

Improving Effectiveness by Clustering of BioMedical Documents Using Semi supervised Clustering Method

M.Mounika

M.Tech Student,
Department of CSE,
School of Technology,
GITAM University, Hyderabad.

Dr. Sudarson Jena

Associate Professor,
Department of IT,
School of Technology,
GITAM University, Hyderabad.

Abstract:

The Global content and Mesh Semantic information are considered for clustering the biomedical Documents from whole MEDLER collection and Mesh Semantic information, topic extraction and fast information retrieval or filtering. Document clustering involves the use of descriptors and descriptor extraction. Descriptors are sets of words that describe the contents within the cluster. Document clustering is generally considered to be a centralized process. Clustering high-dimensional data is the cluster analysis of data with anywhere from a few dozen to many thousands of dimensions.

Such high-dimensional data spaces are often encountered in areas such as medicine, where DNA microarray technology can produce a large number of measurements at once, and the clustering of text documents, where, if a word-frequency vector is used, the number of dimensions equals the size of the vocabulary. In this paper, We are considering: local-content (LC) information, global-content (GC) information from PubMed and MESH (medical subject heading-MS) for the clustering of bio-medical documents. The performances of MEDLINE document clustering are enhanced from previous methods by combining both the LC and GC. We propose a semi-supervised spectral clustering method to overcome the limitations of representation space of earlier methods.

Keywords:

Data Mining, Biomedical Documents, Text Mining, Document Clustering, Semi Supervised Clustering, Constraints.

Introduction:

For scientific researchers, the most important is Literature reading to vestige scientific progress and hypothesis. This MEDLINE database [1] contains over 12 million references to scientific literature with about $\frac{3}{4}$ of recent articles including an abstract of the publication. We tested the capabilities of our system to retrieve MEDLINE references which are relevant to the subject of stem cells. The Document clustering [2] is a fundamental operation used in unsupervised document organization, automatic topic extraction and the information retrieval. Clustering involves two techniques Agglomerative hierarchical clustering and k means, these are commonly used for document clustering. The database contains more than 21.6 million records from 5,639 selected publications covering biomedicine and health from 1950 to the present. Originally the database covered articles starting from 1965, but this has been enhanced, and records as far back as 1950/51 are now available within the main index. The database is freely accessible on the Internet via the PubMed interface and new citations are added Tuesday through Saturday. For citations added during 1995-2003: about 48% are for cited articles published in the U.S., about 88% are published in English, and about 76% have English abstracts written by authors of the articles.

MEDLINE uses Medical Subject Headings (MeSH) for information retrieval. Engines designed to search MEDLINE (such as Entrez and PubMed) generally use a Boolean expression combining MeSH terms, words in abstract and title of the article, author names, date of publication, etc. Entrez and PubMed can also find articles similar to a given one based on a mathematical scoring system that takes into account the similarity of word content of the abstracts and titles of two articles. A service such as MEDLINE strives to balance usability with power and comprehensiveness. Keeping with the

fact that MEDLINE's primary user community is professionals (medical scientists, health care providers), searching MEDLINE effectively is a learned skill; untrained users are sometimes frustrated with the large numbers of articles returned by simple searches. Counterintuitively, a search that returns thousands of articles is not guaranteed to be comprehensive. Unlike using a typical Internet search engine, PubMed searching of MEDLINE requires a little investment of time. Using the MeSH database to define the subject of interest is one of the most useful ways to improve the quality of a search. Using MeSH terms in conjunction with limits (such as publication date or publication type), qualifiers (such as adverse effects or prevention and control), and text-word searching is another. Finding one article on the subject and clicking on the "Related Articles" link to get a collection of similarly classified articles can expand a search that otherwise yields few results.

Document clustering involves the use of descriptors and descriptor extraction. Descriptors are sets of words that describe the contents within the cluster. Document clustering is generally considered to be a centralized process. Examples of document clustering include web document clustering for search users. The application of document clustering can be categorized to two types, online and offline. Online applications are usually constrained by efficiency problems when compared offline applications. In general, there are two common algorithms. The first one is the hierarchical based algorithm, which includes single link, complete linkage, group average and Ward's method. By aggregating or dividing, documents can be clustered into hierarchical structure, which is suitable for browsing. However, such an algorithm usually suffers from efficiency problems. The other algorithm is developed using the K-means algorithm and its variants. These algorithms can further be classified as hard or soft clustering algorithms.

Hard clustering computes a hard assignment – each document is a member of exactly one cluster. The assignment of soft clustering algorithms is soft – a document's assignment is a distribution over all clusters. In a soft assignment, a document has fractional membership in several clusters. Dimensionality reduction methods can be considered a subtype of soft clustering; for documents, these include latent semantic indexing (truncated singular value decomposition on term histograms) and topic models.

Four problems need to be overcome for clustering in high-dimensional data:

- Multiple dimensions are hard to think in, impossible to visualize, and, due to the exponential growth of the number of possible values with each dimension, complete enumeration of all subspaces becomes intractable with increasing dimensionality. This problem is known as the curse of dimensionality.
- The concept of distance becomes less precise as the number of dimensions grows, since the distance between any two points in a given dataset converges. The discrimination of the nearest and farthest point in particular becomes meaningless:
- A cluster is intended to group objects that are related, based on observations of their attribute's values. However, given a large number of attributes some of the attributes will usually not be meaningful for a given cluster. For example, in newborn screening a cluster of samples might identify newborns that share similar blood values, which might lead to insights about the relevance of certain blood values for a disease. But for different diseases, different blood values might form a cluster, and other values might be uncorrelated. This is known as the local feature relevance problem: different clusters might be found in different subspaces, so a global filtering of attributes is not sufficient.

- Given a large number of attributes, it is likely that some attributes are correlated. Hence, clusters might exist in arbitrarily oriented affine subspaces. Recent research indicates that the discrimination problems only occur when there is a high number of irrelevant dimensions, and that shared-nearest-neighbor approaches can improve results.

EXISTING SYSTEM:

There are two categories named constraint-based and distance based in the existing method. Constraint-based methods have user-provided labels or constraints to guide the algorithm towards a more appropriate data partitioning. By modifying the objective function for evaluating clustering's, it is done. Thus it includes satisfying constraints, enforcing constraints during the clustering process, or initializing and constraining the clustering based on labeled examples.

Existing System Technique:

K-mean's clustering

1. Choose the number of different clusters, k.
2. Generate k clusters randomly and determine where the cluster centers.
3. Assign each point to the nearest cluster center, where we can define "nearest" wrt one of the distance measures discussed.
4. Recompute the new cluster centers.
5. Repeat the previous steps until some convergence criterion is met.

Existing System Drawbacks:

1. True similarity would not be a simple linear relationship between different similarities.
2. The quality of similarity in a data set may not be same for all document pairs. Some pairs may be more reliable and need more attention.
3. Existing system couldn't manage with a suitable weighting configuration to balance three or more different types of similarities in integrating them.

PROPOSED SYSTEM:

To improve the clustering performance, Semi supervised spectral clustering algorithms are used. The prior knowledge to improve clustering is usually provided by labeled instances or, more typically, by two types of constraints, i.e., must-link (ML) and cannot-link (CL), where ML means that the two corresponding examples should be in the same cluster and CL means that the two corresponding examples which we are considering should not be in the same cluster.

We know that the Spectral clustering is a well accepted method for clustering nodes over a graph or an adjacency matrix, where clustering is a graph cut problem that can be solved by matrix trace optimization.

Scope of the Project:

By improving the performance, we have gone for alternative methods where user can search Biomedical text in our project. Usually, when user will search any text, it has to follow online databases. For searching about biomedical text, user can search documents from PubMed, Medline, PMC, Mesh, etc. These databases contain bulk amount of data. The retrieving of documents from these databases makes the performance slow. For this, we can provide option where to get documents, either from online databases or from our local database. We will make clustering of all our local database documents and can get documents from different clusters with the rank.

Proposed System Technique:

Semi-supervised spectral clustering We usually use Medline, PubMed or some other databases for searching biomedical related documents. In all these databases huge number of documents are available. While retrieving those documents, performance will get slow. Hence we can retrieve some selected documents in our local database. Thus the performance could be increased. And if we go for second time search, No need to go for online Database. Get it from our local data-base only. In our proposed algorithm, set of documents $V (= \{v_1, v_2, \dots, v_N\})$ has to be clustered. Let $\text{Sim}(\bullet, \bullet)$ be the function showing similarity between two inputs, and for example, $\text{Sim}(M, M_)$ outputs similarity between two MeSH main headings M and $M_$. We denote the LC similarity matrix by W_l with the (i, j) -element W_{lij} , the GC similarity matrix by W_g with the (i, j) -element W_{gij} , and the semantic similarity matrix by W_s with the (i, j) -element W_{sij} .

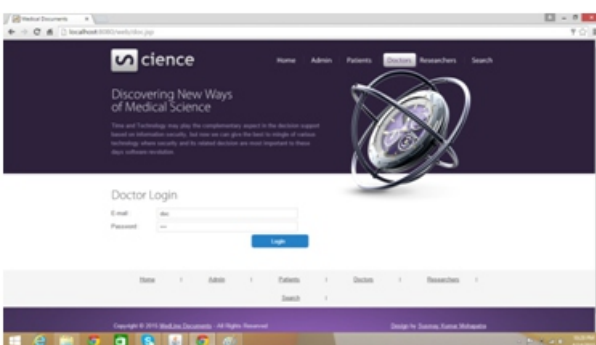
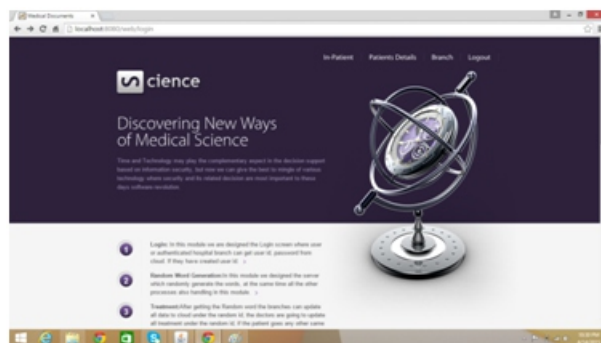
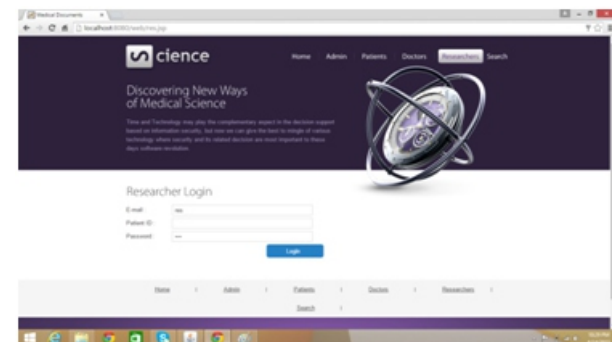
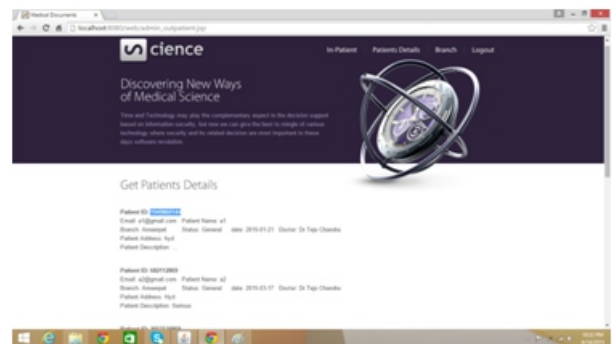
1. Get the url for service given by the PubMed.
2. Right click on solution Explorer. Click add Service Reference.
3. Paste the url taken from web browser or the service url of PubMed
4. Click on go Button and in the namespace textbox, change the name as eUtils.
5. Now the proxy of service will get added in project. By using that proxy, we can call all the methods needed to retrieve the Biomedical Documents.

Proposed System Advantages:

1. Proposed system made the most of the noisy constraints to improve the clustering performance.
2. It was viewed that ML constraints were highly powerful and CL constraints were very promising.

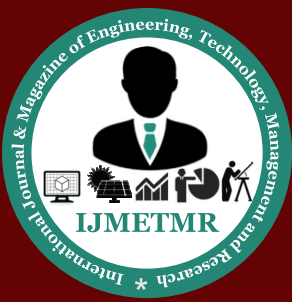
CONCLUSIONS:

We have presented a semi supervised spectral clustering method, which can incorporate both ML and CL constraints, for integrating different information for biomedical document clustering. We have emphasized that our idea behind this project is to incorporate different type of similarities, i.e., the LC, MS and GC similarities. Semi-supervised clustering realizes this new idea, providing a more flexible framework than a method of linearly combining different similarities.



References:

- [1] Jun Gu, Wei Feng, Member, IEEE, JiaZeng, Hiroshi Mamitsuka, and Shanfeng Zhu, Efficient Semisupervised MEDLINE Document Clustering With MeSH-Semantic and Global-Content Constraints, IEEE TRANSACTIONS ON CYBERNETICS, VOL. 43, NO. 4, AUGUST 2013



- [2]. M. Krallinger, A. Valencia, and L. Hirschman, "Linking genes to literature: Text mining, information extraction, and retrieval applications for biology," *Genome Biol.*, vol. 9, no. S2, pp. S8–S14, Sep. 2008.
- [3]. D. Saravanan, Dr. S. Srinivasan, "Matrix Based Indexing Technique for Video Data", *International journal of Computer Science*, 9(5): 534-542, 2013, pp 534-542.
- [4] Bader Aljaber, Nicola Stoke, James Bailey and Jian Pei, "Document clustering of scientific texts using citation contexts," 2009 .
- [5]. A. Rzhetsky, M. Seringhaus, and M. Gerstein, "Seeking a new biology through text mining," *Cell*, vol. 134, no. 1, pp. 9–13, Jul. 2008.
- [6]. R. Baeza-Yates and B. Ribeiro-Neto, *Modern Information Retrieval*. Reading, MA: Addison-Wesley, 1999. Number 9, July 2011. Pages 569.
- [7]. M. Lee, W. Wang, and H. Yu, "Exploring supervised and unsupervised methods to detect topics in biomedical text," *BMC Bioinformat.*, vol. 7, no. 1, p. 140, Mar. 2006.
- [8]. G. Salton and M. McGill, *Introduction to Modern Information Retrieval*. New York: McGraw-Hill, 1983.
- [9]. J. Lin and W. Wilbur, "PubMed related articles: A probabilistic topic based model for content similarity," *BMC Bioinformat.*, vol. 8, no. 1, p. 423, Oct. 2007.
- [10]. T. Theodosiou, N. Darzentas, L. Angelis, and C. Ouzounis, "PuReDMCL: A graph-based PubMed document clustering methodology," *Bioinformatics*, vol. 24, no. 17, pp. 1935–1941, Sep. 2008.