

## Search Engines (Information Retrieval)

## Motivation and Outline

- **Background**
  - Definitions, etc.
- **The Problem**
  - 100,000+ pages
- **The Solution**
  - Indexing docs
  - Ranking results
- **Extensions**
  - Relevance feedback, **clustering**, **query expansion**, etc.

## Search Engine Architecture

- **Spider**
  - Crawls the web to find pages. Follows hyperlinks. Never stops
- **Indexer**
  - Produces data structures for fast searching of all words in the pages
- **Retriever**
  - Query interface
  - Ranking (!!)

## Indexing

- Arrangement of data (data structure) to permit fast searching.
- Index in the back of a text book.
  - Author names --- page numbers.
  - Key words --- page numbers.
- The web is a large "book".

## Inverted Files

**A file is a list of words by position**

- First entry is the word in position 1 (first word)
- Entry 4562 is the word in position 4562 (4562<sup>nd</sup> word)
- Last entry is the last word

**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 Files for Multiple Documents

WORD	NDOCS	PTR	LEXICON			
			DOCID	OCCUR	POS 1	POS 2
jezebel	20	→	34	6	1	118
			44	3	215	2291
			56	4	5	22
jezerit	1	→	3922	3981	5002	
			566	3	203	245
jeziah	1	→	287			
jeziel	1	→				
jeziah	1	→	67	1	132	
jezoar	1	→				
jezrahiliah	1	→				
jezreel	39	→	107	4	322	354
			232	6	15	195
			677	1	481	248
			713	3	42	312
						802

**WORD INDEX**

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

## Ranking (Scoring) Hits

- Hits must be presented in some order
- What order?
  - Relevance, recency, popularity, reliability?
- Some ranking methods
  - Presence of keywords in title of document
  - Closeness of keywords to start of document
  - Frequency of keyword in document
  - Link popularity (how many pages point to this one)
- Can the page owner control?
- Can you find out what order is used?

18-Jan-01

CSE 490i

7

## Link Popularity

- How many pages link to this page?
  - Popular sites have many links pointing at them.
  - The answer is .... Yahoo!
  - (but what was the question again?)
- You can't just rely on number of links!
  - Need to factor in query words.
  - Need to weight links differently.

18-Jan-01

CSE 490i

8

## Authorities & Hubs

- **Authorities:** sites (pages) that are pointed to by many links! [in-degree]
- **Hubs:** sites (pages) that point to many links! [out-degree]
- Assign each page an authority weight & a hub weight. I.e, continuous parameters.
- Bigger weights are "better".

18-Jan-01

CSE 490i

9

## Authorities & Hubs II

- Authority(P)= sum of hub weights of pages that point to P.
- Hub(P)= sum of authority weights of pages that it points to.
- Isn't this definition circular?

18-Jan-01

CSE 490i

10

## Algorithm Outline

- Start w/ seed set of core docs.
- Initialize weights to be equal. Then,
  - Recompute authority weights based on hub weights.
  - Recompute hub weights based on authority weights.
- Repeat until ranking doesn't change much.

18-Jan-01

CSE 490i

11

## Information Retrieval (IR)

- Given a large repository of documents, how do I get at the ones that I want
  - Examples: Lexus/Nexus, Medical reports, AltaVista
- Different from databases
  - Unstructured (or semi-structured) data
  - Information is (typically) text
  - Requests are (typically) word-based

18-Jan-01

CSE 490i

12

## Information Retrieval Task

- Start with a set of documents
- User specifies *information need*
  - Keyword query, boolean expression, high-level description
- System returns a list of documents
  - Ordered according to relevance
- Known as the *ad-hoc retrieval problem*

18-Jan-01

CSE 490i

13

## Basic IR System

- Use word overlap to determine relevance
  - Word overlap alone is inaccurate
- Rank documents by similarity to query
- Computed using *Vector Space Model*

18-Jan-01

CSE 490i

14

## Vector Space Model

- Represent documents as a matrix
  - Words are rows
  - Documents are columns
  - Cell  $i,j$  contains the number of times word  $i$  appears in document  $j$
  - Similarity between two documents is the cosine of the angle between the vectors representing those words

18-Jan-01

CSE 490i

15

## Vector Space Example

- a: System and human system engineering testing of EPS
- b: A survey of user opinion of computer system response time
- c: The EPS user interface management system
- d: Human machine interface for ABC computer applications
- e: Relation of user perceived response time to error measurement
- f: The generation of random, binary, ordered trees
- g: The intersection graph of paths in trees
- h: Graph minors IV: Widths of trees and well-quasi-ordering
- i: Graph minors: A survey

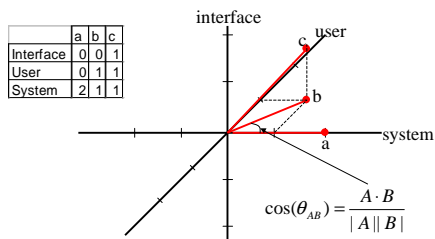
	a	b	c	d	e	f	g	h	i
Interface	0	0	1	0	0	0	0	0	0
User	0	1	1	0	1	0	0	0	0
System	2	1	1	0	0	0	0	0	0
Human	1	0	0	1	0	0	0	0	0
Computer	0	1	0	1	0	0	0	0	0
Response	0	1	0	0	1	0	0	0	0
Time	0	1	0	0	1	0	0	0	0
EPS	1	0	1	0	0	0	0	0	0
Survey	0	1	0	0	0	0	0	0	1
Trees	0	0	0	0	0	1	1	1	0
Graph	0	0	0	0	0	0	1	1	1
Minors	0	0	0	0	0	0	0	1	1

18-Jan-01

CSE 490i

16

## Vector Space Example cont.



18-Jan-01

CSE 490i

17

## Similarity in Vector Space

$$A \cdot B = A_1 B_1 + A_2 B_2 + \dots + A_n B_n$$

Measures word overlap

$$\cos(\theta_{AB}) = \frac{A \cdot B}{|A| |B|}$$

Normalizes for different length vectors

$$|A| = \sqrt{\sum_{i=1}^n A_i^2}$$

Other metrics exist

18-Jan-01

CSE 490i

18

## Answering a Query Using Vector Space

- Represent query as vector
- Compute distances to all documents
- Rank according to distance
- Example
  - “computer system”

Query	a	b	c	d	e	f	g	h	i
Interface	0	0	1	0	0	0	0	0	0
User	0	1	1	0	1	0	0	0	0
System	1	2	1	1	0	0	0	0	0
Human	0	1	0	0	1	0	0	0	0
Computer	1	0	1	0	1	0	0	0	0
Response	0	0	1	0	0	1	0	0	0
Time	0	0	1	0	0	1	0	0	0
EPS	0	1	0	1	0	0	0	0	0
Survey	0	0	1	0	0	0	0	0	1
Trees	0	0	0	0	0	0	1	1	0
Graph	0	0	0	0	0	0	1	1	1
Minors	0	0	0	0	0	0	0	0	1

18-Jan-01

CSE 490i

19

## Common Improvements

- The vector space model
  - Doesn't handle morphology (eat, eats, eating)
  - Favors common terms
- Possible fixes
  - Stemming
    - Convert each word to a common root form
  - Stop lists
  - Term weighting

18-Jan-01

CSE 490i

20

## Handling Common Terms

- Stop list
  - List of words to ignore
    - “a”, “and”, “but”, “to”, etc.
- Term weighting
  - Words which appear everywhere aren't very good discriminators
  - Apply a scaling factor to frequencies

18-Jan-01

CSE 490i

21

## Term Weighting

- tf.idf weighting
  - Term frequency-inverse document frequency
  - $tf$  = term frequency     $df$  = document frequency
- Family of schemes
  - Depends on details of  $tf$  and  $df$  scaling
    - Straight frequencies as above
    - log-based weightings
    - $1 + \log(tf)$  for  $tf > 0$        $\log(\frac{N}{df})$

18-Jan-01

CSE 490i

22

## Extensions

- Meet demands of web-based systems
- Modified ranking functions for the web
- Relevance feedback
- Query expansion
- Document clustering
- Latent Semantic Indexing
- Other IR tasks

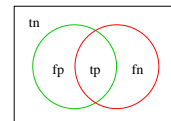
18-Jan-01

CSE 490i

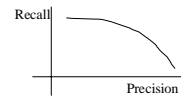
23

## Measuring Performance

- Precision  $\frac{tp}{tp + fp}$ 
  - Proportion of selected items that are correct
- Recall  $\frac{tp}{tp + fn}$ 
  - Proportion of target items that were selected
- Precision-Recall curve
  - Shows tradeoff



System returned these  
Actual relevant docs



18-Jan-01

CSE 490i

24

## IR on the Web

- Query AltaVista with “Java”
  - Almost  $10^7$  pages found
- Avoiding latency
  - User wants (initial) results **fast**
- Solution
  - Rank documents using word-overlap
  - Use special data structure - *inverted index*

18-Jan-01

CSE 490i

25

## Spamdexing

- Attempting to influence retrieval ranking by altering a web page
- Repeating words many times
- Adding words not present on the page
- Adding non-printing words to web pages through HTML <META> tags
- Some engines (Lycos) ignore META tags

18-Jan-01

CSE 490i

26

## Improved Ranking on the Web

- Not just arbitrary documents
- Can use HTML tags and other properties
  - Query term in <TITLE></TITLE>
  - Query term in <IMG>, <HREF>, etc. tag
  - Check date of document (prefer recent docs)
  - Search in a particular IP domain

18-Jan-01

CSE 490i

27

## Relevance Feedback

- System returns initial set of documents
- User identifies relevant documents
- System refines query to get documents more like those identified by user
  - Add words common to relevant docs
  - Reposition query vector closer to relevant docs
- repeat...

18-Jan-01

CSE 490i

28

## Query Expansion

- Given query, add words to improve recall
  - Workaround for synonym problem
- Example
  - boat → boat OR ship
- Can involve user feedback or not
- Can use thesaurus or other online source
  - WordNet

18-Jan-01

CSE 490i

29

## Document Clustering

- Group similar documents
  - Similar means “close in vector space”
- If a document is relevant, return whole cluster
- Can be combined with relevance feedback
- GROUPER
  - <http://www.cs.washington.edu/research/clustering>

18-Jan-01

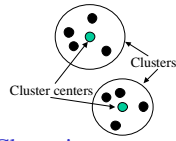
CSE 490i

30

## Clustering Algorithms

- **K-means**

Initialize k cluster centers  
 Loop  
 Assign all document to closest center  
 Move cluster centers to better fit assignment  
 Until little movement



- **Hierarchical Agglomerative Clustering**

Initialize each document to a singleton cluster  
 Loop  
 Merge two closest clusters  
 Until k clusters exist

*Many ways to measure distance between clusters*

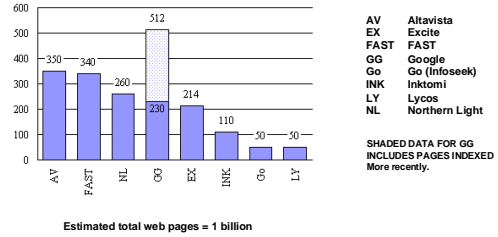
18-Jan-01

CSE 490i

31

## Search Engine Sizes (June 2, 2000)

Millions of Web Pages Indexed



Estimated total web pages = 1 billion

SOURCE: SEARCHENGINEWATCH.COM

18-Jan-01

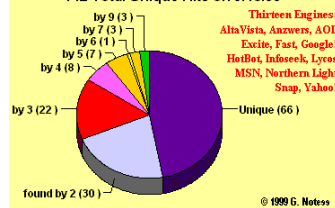
CSE 490i

32

## Search Engines Disjointness

### Overlap of Five Searches

142 Total Unique Hits on 9/10/99



SOURCE: SEARCHENGINESHOWDOWN

18-Jan-01

CSE 490i

33

## Conclusions

- Web search engines use IR + web-specific features such as link popularity.
- Search engines are limited in their coverage (I.e., limited recall). **Meta-search!**
- Search engines often don't get the 'right' results (limited precision). **Clustering!**

18-Jan-01

CSE 490i

34