	term	SaS	PaP	WH
SaS: Sense and	affection	115	58	20
Sensibility	jealous	10	7	11
PaP: Pride and	gossip	2	0	6
<i>Prejudice,</i> and WH: <i>Wuthering</i>	wuthering	0	0	38
Heights?	Terr	n freque	encies (c	ounts)
Note: To simplify th	nis examp	le, we do	n't do idf	weightin





