

types of information systems

- informal
  - evolve from patterns of human behavior (can be complex)
  - not formalized (i.e., designed)
  - rely on "word of mouth" ("the grapevine")

#### • manual

- $-\ensuremath{\mathsf{formalized}}\xspace$  but not computer based
- historical handling of information in organizations, before computers (i.e., human "clerks" did all the work)
- some organizations still use aspects of manual IS (e.g., because computer systems are expensive or don't exist to relace specialized human skills)
- computer-based
  - automated, technology-based systems
  - typically run by an "IT" (information technology) department within a company or organization (e.g., ITS at BC)

cis20.2-spring2010-sklar-lecl1.1

# computer-based information systems

- data processing systems (e.g., accounting, personnel, production)
- office automation systems (e.g., document preparation and management, database systems, email, scheduling systems, spreadsheets)
- management information systems (MIS) (e.g., produce *information* from *data*, data analysis and reporting)
- decision support systems (DSS) (e.g., extension of MIS, often with some intelligence, allow prediction, posing of "what if" questions)
- executive information systems (e.g., extension of DSS, contain strategic modeling capabilities, data abstraction, support high-level decision making and reporting, often have fancy graphics for executives to use for reporting to non-technical/non-specialized audiences)



cis20.2-spring2010-sklar-lecll.1

social effects of IS

- change management
- broad implementation (not just about software)
- education and training
- skill change
- societal and cultural change

integrative models
• computers in society
• the internet revolution (internet 2, web 2.0)
• "big brother"
• ubiquitous computing

# what is information?

- definition comprises ideas from philosophy, psychology, signal processing, physics...
- OED:
  - information = "informing, telling; thing told, knowledge, items of knowledge, news"
  - knowledge = "knowing familiarity gained by experience; person's range of information; a theoretical or practical understanding of; the sum of what is known"
- other ideas:

cis20.2-spring2010-sklar-lecll.1

- relating data to *context*
- must be recorded
- $\ensuremath{\mathsf{has}}$  potential to become knowledge
- what is the relationship between data and information and knowledge and intelligence???

### types of information

- can be differented by:
  - form
  - $-\operatorname{content}$
  - quality
  - $-\ensuremath{\mathsf{associated}}\xspace$  information
- properties
  - can be communcated electronically (methods: broadcasting, networking)
  - $-\ensuremath{\mathsf{can}}$  be duplicated and shared (issues: ownership, control, maintenance, correction)

#### 9 cis20.2-spring2010-sklar-lecll.1

intuitive notion of information (from Losee, 1997)

- information must be something, although its exact nature is not clear
- information must be "new" (repeating something old isn't considered "information"... or is it?)
- information must be true (i.e., not "mis-information")
- information must be about something
- $\bullet$  note human-centered definition that emphasizes meaning and message

#### human perspective

- cognitive processing
  - perception, observation, attention
  - reasoning, assimilating, interpreting, inferring
  - communicating
- knowledge, belief
- belief = "an idea held on some support; an internally accepted statement, result of inductive processes combining observed facts with a reasoning process"
- does "information" require a human mind?

#### cis20.2-spring2010-sklar-lecl1.1

## meaning versus form

- is the form of information the information itself? or another kind of information?
- is the meaning of a signal or message the signal or message itself?
- representation (from Norman 1993)
  - why do we write things down?
    - \* Socrates thought writing would obliterate serious thought
  - \* sound and gestures fade away
  - artifacts help us reason
  - anything not present in a representation can be ignored (do you agree with that?)
  - things left out of a representation are often those things that are hard to represent, or we don't know how to represent them

# The Library of Babel, by Jorge Luis Borges (1941)

- a story about a universe comprised of an indefinite (possibly infinite) number of hexagonal rooms, each containing walls of bookshelves that contain books which, in turn contain all possible combinations of letters
- is this information? data? knowledge? intelligence?
- how is the internet like (or unlike) the library of babel?

cis20.2-spring2010-sklar-lecl1.1

information theory

• Claude Shannon, 1940's, IBM

cis20.2-spring2010-sklar-lecll.1

- studied communication and ways to measure information
- communication = producing the same message at its destination as at its source
- problem: *noise* can distort the message
- message is *encoded* between source (transmitter) and destination (receiver)

#### communication theory

- many disciplines: mass communication, media, literacy, rhetoric, sociology, psychology, linguistics, law, cognitive science, information science, engineering, medicine...
- human communication theory: do you understand what I mean when I say something?
- what does it mean to say a message is received? is received the same as understood?
- the *conduit metaphor*
- meaning: syntactic versus semantic



- answer list is apparently disorganized

cis20.2-spring2010-sklar-lecll.1

# information seeking behavior "berry-picking" model (from Bates 1989) • two parts of a process: • interesting information is scattered like berries in bushes - search and retrieval • the eye of the searcher is continually moving - analysis and synthesis of search results • new information may trigger new ideas about where to search • search tactics and strategies • searching is generally not satisfied by one answer - tactics $\Rightarrow$ short-term goals, single actions, single operators - strategies $\Rightarrow$ long-term goals, complex actions, combinations of operators (macros) • need to keep search on track by *monitoring* search - check: compare next move with current "state" - weigh: evaluate cost/benefit of next move/direction - pattern: recognize common actions - correct: fix mistakes - record: keep track of where you've been (even wrong directions) search tactics - specify: be as specific as possible in terms you are looking for cis20.2-spring2010-sklar-lecll.1 cis20.2-spring2010-sklar-lecll.1

- exhaust: use all possible elements in a query
- $-\ensuremath{\,\mbox{reduce:}}$  subtract irrelevant elements from a query
- parallel: use synonyms ("term" tactics)
- $\operatorname{pinpoint:} focus query$
- block: reject terms
- relevance how can a retrieved document be considered relevant?
  - it can answer original question exactly and completely
  - it can partially answer the question
  - $-\operatorname{it}$  can suggest another source for more information
  - $-\operatorname{it}$  can provide background information for answering the question
  - it can trigger the user to remember other information that will help answer the question and/or retrieve more information about the question

- parametric search
- most documents have "text" and "meta-data", organized in "fields"
- in parametric search, we can associate search terms with specific fields
- example: search for apartments in a certain **geographic neighborhood** within a certain **price range** of a certain **size**
- the data set can be organized using *indexes* to support parametric search

cis20.2-spring2010-sklar-lecl1.1



• assign a weight to each zone (or field) and evaluate:

 $score = 0.6 * (Brooklyn \in neighborhood) + 0.5 * (3 \in bedrooms) + 0.4 * (1000 = price)$ 

• problem:

it is frequently hard for a user to assign a weighting that adequately or accurately reflects their needs/desires

ncide	ence	matrices		

- recall = document (or a zone or field in the document) is a binary vector  $X \in \{0, 1\}^v$
- query is a vector
- score is overlap measure:  $|X \cap Y|$
- example:

## Julius Caesar The Tempest Hamlet Othello Macbeth

Antony	1	0	0	0	1
Brutus	1	0	1	0	0
Caesar	1	0	1	1	1
Calpurnia	1	0	0	0	0
Cleopatra	0	0	0	0	0

score is sum of entries row (or column, depending on what the query is)

cis20.2-spring2010-sklar-lecll.1



- term frequency (how often does a term occur in a document)
- term scarcity in collection (how infrequently does the term occur in all documents in the colletion)
- length of documents searched
- what about density?
- if a document talks about a term more, then shouldn't it be a better match?
- what if we have more than one term? this leads to *term weighting*

cis20.2-spring2010-sklar-lecll.1

- By this scoring measure, the top-scoring play is likely to be the one with the most "of" s is this what we want?
- NOTE that in the IR literature, "frequency" typically means "count" (not really "frequency" in the engineering sense, which would be count normalized by document length...)
- term frequency (tf)
  - $\mbox{ somehow we want to account for the length of the documents we are comparing }$
- collection frequency (cf)
  - the number of occurrences of a term in a collection (also called *corpus*)
- document frequency (df)
  - the number of documents in a collection (corpus) containing the term
- tf x idf or tf.idf
  - tf = term frequency
  - -idf = inverse document frequency; could be 1/df, but more commonly computed as:

$$df_i = log\left(\frac{n}{df_i}\right)$$

cis20.2-spring2010-sklar-leclI.1



- cis20.2-spring2010-sklar-lecl1.1
  - "weight" of term i occurring in document d  $(w_{i,d})$  is then:  $w_{i,d} = tf_{i,d} \times idf_i$

# $= tf_{i,d} \times \log(n/df_i)$

- where
- $tf_{i,d} = frequency of term \ i \ in \ document \ d$
- $n = {\rm total} \ {\rm number} \ {\rm of} \ {\rm documents} \ {\rm in} \ {\rm collection}$
- $df_i =$  number of documents in collection that contain term i
- $-\ensuremath{\mathsf{weight}}$  increases with the number of occurrences within a document
- weight increases with the rarity of the term across the whole collection
- so now we recompute the matrix using the  $w_{i,d}$  formula for each entry in the matrix, and then we can do our ranking with a query

### references

- Dr Phil Trinder, Heriot-Watt University, UK F29IF1 Database and Information Systems (Fall 2007) http://www.macs.hw.ac.uk/~trinder/DbInfSystems/
- Prof Ray Larson and Prof Marc Davis, UC Berkeley IS 202 Information Organization and Retrieval (Fall 2003) http://www.sims.berkeley.edu/academics/courses/is202/f03/
- Prof Christopher Manning and Prof Prabhakar Raghavan, Stanford University CS276 / LING 286 Information Retrieval and Web Mining (Fall 2006) http://www.stanford.edu/class/cs276/

cis20.2-spring2010-sklar-lecll.1

33