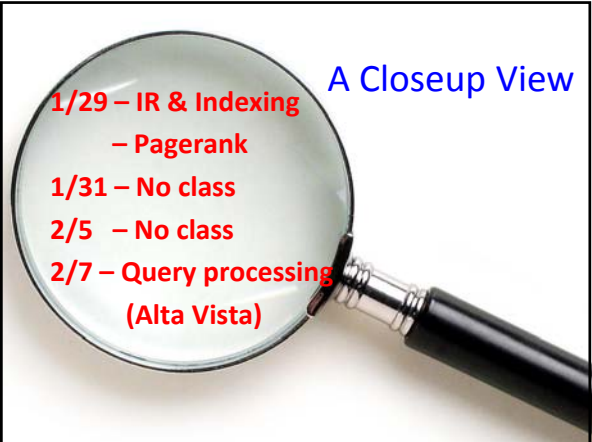


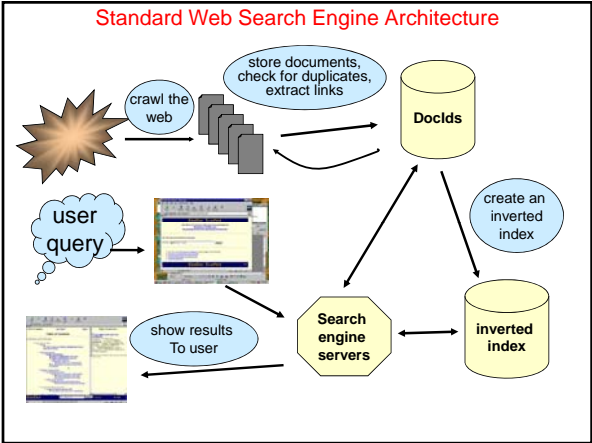
**CSE 454**

**Information Retrieval & Indexing**

**A Closeup View**



1/29 – IR & Indexing  
 – Pagerank  
 1/31 – No class  
 2/5 – No class  
 2/7 – Query processing  
 (Alta Vista)



- Relevance**
- **Complex concept – extensive study**
    - Less consensus
    - People often disagree on relevance
    - Many factors...
  - **Retrieval models make various assumptions about relevance to simplify problem**
    - e.g., *topical* vs. *user* relevance
    - e.g., *binary* vs. *multi-valued* relevance
- from Croft, Metzler, Strohan. © Addison Wesley

- Retrieval Model Overview**
- **Older models**
    - Boolean retrieval
    - Overlap Measures
    - Vector Space model
  - **Probabilistic Models**
    - BM25
    - Language models
  - **Combining evidence**
    - Inference networks
    - Learning to Rank
- from Croft, Metzler, Strohan. © Addison Wesley

**Test Corpora**

TABLE 4.3 Common Test Corpora

Collection	NDocs	NQrys	Size (MB)	Term/Doc	Q-D RelAss
ADI	82	35			
AIT	2109	14	2	400	>10,000
CACM	3204	64	2	24.5	
CISI	1460	112	2	46.5	
Crowfield	1400	225	2	53.1	
LISA	5872	35	3		
Medline	1033	30	1		
NPL	11,429	93	3		
OSHMED	34,8566	106	400	250	16,140
Reuters	21,578	672	28	131	
TREC	740,000	200	2000	89-3543	≈ 100,000

slide from Raghavan, Schütze, Larson

## Standard Benchmarks

- **National Institute of Standards +Testing (NIST)**
  - Has run large IR testbed for many years (TREC)
- **Reuters and other benchmark sets used**
- **“Retrieval tasks” specified**
  - sometimes as queries
- **Human experts mark, for each query and for each doc, “Relevant” or “Not relevant”**
  - or at least for subset that some system returned

slide from Raghavan, Schütze, Larson

## Precision & Recall

### Precision

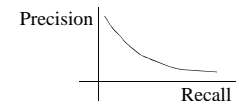
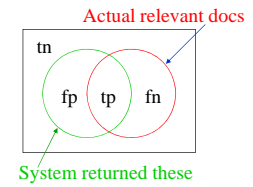
$\frac{tp}{tp + fp}$   
Proportion of selected items that are correct

### Recall

$\frac{tp}{tp + fn}$   
% of target items that were selected

### Precision-Recall curve

Shows tradeoff



## Boolean Retrieval

- **Advantages**
  - Results are predictable, relatively easy to explain
  - Many different features can be incorporated
  - Efficient processing since many documents can be eliminated from search
- **Disadvantages**
  - Effectiveness depends entirely on user
  - Simple queries usually don't work well
  - Complex queries are difficult
  - Brittle with user errors (eg misspelling)

from Croft, Metzler, Strohan. © Addison Wesley

## Interlude

- **Better Models Coming Soon:**
  - Vector Space model
  - Probabilistic Models
    - BM25
    - Language models
- **Shared Issues – What to Index**
  - Punctuation
  - Case Folding
  - Stemming
  - Stop Words
  - Numbers
  - Font size, titles, anchor text

## Punctuation

- **Ne'er:** use language-specific, handcrafted “locale” to normalize.
- **State-of-the-art:** break up hyphenated sequence.
- **U.S.A. vs. USA** - use locale.
- **a.out**

slide from Raghavan, Schütze, Larson

## Numbers

- **3/12/91**
- **Mar. 12, 1991**
- **55 B.C.**
- **B-52**
- **100.2.86.144**
  - Generally, don't index as text
  - Creation dates for docs

slide from Raghavan, Schütze, Larson

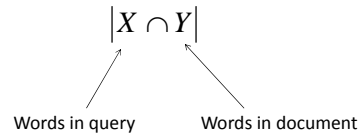
## Case folding

- Reduce all letters to lower case
- Exception: upper case in mid-sentence
  - e.g., *General Motors*
  - *Fed* vs. *fed*
  - *SAIL* vs. *sail*

slide from Raghavan, Schütze, Larson

## Ranking search results

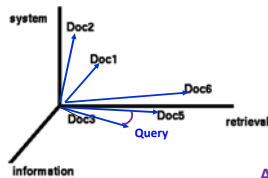
- Boolean queries give inclusion or exclusion of docs.
- Need to assess *quality* of results
  - First attempt: **OVERLAP** between query and document



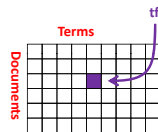
- What's missing?

## Vector Space Model

- Each term defines an axis
  - Even with stemming, may have 20,000+ dimensions
- Each doc is a vector of *frequency* values
  - One "tf" component for each term
- Treat query as just another document
  - How measure *distance* in this vector space??



Equivalently



Are all dimensions equivalent?

## TF x IDF

$$w_{ik} = tf_{ik} * \log(N / n_k)$$

$T_k$  = term  $k$  in document  $D_i$

$tf_{ik}$  = frequency of term  $T_k$  in document  $D_i$

$idf_k$  = inverse document frequency of term  $T_k$  in  $C$

$$idf_k = \log\left(\frac{N}{n_k}\right)$$

$N$  = total number of documents in the collection  $C$

$n_k$  = the number of documents in  $C$  that contain  $T_k$

slide from Raghavan, Schütze, Larson

## BM25

Popular and effective ranking algorithm based on binary independence model

– adds document and query term weights

$$\sum_{i \in Q} \log \frac{(r_i + 0.5) / (R - r_i + 0.5)}{(n_i - r_i + 0.5) / (N - n_i - R + r_i + 0.5)} \cdot \frac{(k_1 + 1) f_i}{K + f_i} \cdot \frac{(k_2 + 1) q f_i}{k_2 + q f_i}$$

- $N$  = number of doc,  $n_i$  = num containing term  $i$
- $R, r_i$  = encode relevance info (if avail, otherwise = 0)
- $f_i$  = freq of term  $i$  in doc;  $q f_i$  = freq in doc
- $k_1, k_2$  and  $K$  are parameters, values set empirically
  - $k_1$  weights  $tf$  component as  $f_i$  increases
  - $k_2$  = weights query term weight
  - $K$  normalizes

adapted from Croft, Metzler, Strohan. © Addison Wesley

## Simple Formulas

But How Process Efficiently?

Copyright © Weld 2002-2007

42

## Thinking about Efficiency

- **Clock cycle: 4 GHz**
  - Typically *completes* 2 instructions / cycle
    - ~10 cycles / instruction, but pipelining & parallel execution
  - Thus: 8 billion instructions / sec
- **Disk access: 1-10ms**
  - Depends on seek distance, published average is 5ms
  - Thus perform 200 seeks / sec
  - (Ignoring rotation and transfer times)
- **Disk is 40 Million times slower !!!**

Copyright © Weld 2002-2007

43

## Retrieval

### Document-term matrix

	$t_1$	$t_2$	$\dots$	$t_j$	$\dots$	$t_m$	$nf$
$d_1$	$w_{11}$	$w_{12}$	$\dots$	$w_{1j}$	$\dots$	$w_{1m}$	$1/ d_1 $
$d_2$	$w_{21}$	$w_{22}$	$\dots$	$w_{2j}$	$\dots$	$w_{2m}$	$1/ d_2 $
$\vdots$	$\vdots$	$\vdots$	$\dots$	$\vdots$	$\dots$	$\vdots$	$\vdots$
$d_i$	$w_{i1}$	$w_{i2}$	$\dots$	$w_{ij}$	$\dots$	$w_{im}$	$1/ d_i $
$\vdots$	$\vdots$	$\vdots$	$\dots$	$\vdots$	$\dots$	$\vdots$	$\vdots$
$d_n$	$w_{n1}$	$w_{n2}$	$\dots$	$w_{nj}$	$\dots$	$w_{nm}$	$1/ d_n $

$w_{ij}$  is the weight of term  $t_j$  in document  $d_i$

Most  $w_{ij}$ 's will be zero.

Copyright © Weld 2002-2007

44

## Naïve Retrieval

Consider query  $Q = (q_1, q_2, \dots, q_j, \dots, q_n)$ ,  $nf = 1/|q|$ .

How evaluate  $Q$ ?

(i.e., compute the similarity between  $q$  and every document)?

**Method 1: Compare  $Q$  with every doc.**

Document data structure:

$d_i : ((t_1, w_{i1}), (t_2, w_{i2}), \dots, (t_j, w_{ij}), \dots, (t_m, w_{im}), 1/|d_i|)$

- Only terms with positive weights are kept.
- Terms are in alphabetic order.

Query data structure:

$Q : ((t_1, q_1), (t_2, q_2), \dots, (t_j, q_j), \dots, (t_m, q_m), 1/|q|)$

Copyright © Weld 2002-2007

45

## Observation

- **Method 1 is not efficient**
  - Needs to access most non-zero entries in doc-term matrix.
- **Solution: Use Index (Inverted File)**
  - Data structure to permit fast searching.
- **Like an Index in the back of a text book.**
  - Key words --- page numbers.
  - E.g, "Etzioni, 40, 55, 60-63, 89, 220"
- **Lexicon**
- **Occurrences**

Copyright © Weld 2002-2007

47

## Search Processing (Overview)

1. **Lexicon search**
  - E.g. looking in index to find entry
2. **Retrieval of occurrences**
  - Seeing where term occurs
3. **Manipulation of occurrences**
  - Going to the right page

Copyright © Weld 2002-2007

48

## Simple Index for One Document FILE

POS	
1	A file is a list of words by position
10	First entry is the word in position 1 (first word)
20	Entry 4562 is the word in position 4562 (4562 <sup>nd</sup> word)
30	Last entry is the last word
36	An inverted file is a list of positions by word!

a (1, 4, 40)  
 entry (11, 20, 31)  
 file (2, 38)  
 list (5, 41)  
 position (9, 16, 26)  
 positions (44)  
 word (14, 19, 24, 29, 35, 45)  
 words (7)  
 4562 (21, 27)

INVERTED FILE

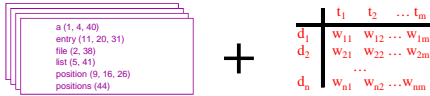
aka "Index"

Copyright © Weld 2002-2007

49

## Requirements for Search

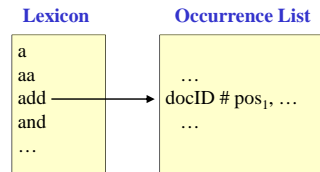
- **Need index structure**
  - Must handle multiple documents
  - Must support phrase queries
  - Must encode TF/IDF values
  - Must minimize disk seeks & reads



Copyright © Weld 2002-2007

50

## How Store Index?

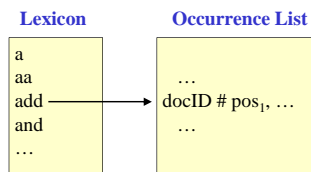


Oracle Database?

Unix File System?

## The Solution

- **Inverted Files for Multiple Documents**
  - Broken into Two Files
- **Lexicon**
  - Hashtable on disk (one read)
  - Nowadays: stored in main memory
- **Occurrence List**
  - Stored on Disk
  - “Google Filesystem”



Copyright © Weld 2002-2007

53

## Inverted Files for Multiple Documents

WORD	NDOCS	PTR	DOCID	OCCUR	POS 1	POS 2	...
jezebel	20		34	6	1	118	2087 3922 3981 5002
jezer	3		44	3	215	2291	3010
jezerit	1		56	4	5	22	134 992 ...
jeziah	1		566	3	203	245	287
jeziel	1						
jeziah	1		67	1	132		
jezoar	1						
jezrahiah	1						
jezreel	39		107	4	322	354 381 405	
			232	6	15	195 248	1897 1951 2192
			677	1	481		
			713	3	42	312	802

**OCURRENCE INDEX**

*Note: "jezebel" occurs 6 times in document 34, 3 times in document 44, 4 times in document 56, ...*

- One method. Alta Vista uses alternative

Copyright © Weld 2002-2007

54

## Many Variations Possible

- **Address space (flat, hierarchical)**
- **Record term-position information**
- **Precalculate TF-IDF info**
- **Stored header, font & tag info**
- **Compression strategies**

Copyright © Weld 2002-2007

55

## Other Features Stored in Index

- **Page Rank**
- **Query word in color on page?**
- **# images on page**
- **# outlinks on page**
- **URL length**
- **Page edit recency**
- **Page Classifiers (20+)**
  - Spam
  - Adult
  - Actor
  - Celebrity
  - Athlete
  - Product / review
  - Tech company
  - Church
  - Homepage
  - ....

Amit Singhai says Google uses over 200 such features [NY Times 2008-06-03]

## Using Inverted Files

Some data structures:

Lexicon: a hash table for all terms in the collection.

	.....
$t_j$	pointer to $I(t_j)$
	.....

- Inverted file lists previously stored on disk.
- Now fit in main memory

Copyright © Weld 2002-2007

57

## The Lexicon

- Grows Slowly (Heap's law)
  - $O(n^\beta)$  where  $n$ =text size;  $\beta$  is constant  $\sim 0.4 - 0.6$
  - E.g. for 1GB corpus, lexicon = 5Mb
  - Can reduce with stemming (Porter algorithm)
- Store lexicon in file in lexicographic order
  - Each entry points to loc in occurrence file (aka inverted file list)

Copyright © Weld 2002-2007

58

## Using Inverted Files

Several data structures:

2. For each term  $t_j$ , create a list (occurrence file list) that contains all document ids that have  $t_j$ .

$$I(t_j) = \{ (d_1, w_{1j}), \\ (d_2, \dots \\ \dots \}$$

- $d_i$  is the document id number of the  $i^{\text{th}}$  document.
- Weights come from freq of term in doc
- Only entries with non-zero weights are kept.

Copyright © Weld 2002-2007

59

## More Elaborate Inverted File

Several data structures:

2. For each term  $t_j$ , create a list (occurrence file list) that contains all document ids that have  $t_j$ .

$$I(t_j) = \{ (d_1, \text{freq}, \text{pos}_1, \dots, \text{pos}_k), \\ (d_2, \dots \\ \dots \}$$

- $d_i$  is the document id number of the  $i^{\text{th}}$  document.
- Weights come from freq of term in doc
- Only entries with non-zero weights are kept.

Copyright © Weld 2002-2007

60

## Inverted files continued

More data structures:

3. Normalization factors of documents are pre-computed and stored similarly to lexicon

$nf[i]$  stores  $1/|d_i|$ .

Copyright © Weld 2002-2007

61

## Retrieval Using Inverted Files

initialize all  $\text{sim}(q, d_i) = 0$

for each term  $t_j$  in  $q$

find  $I(t)$  using the hash table

for each  $(d_i, w_{ij})$  in  $I(t)$

$$\text{sim}(q, d_i) += q_j * w_{ij}$$

for each (relevant) document  $d_i$

$$\text{sim}(q, d_i) = \text{sim}(q, d_i) * nf[i]$$

sort documents in descending similarities and display the top  $k$  to the user;

Copyright © Weld 2002-2007

62

## Observations about Method 2

- If doc  $d$  **doesn't contain** any term of query  $q$ , then  $d$  **won't be considered** when evaluating  $q$ .
- Only **non-zero** entries in the columns of the document-term matrix which correspond to query terms ... are used to evaluate the query.
- Computes the similarities of multiple documents simultaneously (w.r.t. each query word)

Copyright © Weld 2002-2007

63

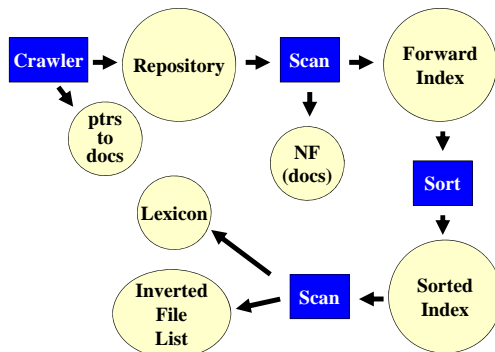
## Efficiency versus Flexibility

- Storing computed document weights is good for efficiency, but bad for flexibility.
  - **Recomputation needed if TF and IDF formulas change and/or TF and DF information changes.**
- Flexibility improved by storing raw TF, DF information, but efficiency suffers.
- A compromise
  - **Store pre-computed TF weights of documents.**
  - **Use IDF weights with query term TF weights instead of document term TF weights.**

Copyright © Weld 2002-2007

66

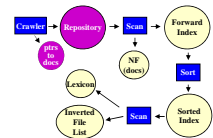
## How Inverted Files are Created



Copyright © Weld 2002-2007

67

## Creating Inverted Files



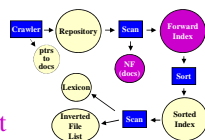
### Repository

- File containing all documents downloaded
- Each doc has unique ID
- Ptr file maps from IDs to start of doc in repository

Copyright © Weld 2002-2007

68

## Creating Inverted Files



NF ~ Length of each document

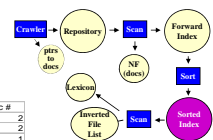
### Forward Index

Term	Doc #	Pos
i	1	1
did	1	2
eract	1	3
julus	1	4
caesar	1	5
i	1	6
was	1	7

Copyright © Weld 2002-2007

69

## Creating Inverted Files



### Sorted Index

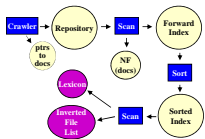
(positional info as well)

Term	Doc #	Term	Doc #
i	1	ambitious	2
did	1	be	2
eract	1	brutus	1
julus	1	capitol	1
caesar	1	caesar	1
i	1	caesar	2
was	1	caesar	2
killed	1	did	1
i	1	eract	1
the	1	hath	1
capitol	1	i	1
brutus	1	killed	1
i	1	me	1
killed	1		
me	1		

Copyright © Weld 2002-2007

70

## Creating Inverted Files



### Lexicon

WORD	NDOCS	PTR	DOCID	OCCUR	POS 1	POS 2	...
jezebel	20		34	6	1	118	2087 3922 3981 5002
jezer	3		44	3	215	2291	3010
jezerit	1		56	4	5	22	134 992
jeziah	1		566	3	203	245	287
jeziel	1						
jeziah	1		67	1	132		
jezoar	1						
jezahiah	1						
jezreel	39						

### Inverted File List

Copyright © Weld 2002-2007

71

## Compression

### • What Should We Compress?

- Repository
- Lexicon
- Inv Index

### • What properties do we want?

- Compression ratio
- Compression speed
- Decompression speed
- Memory requirements
- Pattern matching on compressed text
- Random access

Copyright © Weld 2002-2007

77

## Inverted File Compression

Each inverted list has the form  $\langle f_i; d_1, d_2, d_3, \dots, d_{f_i} \rangle$

A naïve representation results in a storage overhead of  $(f + n) * \lceil \log N \rceil$

This can also be stored as  $\langle f_i; d_1, d_2 - d_1, \dots, d_{f_i} - d_{f_i-1} \rangle$

Each difference is called a **d-gap**. Since  $\sum(d - \text{gaps}) \leq N$ ,

each pointer requires fewer than  $\lceil \log N \rceil$  bits.

Trick is encoding .... since worst case ....

➡ Assume **d-gap representation for the rest of the talk, unless stated otherwise**

Slides adapted from Tapas Kanungo and David Mount, Univ Maryland

Copyright © Weld 2002-2007

78

## Text Compression

Two classes of text compression methods

### • Symbolwise (or statistical) methods

- Estimate probabilities of symbols - modeling step
- Code one symbol at a time - coding step
- Use shorter code for the most likely symbol
- Usually based on either arithmetic or Huffman coding

### • Dictionary methods

- Replace fragments of text with a single code word
- Typically an index to an entry in the dictionary.
  - eg: Ziv-Lempel coding: replaces strings of characters with a pointer to a previous occurrence of the string.
- No probability estimates needed

➡ **Symbolwise methods are more suited for coding d-gaps**

Copyright © Weld 2002-2007

79