

COS 435: Information Retrieval, Discovery, & Delivery

Questions about how we **find, organize, evaluate** and **deliver** information

1

Historic Goals

“ to organize the world's information and make it universally accessible and useful”

“ an individual stores all his books, records, and communications, and which is mechanized so that it may be consulted with exceeding speed and flexibility. It is an enlarged intimate supplement to his memory.”

2

Historic Goals

“Google's mission is to organize the world's information and make it universally accessible and useful” [Google's mission statement](#), ~ 1998.

“A memex is a device in which an individual stores all his books, records, and communications, and which is mechanized so that it may be consulted with exceeding speed and flexibility. It is an enlarged intimate supplement to his memory.” [Vannervar Bush, As we may think, Atlantic Monthly, July 1945.](#)

3

Concepts

- Data ?
- Information ?
- Content?
- Knowledge ?

4

One definition

- **Data:** 0's and 1's stored, with or without structure
- **Information:** *Data* with semantic interpretation
- **Content:** all *information* in a document or collection
- **Knowledge:** a functional understanding of *information*

These definitions basically match class discussion; *content* and *knowledge* can be used both narrowly and broadly and we had definitions matching each

5

Data help us?

- Structured data : data base
Tagged, typed
- Semi-structured data: tagged – XML
HTML?
- Unstructured:
 - Text
 - Graphics: 2D, 3D
 - Music
 - Video

6

What do you want?

- Know it there – Data Bases - data retrieval
- Know it when see it – Information Retrieval
- Surprise me – Data Mining (COS 424)

7

Information Retrieval vs Search

discovery of content
+
retrieval of content relevant to query

= search

SEARCH ENGINES

8

Delivery of content

- in *digital libraries*, search tool and content repository over one umbrella organization: e.g. Library of Congress
- on *Web*, actual Web pages not provided by search engines (although can get cached copy sometimes)
 - Where Web pages stored affects delivery

9

What do you want, Part 2

- information need v.s. query form
 - *User* has information need
 - *Retrieval system* has query form
- Does query capture information need?
- **Relevance**
 - *A judgment* by user
 - Compare: *no* sense of relevance in data retrieval

10

How do you do it?

- Model
 - Contents
 - Query
 - Matching of contents to query - results
- Algorithms
 - Effectiveness
 - Efficiency

11

What are performance issues?

- Effectiveness: does search return relevant results ?
- Large amounts data – disks I/O! or not?
- Networking
 - Where is data?
 - Should data be somewhere else?
- Web
 - How find information?
 - How use Web structure?

12

Information Delivery

Broadly construed can mean:

- User Interfaces
- Protocols
- Storage Management
- Bandwidth management

Big question: what is model of interaction?
compare handheld wireless, CS Dept machine

13

Information Delivery cont.

Focus on latter two:

- Storage management
 - Distributed storage
 - Permanence
- Bandwidth management
 - Caching
 - Prefetching
 - Content distribution networks

14

Topics 1

- query models for searching (keyword-based)
- models of documents
- Indexing and inverted files
- Ranking documents
- Using linking structure for Web content analysis
- Semantic and feedback techniques
- User behavior-based relevance criteria; privacy issues
- Manipulating search engine results (SEOs)
- Evaluating retrieval systems

15

Topics 2

- Web crawling
- Document similarity
- Clustering
- Non-text media search: e.g. music, images
- adding structure to information: databases, XML, the semantic Web

16

Topics 3

- system design of search engines: distributed storage and computing
- Information caching
- Content distribution networks
- Reliability and permanence of information

17

Course logistics

- Texts
 - For IR will assign reading from new online text *Introduction to Information Retrieval*
- Test – two, expect not in class.
- Homework, approx. every couple of week
- Presentation – one short
- Project – your choosing with approval

18