



## Evaluation of browsing behaviour and automated subject classification: examples from KnowLib

Koraljka Golub,  
Knowledge Discovery and Digital Library Research Group  
<http://www.it.lth.se/knowlib/>

SFIS, 21 November 2006



## Outline

→ Subject browsing

❖ Automated subject classification

❖ Crawling

❖ Demonstrators

SFIS, 21 Nov 2006

K. Golub, KnowLib's Research

2 of 36



## Subject browsing

- seeking for information resources by examining a hierarchical tree of broader and narrower subject classes into which the resources have been classified
- browsing services
  - for academic users
    - e.g. Intute (<http://www.intute.ac.uk/>), Renardus (<http://www.renardus.org>),
  - commercial
    - e.g. Yahoo! directory (<http://dir.yahoo.com/>)
    - Google Directory (<http://www.google.com/dirhp>)
  - collaborative
    - DMOZ (<http://dmoz.org/>)
- browsing vs. searching
  - contradictory claims and research results

SFIS, 21 Nov 2006

K. Golub, KnowLib's Research

3 of 36



## Structures for subject browsing

- traditional: classification schemes, thesauri, subject heading systems
- from the WWW: ontologies, search-engine directories
- some better for browsing than others
  - hierarchical structure
  - document collection
  - names of subjects

SFIS, 21 Nov 2006

K. Golub, KnowLib's Research

4 of 36



## Renardus

- <http://www.renardus.org>
- integrated searching and browsing of ca. 80,000 resources from major European subject gateways
  - simple and advanced searching
  - browsing through Dewey Decimal Classification (DDC)
  - browsing support features

SFIS, 21 Nov 2006

K. Golub, KnowLib's Research

5 of 36



## Research issues

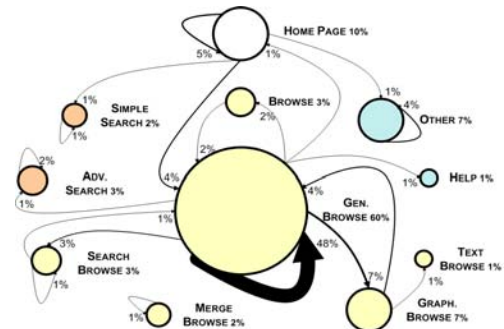
- the balance between browsing, searching and mixed activities
- the degree of usage of the browsing support features
- typical sequences of user activities and transition probabilities in a session, esp. in traversing the hierarchical DDC browsing structure
- typical entry points and referring sites

SFIS, 21 Nov 2006

K. Golub, KnowLib's Research

6 of 36

- log analysis
  - users do not need to be directly involved
  - catches unsupervised behaviour
  - every activity within the system tracked
- cleaned and categorized entries (ca. 460,000) grouped into user sessions (ca. 73,000)
  - all entries from the same address
  - time gap between two entries less than 1 hour
  - one-entry sessions & sessions shorter than 2 seconds removed
- sample
  - 16 months (2002/2003)



- 76% of all activities are browsing
  - majority start using Renardus at a browsing page because directly referred by a search engine
  - layout of Home page "invites" browsing
    - also users starting at Home page predominantly use browsing
- good usage of browsing support features, esp.:
  - graphical overview
  - search entry to browsing pages
- 5% of all activities are searching

- 71% people referred by search engines (mostly Google and Yahoo!)
  - 87% browsing, 2.7% searching
- 22% start at Home page
  - 57% browsing, 12.5% searching
  - more browsing activities per session than the other type
  - use non-browsing activities 3x (Other) and 5x (searching) as often
  - have 2x as many activities per session (ca. 10)
  - they use the service elaborately, in a way system designers intended

- 60% of all activities
- 2/3 are in unbroken browsing sequences
  - up to 86 steps
- keywords
  - good chance of finding browsing pages when using more than one search term

- given proper conditions, browsing is heavily used
  - browsing support features are also heavily used
- it is implied that DDC could serve as a good browsing structure, including terminology

**KnowLib** Outline

---

- ✓ Subject browsing
- Automated subject classification
- ❖ Focused crawling
- ❖ Demonstrators

SFIS, 21 Nov 2006 K. Golub, KnowLib's Research. 13 of 36

**KnowLib** Automated subject classification

---

- subject classification
  - grouping documents that have a property (topic, theme) in common, further sub-grouping of documents based on finer properties
  - establishing relationships between them
- *automated* subject classification
  - machine-based (statistical, NLP techniques)
- application at KnowLib
  - classification of Web pages for browsing
  - classification of Web pages for focused crawling

SFIS, 21 Nov 2006 K. Golub, KnowLib's Research. 14 of 36

**KnowLib** Approaches

---

- text categorization
- document clustering
- string matching

SFIS, 21 Nov 2006 K. Golub, KnowLib's Research. 15 of 36

**KnowLib** Text categorization

---

- machine learning
  - algorithms
- information retrieval
  - vector-space model
  - evaluation measures
- pre-defined browsing structures
  - learning about categories from pre-existing documents in the categories
  - for Web pages, search-engine directories
- e.g. <http://search.thunderstone.com/texis/websearch/>

SFIS, 21 Nov 2006 K. Golub, KnowLib's Research. 16 of 36

**KnowLib** Document clustering

---

- information retrieval
- vector-space model
- browsing structures automatically derived
  - clusters of similar documents and, partially, relationships between them
  - names of the clusters
  - such structures hard to understand
  - rather unstable as well
- e.g. <http://www.kartoo.com/>, <http://www.clusty.com>

\*Automatically-derived structures often result in heterogeneous criteria for category membership and can be difficult to understand\* (Chen, Dumais 2000 <http://research.microsoft.com/~shensh/cv2000.pdf> )

SFIS, 21 Nov 2006 K. Golub, KnowLib's Research. 17 of 36

**KnowLib** String matching

---

- algorithms
  - usually string-to-string matching against a controlled vocabulary
- pre-defined browsing structures
  - controlled vocabularies
    - usu. classification schemes (good for browsing)
- e.g. <http://engine-e.lub.lu.se/>

SFIS, 21 Nov 2006 K. Golub, KnowLib's Research. 18 of 36

## KnowLib Automated classification issues

- automating subject determination
  - logical positivism
    - subject is a string occurring a certain number of times, in a certain location etc.
    - if document 1 is about subject A, and if document 2 is similar to document 1, then document 2 is also about subject A
- evaluation
  - issue of deriving the correct interpretation of a document's subject matter
  - few end-user evaluations

SFIS, 21 Nov 2006

K. Golub, KnowLib's Research

19 of 36

## KnowLib Similarities between approaches

- document pre-processing and indexing
  - removing stop-words
  - extracting relevant words
- utilization of Web-page characteristics
  - structural elements
  - metadata
  - text neighbouring headings and anchor text
  - text from linked pages
- assumption: idea exchange beneficial

SFIS, 21 Nov 2006

K. Golub, KnowLib's Research

20 of 36

## KnowLib Is there an exchange of ideas?

- main research question
  - to what degree the three communities utilize others' ideas, methods, and findings
- direct links
  - do authors from one community cite authors from another
- indirect links
  - bibliographic coupling of papers
- sample
  - 148 papers: 52 ML, 63 IR, 33 LS

SFIS, 21 Nov 2006

K. Golub, KnowLib's Research

21 of 36

## KnowLib Direct links

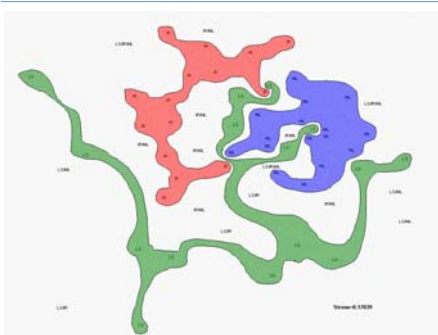
- the ML community uses IR methods and both tended to cite each other to a certain extent
- few cases where LS authors were cited by either of the two other communities and the other way around

SFIS, 21 Nov 2006

K. Golub, KnowLib's Research

22 of 36

## KnowLib Indirect links



SFIS, 21 Nov 2006

K. Golub, KnowLib's Research

23 of 36

## KnowLib Using Web-page elements

- what is the importance of distinguishing between different parts of a Web page?
  - title, headings, main text, metadata
  - what are the appropriate significance indicators?

e.g. <http://froggy.lbl.gov/virtual/>

```
<title>Virtual Frog Dissection Kit Version 2.2</title>
<meta name="description" content="Virtual Frog Dissection Kit">
<meta name="keywords" content="frog dissection K-12 education">
<h2 align="center">Virtual Frog Dissection Kit</h2>
<h2>Frog watch</h2>
main text:
```

"This award-winning interactive program is part of the "Whole Frog" project. You can interactively dissect a (digitized) frog named Fluffy, and play the Virtual Frog Builder Game. The interactive Web pages are available in a number of languages...."

SFIS, 21 Nov 2006

K. Golub, KnowLib's Research

24 of 36

## KnowLib Structural elements and metadata

- collection
  - 1003 Web pages in engineering
- Ei classification scheme
  - 6 main classes
  - decimally subdivided
  - up to 5 hierarchical levels

4 Civil Engineering  
 ...  
 44 Water and Waterworks Engineering  
 441 Dams and Reservoirs  
 ...  
 445 Water Treatment  
 445.1 Water Treatment Techniques  
 445.1.1 Potable Water Treatment Techniques  
 ...

## KnowLib Approach

- algorithm: string matching
  - when a match is found, the corresponding class is assigned, with a relevance score, based on:
    - which term is matched (single word, phrase, Boolean)
    - type of class matched (main or optional)
    - the part of the Web page in which the match is found
- significance indicators
  - derived using various measures of correctness
    - precision and recall
    - semantic distance
    - multiple regression

## KnowLib Major results

- title performs best, followed by headings, metadata, and text
- necessary to use all structural elements and metadata (not all of them occur on every Web page)
- how to combine them not important
  - the best combination was only 3% better than the worst one

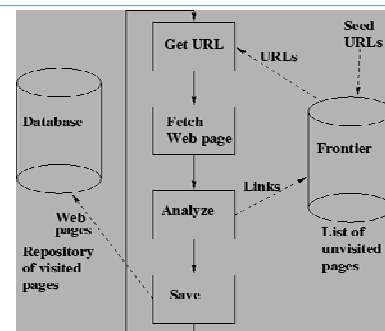
## KnowLib Near-future research

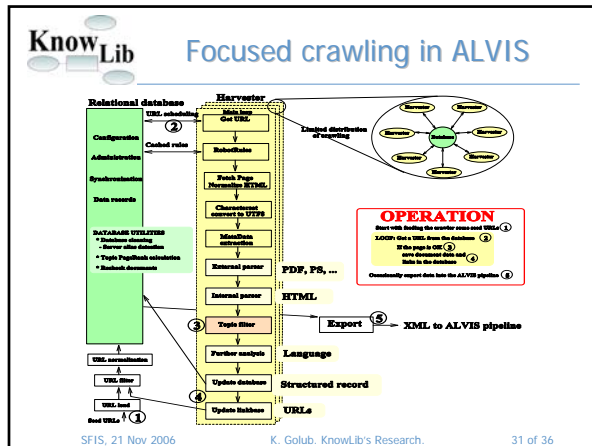
- string-matching
  - termlist expansion (using NLP)
  - adjusting term weighting
  - adjusting cut-offs
- comparison between string-matching and SVM (text categorization)
  - 1) on a test collection, using standard precision and recall
  - 2) with users

## KnowLib Outline

- ✓ Subject browsing
- ✓ Automated subject classification
- ➔ Focused crawling
- ❖ Demonstrators

## KnowLib Simple crawling





- ## KnowLib Focused crawling in ALVIS
- Combine focused crawler
    - availability: <http://combine.it.lth.se/>
    - download, documentation, publications
  - focused crawling in ALVIS:
    - <http://www.it.lth.se/knowlib/publ/ESWC.xfiq.v4.pdf>
  - ALVIS <http://www.alvis.info>
- SFIS, 21 Nov 2006 K. Golub, KnowLib's Research. 32 of 36

- ## KnowLib Outline
- ✓ Subject browsing
  - ✓ Automated subject classification
  - ✓ Focused crawling
- ➔ Demonstrators
- SFIS, 21 Nov 2006 K. Golub, KnowLib's Research. 33 of 36

- ## KnowLib Demonstrators
- <http://www.it.lth.se/knowlib/demos.htm>
  - also, automatic vocabulary mapping
    - <http://dbkit02.it.lth.se/exp/map/>
- SFIS, 21 Nov 2006 K. Golub, KnowLib's Research. 34 of 36

