AHIS Dataset Search Engine: An intelligent approach to EML Data Management Hung V. Nguyen, Corinna Gries, and Hasan Davulcu

Overview

et ide"T' ey

n Number of Sites using a Keywo in any specific source (EML, DTOC etc.)

0% of all keywords are us

Mean=219 9tt Dex =3.19

at 4 or fewer LTER site

Number of LTER Sites

72.2% of all keywords are used at only a single LTER site







School of Computing and Informatics

ARIZONA STATE UNIVERSITY

Screenshot Example flexibility for users to either "water"

narrow down or expand the search for relevant datasets by building boolean search based on "contextual-aware" related terms of current

query term(s).

Our system offers a

"water AND maricopa" The human evaluation shows that our approach yields high accuracy of relevance in terms of both

related keywords and datasets.

"(water AND maricopa) OR

ecosystem"

Our system also offers a set of Web-Based applications that let the users update new data from EML Exist database, prune irrelevant keywords and re-index filtered keyword set for search engine

Contact:

Hung V. Nguyen is a Ph.D. candidate at CSE Department, ASU. His research interests are Web/Text Mining, Web Advertising, Machine Learning and Applications in Data Mining. Email: hung@asu.edu

Web: http://www.public.asu.edu/~vnguve1

Gries received her Ph.D in 1988. Currently, she is information manager at the Central Arizona - Phoenix LTER site and co-chair of the LTER information managers committee. Her research interests are the implementation of community standards for discovery, access, visualization and use of ecological data, development of large natural history collections databases and online collection management tools. She is leading development and implementation of Arizona Hydrologic Information System as a pilot project on the State level and a node to CUAHSI and the USGS Geologic Information System. Email: corinna@asu.edu

Hasan Davulcu is an Asst. Professor at CSE Department, ASU.

Dr. Davulcu directs CIPS Lab. Lab's research focuses on developing novel data mining techniques and tools for structuring and organizing unstructured sources such as text, Web and biological data into semantic machine processable information. His interests also include Workflows, Web Services and Databases Systems.









Overview Architecture and Work Flow of the System

• Text Mining: Study the statistics of terms correlation will help to achieve: Higher accuracy in suggesting related terms for users

· Ecological data is widely organized and documented

in Ecological Markup Language - EML format.

Term Ecological Research (LTER) is collected

Some times a user's search query may be an

when the information need is well described, a search engine or information retrieval system may

not be able to retrieve documents matching the

understand the important keywords and relations

among them. A large collection of terms from Long

imperfect description of their information need. Even

query as stated. In this project, we develop a search

Data extraction: Extract informative keywords from

important parts of a dataset, namely, abstract

engine empowered with text mining techniques for ecological datasets to bridge this gap

Efficient dataset retrieval system needs to

· Higher accuracy in relevance of retrieved datasets

Story:

Technique:

are pruned

keyword list, title. Stop words

(e.g. and, or, the, an, etc)

Search Engine **Data Loading Engine** Process the user's queries Setup connection to Exist Database Return the relevant datasets Update the searching repository with new datasets Return the related keyword list to each query term Extract important keywords from new datasets (using indexing engine) Filtering Tool **Indexing Engine** Construct Term-Document (term-doc) matrix M: $M_{ij} =$ frequency of term *i* in document *j* Allow users to prune irrelevant or meaningless keywords Compute the term correlation matrix S $S = M * M^T$ Update the searching repository with updated keyword lists Normalize S to obtain C $\overline{S_{u,u} + S_{v,v} - S_{u,v}}$