

CLUSTERING OF WEB TEXT

C.Guru Sunanda^{1*}, K. Ishthaq Ahamed^{2*}, Y.Rama Mohan^{3*}

1. M.Tech (CSE) Student, Dept of CSE, G.Pullareddy Engg College, Dist: Kurnool, AP, India
2. Assoc.Prof, Dept of CSE, G.Pullareddy Engg College, Dist: Kurnool, AP, India
3. Asst.Prof, Dept of CSE, G.Pullareddy Engg College, Dist: Kurnool, AP, India

Keywords:

Clustered Data,
Vector Space
model, Feature
Combination

Abstract: Over the last decade, organizations have invested heavily in systems that enable rapid access to structured data. However, this data represents a fraction of all corporate information. A far larger volume exists as text - in documents, web pages, manuals, reports, email, faxes, and presentations etc. In this paper, we propose a new vision of innovation stimulation, targeting both locally based clusters and new forms of innovation hubs in html text. Considering the tuning guidelines in web text document as a starting point the environment will be different from the lab and results may vary. If we run a different mix of applications or if users work with large reports or data extractions, additional tuning may be required to optimize the performance of your implementation. Hence of in this paper, we try to implement the hyperlink based documents which are clustered and give the citation to the web text.

1. INTRODUCTION

In the age of modern technology of changing technology, A valuable Web site provides information not just data. There is a continued explosion in both content volume and content types (documents, images, streaming media, instant messages, e-mail and so on). The vector space model does not regard word order. We have tried to extend it with nominal phrases in different ways. We have also tried to differentiate between homographs, words that look alike but mean different things, by augmenting all words with a tag indicating their part of speech. None of our experiments using phrases or part of speech information have shown any improvement over using the ordinary model.

Many text clustering methods use the same theoretical foundation as search engines, the vector space model. It is a model for representing (the content of) texts. The following sections give a brief introduction to it. In the vector space model each text in a set of texts is represented by a vector in a high-dimensional space, with as many dimensions as the number of different words in

* C.Guru Sunanda

M.Tech (CSE) Student, Dept of CSE, G.Pullareddy Engg College, Dist: Kurnool, AP, India

the set. Each text is assigned weights (values) in the indices (dimensions) based on what words appear in them. These weights can be thought of as modeling how important the corresponding word is deemed to be to explain the content of the text. Each weight is dependent on whether (and how often) the word appears in the text and in the entire set. Texts whose vectors are close to each other in this space are considered to be similar in content. If two documents share one or more words, then we consider them to be semantically similar. Extending this notion to links, if two documents share one or more links or in-links, then we consider them to be similar as well. This simple observation is the key to the present paper. We propose a precise notion to capture the similarity between two hypertext documents along all the three features in an unified fashion.

2. RELATED WORK

Internet searches have shown that significant information uplift can accrue from search technology. Without search engines, the Internet would still have billions of web pages, but surfers would have to know URLs a priori, or navigate through directories, to locate pages of interest. Clearly it is Search that makes Google popular,

and the Internet more useful even as the amount of information on it grows at a rapid pace. Text clustering can be used to explore the contents of a text set. We have developed a visualization method that aids such exploration, and implemented it in a tool, called in-format. It presents the representation matrix directly in two dimensions. When the order of texts and words are changed, by for instance clustering, distributional patterns that indicate similarities between texts and words appear.

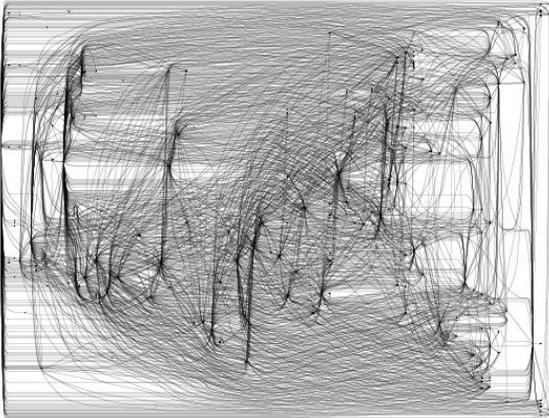


Fig.2.1 showing the clustered structure of data in graphical representation of net web model

The motive for creating a partition, be it a categorization or a clustering, can be of several different kinds. It could, for instance, be beneficial to have a set of texts divided into groups of different levels of readability, so that you can choose which text to read depending on how good a reader you are. Texts can also be categorized into genres in several ways. There are many different clustering algorithms. Most of them need to know the similarity between the objects (texts). Some of them need a representation for each object, and a definition of similarity, so they can calculate it when necessary. How the objects are represented and the definition of similarity differs between applications. It is usually convenient to build a representation and define similarity in terms of it. There are many ways to achieve this. If, for instance, the objects can be represented as points in an n-dimensional vector space, dissimilarity could be defined as the distance between them. When given a set of objects and the similarity between them, a clustering algorithm outputs a partition that tries to satisfy some criteria. It could be that the objects in each cluster should be as similar as possible. However, in order for the result to be useful at all, we must have reasons to believe that the similarity definition reflects the similarity between the actual objects.

3. METHODS

Text clustering divides a set of texts into clusters. It may be used to uncover the structure and content of unknown text sets as well as to give new perspectives on familiar ones. The main contributions of this thesis are an investigation of text representation for Swedish and some extensions of the work on how to use text clustering as an exploration tool. We have also done some work on synonyms and evaluation of clustering results.

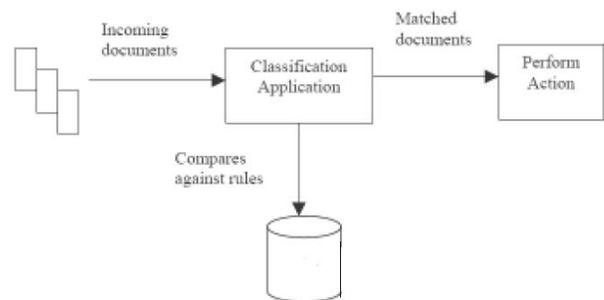


Fig.3.1 Showing the classification of Data (Web contents)

In the fig. 3.1; the data flow in one direction get classified on the middle base then the filtration or mining of clustered show a in a graphical based data. Information on the Internet consists overwhelmingly of web pages. In the Intranet, information – data and content – is spread across web pages, databases, mail servers or other collaboration software, document repositories, file servers, and desktops. An Intranet search engine must be able to search an organization's web content, its applications, databases, and mail through the same interface.

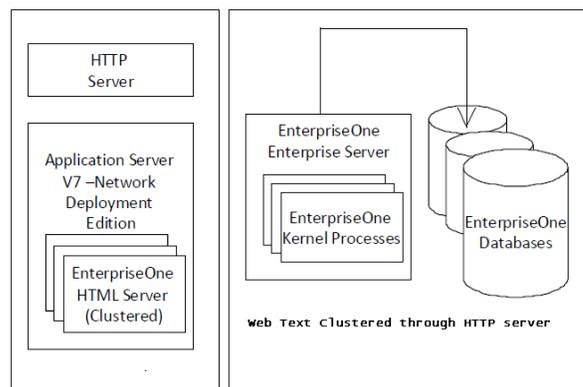


Fig.3.2 Showing the HTTP server based clustering Architecture

In the fig. 3.2; the http server which in one layer contact another layer of MVC architecture based cluster based approach where the data or information itself get clustered. It depends upon how much a system would improve on several factors, among others the length of the queries and the documents. This has probably a strong connection to the collocation effect; a query with many words, at least partly defines the meaning of a homograph in it, since the documents that are retrieved contain most of the words and these tend to come from the same domain. Most queries put to search engines are short. Many relevant documents not containing the few query words are not retrieved. This can, at least theoretically, be remedied by query expansion, in which the query is expanded with words that are related to those already in it. This may be accomplished in many ways. In relevance feedback the user marks some of the retrieved documents as relevant. The system then uses these to reformulate the query, by for instance expanding it with frequent words in the relevant documents. Most methods for query expansion not involving a user utilize some sort of thesaurus, which may be either manually constructed or built using statistics of word co-occurrences. Manually constructed thesauri with word relations, such as WordNet10, are often elaborate and provide many kinds of relations. They are normally constructed for general purposes and the many relations may be hard to adapt to a specific task. Lastly the method may be beneficial in the case of high or huge data in the perspective of cost and economical value to the client.

4. CONCLUSION

As of technology and the computer science is considered where web plays an important role. In this conclusion phase of our paper we stress on the concept of vector space clustering. Internet search engines like Google use the links that URLs provide between web pages to deduce the importance or relevance of a document in a given search. Unfortunately, Intranet resources do not invariably vote for each other by URL links: a document authored in PDF may not link to the database record of a customer that it describes. Consequently, different techniques are needed for high relevance when it comes to Intranet search. In this paper we try to concentrate on clustered data rather than the dimension of data. The vector space model gives the algorithmic approach may of k the element. Hence in the future work we likely to develop the base tree structured algorithmic approach to find the highest optimized clustered data in the web based template extraction.

5. REFERENCE

- [1] BOTAFAGO, R. A. Cluster analysis for hypertext systems. In ACM SIGIR (1993).
- [2] BRADLEY, P., AND FAYYAD, U. Refining initial points for k-means clustering. In ICML (1998), pp. 91–99.
- [3] CHAKRABARTI, S., DOM, B. E., AND INDYK, P. Enhanced hypertext categorization using hyperlinks. In ACM SIGMOD (1998).
- [4] CHAKRABARTI, S., DOM, B. E., KUMAR, S. R., RAGHAVAN, P., RAJAGOPALAN, S., TOMKINS, A., KLEINBERG, J. M., AND GIBSON, D. Hypersearching the web. *Scientific American* (June 1999).
- [5] CHAKRABARTI, S., DOM, B. E., RAGHAVAN, P., RAJAGOPALAN, S., GIBSON, D., AND KLEINBERG, J. Automatic resource compilation by analyzing hyperlink structure and associated text. In WWW7 (1998).
- [6] CHEN, C. Structuring and visualizing the www by generalized similarity analysis. In ACM Hypertext (1997).
- [7] CHEN, C., AND CZERWINSKI, M. From latent semantics to spatial hypertext—An integrated approach. In ACM Hypertext (1998).
- [8] CROFT, W. B., AND TURTLE, H. R. A retrieval model for incorporating hypertext links. In ACM Hypertext (1989).
- [9] DHILLON, I. S., AND MODHA, D. S. Concept decompositions for large sparse text data using clustering. Tech. Rep. RJ 10147 (95022), IBM Almaden Research Center, 1999.
- [10] FRAKES, W. B., AND BAEZA-YATES, R. *Information Retrieval: Data Structures and Algorithms*. Prentice Hall, Englewood Cliffs, New Jersey, 1992.
- [11] FREI, H. P., AND STEIGER, D. Making use of hypertext links when retrieving information. In ACM European Conference on Hypertext (1992).
- [12] C.-N. Hsu and M.-T. Dung, “Generating Finite-State Transducers for Semi-Structured Data Extraction from the Web,” *Information Systems*, vol. 23, no. 8, pp. 521-538, 1998.
- [13] <http://daisen.cc.kyushu-u.ac.jp/TBDW/>, 2009.
- [14] <http://www.w3.org/html/wg/html5/>, 2009.

[15] N. Kushmerick, "Wrapper Induction: Efficiency and Expressiveness," *Artificial Intelligence*, vol. 118, nos. 1/2, pp. 15-68, 2000.

[16] A. Laender, B. Ribeiro-Neto, A. da Silva, and J. Teixeira, "A Brief Survey of Web Data Extraction Tools," *SIGMOD Record*, vol. 31, no. 2, pp. 84-93, 2002.

[17] B. Liu, R.L. Grossman, and Y. Zhai, "Mining Data Records in Web Pages," *Proc. Int'l Conf. Knowledge Discovery and Data Mining (KDD)*, pp. 601-606, 2003.

[18] W. Liu, X. Meng, and W. Meng, "Vision-Based Web Data Records Extraction," *Proc. Int'l Workshop Web and Databases (WebDB '06)*, pp. 20-25, June 2006.

[19] L. Liu, C. Pu, and W. Han, "XWRAP: An XML-Enabled Wrapper Construction System for Web Information Sources," *Proc. Int'l Conf. Data Eng. (ICDE)*, pp. 611-621, 2000.

[20] Y. Lu, H. He, H. Zhao, W. Meng, and C.T. Yu, "Annotating Structured Data of the Deep Web," *Proc. Int'l Conf. Data Eng. (ICDE)*, pp. 376-385, 2007.

[21] J. Madhavan, S.R. Jeffery, S. Cohen, X.L. Dong, D. Ko, C. Yu, and A. Halevy, "Web-Scale Data Integration: You Can Only Afford to Pay As You Go," *Proc. Conf. Innovative Data Systems Research (CIDR)*, pp. 342-350, 2007.