

A KNOWLEDGE MANAGEMENT PROTOTYPE

Mary S. Neff and James W. Cooper
IBM Thomas J. Watson Research Center

1. Background

As organizations become increasingly competitive, they recognize the capital inherent in the accumulated knowledge of their members. In a rapidly growing organization, it becomes ever more important to encapsulate and store this knowledge and invent ways to make it available to newer members of the team. We describe a Knowledge Management prototype, dubbed *Avocado*, a project that grew out of the need to provide sophisticated document analysis as an aid to retrieving more useful data from a technical document collection. It has evolved from earlier incarnations with different names [3, 4].

Avocado combines a number of components, several of which are already available in IBM products, e.g., Intelligent Miner for Text. Of interest here is its use of natural language processing technology and corpus-based NLP techniques in the foreground, and databases constructed using NLP technology in the background.

We applied previous versions of the prototype to several collections of information technology documents and developed a user interface for representing the results. Avocado represents a more ambitious effort in that it incorporates some new components and is applied to a collection of “real data” that is more heterogeneous in form and content than any of our previous collections.

2. Avocado

The problem of finding important and relevant documents in an online collection is increasingly difficult as documents proliferate. Typically, searching for information in a document collection involves refining a query, scanning a large list of documents returned by the search engine to determine their relevance, and then searching the documents for the desired information. Avocado’s NLP and corpus-based NLP techniques assist the user both on the “query end” of the process, with Prompted Query Refinement ([2, 3]) and Lexical Navigation and on the “document viewing” end with Automatic Summarization, Keyword Extraction, and Active Markup ([4]).

These several components rely on some more basic technologies, which run in the background at the same time that the documents are indexed for the search engine. They share the same data structures and work together to identify and index names, multiword terms, abbreviations, and relations (named and unnamed), and then count their frequency in the collection.

3. Prompted Query Refinement and Lexical Navigation

Once the user enters a query, Prompted Query Refinement uses the index of names, terms, and relations (the collection vocabulary and the Context Thesaurus) to display to the user other related or possibly related terms that co-occur in the collection with the terms in the query. By selecting some of these related terms to be added to the query, the user can refine the query directly, without having to think up or type in any more items.

In addition to the terms proposed by the Context Thesaurus, the system also retrieves terms from the collection database that are related to terms in the query and can display the nature of the relations. The relationships may be named relations (“CEO of,” “makes,” “is located in”) [1] or unnamed relations, where terms have a strong bi-directional relationship. These relationships can be viewed in lists or plotted graphically; all displayed terms are available for adding to a query.

4. Automatic Summarization and Keyword Extraction

For document summarization, we use a shallow summarization by sentence extraction method. Relying on the statistics of items in the document vocabulary and comparing the relative frequency of items in the document with that of items in the collection, we arrive at a notion of *salience* of items in the document. Terms in title and headings are also considered salient. The most salient items become the keywords. Sentences are scored according to the salience of the terms in them and their position in the document or discourse structure, and the most salient sentences are extracted for a summary. Such a summary is not necessarily coherent, but we try to minimize this problem in the way that the summary is displayed and used. In Avocado, we use a short, 4-sentence extract in the document hit list, and a longer one of user-specified length in a frame displayed above the document.

5. Active Markup

Active Markup is a method of navigation through a group of documents. A document is displayed together with an upper frame that contains a list of the most salient terms (highlighted in both summary and document by category or salience) and a summary of the desired length. The keywords are active components that can cause the server to return related information. In this

implementation, the active components are used to query the server for a list of related terms. The related terms can then be used to compose another query for another list of documents. The summary sentences in the upper frame are hot-linked to their sources in the document, enabling the user to click to skip down to important information. This active markup approach coupled with the computer-generated summaries provides a form of “query-free” searching.

6. Implementation

The system is a Java client running in a frame of a web browser, which connects to a Java-based server running on Windows NT using Java RMI. This server in turn connects to the database using JDBC and launches programs for carrying out the initial search and for producing the final summary as a pair of linked HTML documents which are displayed with the most salient keywords and summary in an upper frame and the complete document in the lower frame. The keywords appear with a set of JavaScript form. Clicking on these form buttons launches a Java applet is used to display related keywords and the documents containing them. This constitutes “Active Markup” of the document and provides an approach for query-free searching of the lexical neighborhood of the document.

7. Status

We indexed the IBM Global Services consultants’ reports on customer engagements for the Avocado prototype. The data are much more problematic than what is found in well-edited news story or article formats. There are 50 large Lotus Notes databases, each for a different industry, with different editing, keyword, and submission criteria. Most of the documents have attachments in Word, WordPro, AmiPro, Freelance, PowerPoint, PS and PDF. Further, the interesting parts of the documents are the attachments. Not all the documents are in English, and some have no significant content (outlines, templates, management-speak). We will report on our experiences with this “real world” collection. A demonstration or ScreenCam movie will also be available.

References

[1] BYRD, R. J. and Y. RAVIN, Identifying and Extracting Relations in Text, NLDB99, 1999.

[2] COOPER, J. W. and R. J. BYRD, Lexical Navigation: Visually Prompted Query Expansion and Refinement, in: Proceedings of DIGLIB97, Philadelphia, PA, July, 1997.

[3] COOPER, J. W. and R. J. BYRD, OBIWAN - A Visual Interface for Prompted Query Refinement, in: Proceedings of HICSS-31, Kona, Hawaii, 1998.

[4] NEFF, M. S. and J. W. COOPER., ASHRAM: Active Summarization and Markup, in: Proceedings of HICSS-32, Maui, Hawaii, 1999.