

MULTILINGUAL INFORMATION RETRIEVAL BASED ON KNOWLEDGE CREATION TECHNIQUES

Archana.M¹ and Dr. Sumithra Devi K.A²

¹ Asst.Professor, Dept of MCA, RV College of Engineering, Bangalore-59

E-mail: archanams_m@yahoo.com

² Director, Dept of MCA, RV College of Engineering, Bangalore-59

E-mail: sumithraka@gmail.com

ABSTRACT

As the information access across languages increases, the importance of a system that supports query-based searching with the presence of multilingual also grows. Gathering the information in different natural language is the most difficult task, which requires huge resources like database and digital libraries. Cross language information retrieval (CLIR) enables to search in multilingual document collections using the native language which can be supported by the different data mining techniques. This paper deals with various data mining techniques that can be used for solving the problems encountered in CLIR.

KEYWORDS

Information Retrieval, Cross Language Information Retrieval (CLIR), Data Mining

1. INTRODUCTION

In conventional world, information retrieval was mainly concerned with indexing the terms and search for the useful documents. Nowadays, research in IR includes modeling, document classification, categorization, search engines, user interfaces, data visualization, information filtering, natural language processing or query language and systems architecture. Also from the digital resource perspective, IR research includes text mining, multimedia retrieval, and digital libraries. The representation of IR system in different principles is given below:

| | | | |
|--|---|---|---|
| | Finding answers and information that already exist in a system | | Creating answers and new information by analysis and inference – based on query |
| | Search by navigation (following links, as in a subject directory and the Web generally) | Search by query (as in Google) | |
| Unstructured information (text, images, sound) | Hypermedia systems (Many small units, such as paragraphs and single images, tied together by links) | IR systems (Often dealing with whole documents, such as books and journal articles) | |
| Structured information | | Database management systems (DBMS) | Data analysis systems Expert systems |

Figure 1: Information Retrieval System ^[14]

The rapid growth of communication technologies has immense impact on information retrieval (IR) technique, which has allowed people worldwide to access previously unavailable information. With these advances, however, it has become clear that there is a growing need for retrieval of information in many languages. IR helps the user get useful information from digital resources including digital libraries, WWW and documents.

Some of the problems encountered while retrieving the document that gave rise to the cross language information retrieval are [1]:

- i. The collection of document in different languages, where query formulation for each language would be extremely inefficient.
- ii. Documents that contains text in different language (more than one language).
- iii. User is not able to write the query apart from the native language, but able to make use of documents retrieved in different language that contains images or names – not requiring much of the efforts.

Cross-Language information retrieval would be very useful in the fields of research, where lot of information can be accessed that are present in different language related to the required query. The main goal of cross language information retrieval (CLIR) is exploring the document in the foreign languages. It tries to identify relevant documents in different language from that of the query. Simple knowledge structures such as bilingual term lists have proven to be a useful basis for bridging the language gap. Since the queries and documents are in different language, they have to be translated before they are matched.

The barrier of the CLIR is what should be translated:

- i. Query may be translated
- ii. Document may be translated
- iii. Both query and the document may be translated

Knowledge Creation (Data mining) is a process of automatically discovering useful information in the large data set. Data mining techniques are deployed to search large databases in order to find useful patterns. It is classified into two types:

- i. Direct – target field in explained in this method. The best examples of this method are classification, estimation and predication.

- ii. Undirected – without the particular target field, this method try to fine the pattern or similarities among the groups and best examples are affinity

Information Retrieval and Data Mining are technologies for searching, analyzing and automatically organizing text documents, multi-media documents, and structured or semi structured data.

2. CLIR APPROACHES

Some of the basic approaches to CLIR are Machine translation, controlled vocabulary and dictionary-based approach

2.1 Machine Translation

One of the common methodologies to CLIR is the use of machine translation (MT). This approach can be implemented in two ways:

- i. To translate all available document in a foreign language into the language of user query. This translation has to be done beforehand.
- ii. The query can be translated from the language of the source to another language for search and translate the result back into the source language for viewing.

The only snag with MT system is that it often makes translation errors because of missing information in the term index or ambiguous definitions. It can only produces high quality translations for specific domains, such as those containing specific technical terminology, possibly because semantic accuracy suffers when insufficient domain knowledge is incorporated into a translation system.

2.2 Controlled Vocabulary

The controlled vocabulary is the most traditional approach to CLIR. It is mainly used for indexing and retrieval. In this method a documentlist selects for each document a few descriptors taken from a closed list of authorized terms. A multilingual thesaurus of some sort is created to hold a list of descriptors for each document in a collection and the semantic relations between them, and each term in the thesaurus must be translated for each language involved. The descriptors can be added to the thesaurus manually or automatically if the system can learn from previous indexing which terms are likely to be important. The hitch in this approach is that a query must be generated using only vocabulary from the thesaurus, in which case it may be difficult to search for specific terms that are not included. Larger the size of the vocabulary in the thesaurus, the less effective it becomes.

2.3 Dictionaries-Based

Dictionary Based method can be divided into four logical steps [12]:

- i. Pre translation query modification
- ii. Dictionary lookup
- iii. Equivalent selection and weighting
- iv. Post-translation query translation modification

Each term (semantic unit, single word or a phrase) in the user query is looked up in the machine-readable bilingual dictionary. Some form of ambiguity resolution or equivalent selection is applied to pick the best translation of that term from the list given by the dictionary. This translation is then added to the document language semantic mapping of the bag of words query

and then matching against the document collection as if it had been directly derived from the initial user request.

The advantage of applying this method is that, dictionaries and wordlists covering a wide range of subject areas and language pairs are readily available. In addition, the time needed to implement and set up a dictionary-based system from a printed or electronic source is considerably less comparative to an MT engine for a new language pair. A machine-readable bilingual dictionary can be considered as a data structure which contains a list of dictionary entries for a given set of terms, and a lookup mechanism which, given a source-language query term, consults this data structure to obtain a bag of one or more possible translations or equivalents of the term in question. An entry in a machine-readable bilingual dictionary is a data structure within a dictionary containing all of the necessary information for a given spelling of a source-language query term.

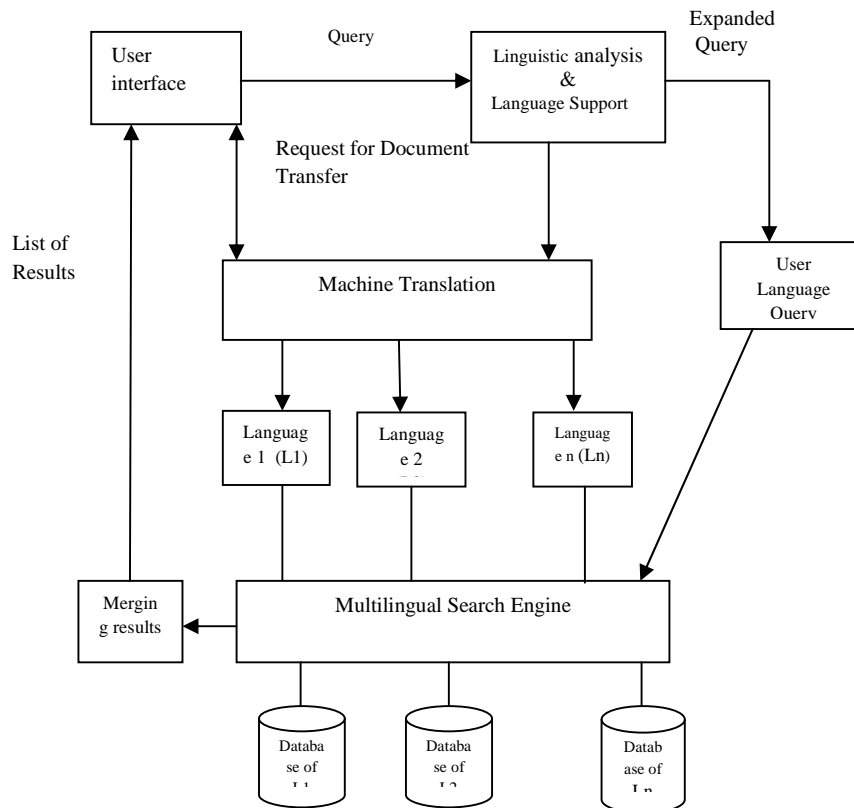


Figure2: Cross Language Information Retrieval System

3. KNOWLEDGE CREATION TECHNIQUES FOR CLIR

Cross language information retrieval involves basically three problems [3]:

- i. Crossing the language barrier i.e. Find the way to translate the term expressed in one language might be written in another.
- ii. Determining which of the translations method to be choose. Selecting more than one translation methods helps in recall.

iii. Deciding the proper weight, if more than one translation is chosen.

Data mining automates the process of sifting through historical data in order to discover new information. Data mining techniques can yield the benefits of automation on existing software and hardware platforms to enhance the value of existing information resources, and can be implemented on new products and systems as they are brought on-line. The core components of data mining technology have been under development for decades, in research areas such as statistics, artificial intelligence, and machine learning.

A data mining operation is achieved using one of a number of techniques or methods. Each technique can itself be implemented in different ways, using a variety of algorithms. Some of the classification algorithms that can be used on cross-language information retrieval platform are:

- (i) Neural Networks
- (ii) Decision Tree
- (iii) K-nearest neighbor
- (iv) Naïve Bayesian
- (v) Cluster analysis.

3.1 Neural Networks

Neural networks consist of links, nodes and these nodes are consisting of output values and input values. Neural network technique can be explained in terms of layered subsystem, where data received in multiple phases can be fed to the next layer as a single phase.

When a query is issued by the user, the subsystem processes it and assigns a keyword to it. Intermediate subsystem indexes the given query and compares with document index. Based on the comparisons, the database system retrieves the relevant documents as a result [4].

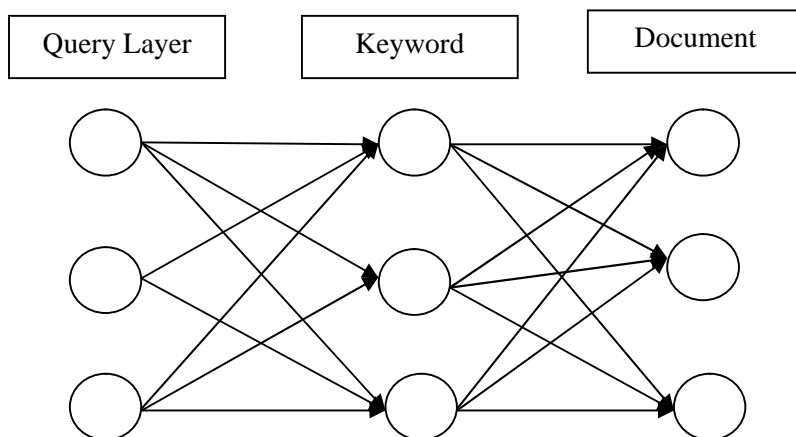


Figure3: Three layer neural network techniques for CLIR [4]

3.2 Decision Tree

A decision tree is a predictive model in which each branch can be viewed as classification of questions and leaves as partitions of the dataset with their classification. When a given object is subjected to a series of tests, in which the outcome contains class label to which object it has to be associated. In a decision tree, branch (non terminal) nodes are tests and leaf (terminal) nodes are class labels. Each branch node has number of child nodes [7].

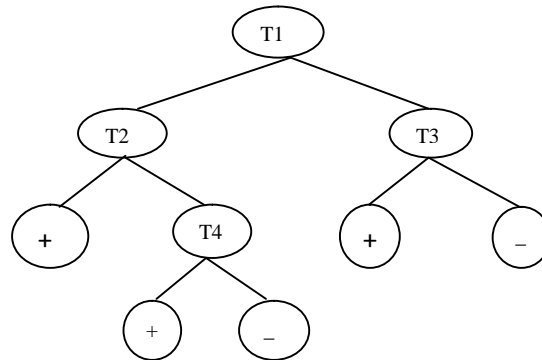


Figure4: Decision Tree for CLIR ^[7]

Classifications of an object in the decision tree begin from the root node and associated testes are applied. Depending on the obtained result an arc is traversed to the appropriate child node. If the subsequent node is a branch node, then its associated test is applied and arc is again traversed. This iterative process is continued until a leaf node is reached.

3.3 K-nearest neighbor

In K-nearest neighbor approach given a test document d, the system finds the K-nearest neighbors among training documents, and weight is assigned to the candidates using their classes. The total weight of the class is taken as the similarity score of each nearest neighbor document to the test document. If number of documents sharing the class is more in number, then per-neighbor weight of that class is added and the resulted weight is used as the likelihood score of that class with reference to the test document. A rank list can be obtained by sorting the scores of candidate classes. The decision rule in KNN classification is given as [10]:

$$\text{score}(d, c_i) = \sum_{d_j \in \text{KNN}(d)} \text{Sim}(d, d_j) \delta(d_j, c_i)$$

Where

- i. KNN (d) indicates the set of K-nearest neighbors of document d.
- ii. (d_j, C_i) is the classification for document d_j with respect to class C_i , that is

$$\delta(d_j, c_i) = \begin{cases} 1 & d_j \in c_i \\ 0 & d_j \notin c_i \end{cases}$$

For test document d, it should be assigned the class that has the highest resulting weighted sum.

3.4 Naïve Bayesian

Naïve-Bayes technique is the organization of both predictive and descriptive. The conditional probability for each relationship is derived by analyzing the relationship between the dependent and independent variables. To generate a classification model only one pass through the training set makes this concept more efficient data mining technique. One of the drawbacks is that, it will not handle the continuous data, so independent or dependent variables that contain continuous values must be cased.

Prior probability is calculated by counting the number of occurrences of the dependent variable in the training dataset. Apart from this, naive-bayes is used to compute how frequently each independent variable value occurs in combination with each dependent variable value. These frequencies are then used to compute conditional probabilities that are combined with the prior probability to make the predictions.

3.5 Cluster Analysis

Clustering or Cluster analysis classifies the given data into groups (clusters) that are meaningful and useful. When a query is placed to the search engine, retrieved result may contain numerous pages. Clustering can be used to group these search results into a small number of clusters, each of which captures a particular aspect of the query. For instance, a query of “newspaper” might return web pages grouped into categories such as national new, international news, sports news, advertisements and entertainments. Each category (cluster) can be broken into subcategories (sub-clusters), producing a hierarchical structure that further assists a user’s exploration of the query results [13].

Figure 5: Clustering analysis for CLIR^[13]



The following table [1] summarizes different types of algorithms that can be implemented for CLIR

| | Techniques | Algorithms |
|---|--------------------|---|
| 1 | Neural Networks | i. Perceptron ii. Backpropagation |
| 2 | Decision Tree | i. Quick reduct ii. Rough set based decision tree ensemble (RSDTE) |
| 3 | K-nearest neighbor | i. Neighbor-weighted K-nearest neighbor |
| 4 | Naïve Bayesian | i. Naïve Bayes Metiore (NBM) ii. Weighted Naïve Bayes Metiore (WNBM) |
| 5 | Cluster analysis | i. Hierarchical Agglomerative ii. Clustering without a recomputed matrix |

Table 1: List of algorithms for different techniques

4. CONCLUSION

The methodology described above makes it possible to bridge the gap between the query and document language. It also handles the issues that are applicable to the CLIR, such as translation ambiguity, lack of translation resource and untranslated terms.

Neural networks can be viewed as a subset of different layer, where each layer has its own set of contribution to the next layer. Decision tree is viewed as a tree with non terminal nodes as branches and terminal nodes as leaf. In K-nearest neighbor, a weight is assigned to each neighbor and depending on the total weight the related document is selected. Naïve Bayesian is organized around predictive and descriptive concept. Finally, Clustering view data as the groups. Perhaps this paper has looked only at CLIR in terms of data mining techniques.

There are variations of this being researched as well, such as elaborating the study on the algorithms of each technique with their pro and cons.

5. REFERENCES

- [1]. Oard and Dorr, "A Survey of Multilingual Text Retrieval", *Technical Report UMIACS-TR-9619*, University of Maryland, 1996
- [2]. Douglas W. Oard, "Alternative Approaches for Cross-Language Text Retrieval," in *Cross-Language Text and Speech Retrieval*, AAAI Technical Report SS-97-05.
- [3]. Grefenstette, G., *Cross-Language Information Retrieval*. 1998, Boston: Kluwer Academic Publishers. 2000.
- [4]. Igor Mokris, Lenka Skovajsova " Neural network model of system for information retrieval from text documents in Slovak language". *Acta Electrotechnica et informatica* no-3 vol 5,2005
- [5]. Ballesteros, L., & Croft, W. B. (1997). Phrasal translation and query expansion techniques for cross-language information retrieval. In *Proceedings of the 20th ACM SIGIR conference on research and development in information retrieval* (pp. 84–91).
- [6]. Chu, H. (2003). *Information representation and retrieval in the digital age*. Medford, NJ: Information Today.
- [7]. Ronnie Fanguy, Vijay Raghavan " Generating rule-based tree from decision tree for concept-based information retrieval", *wss03 applications, products and services of web-based supported systems*.
- [8]. Soergel, D. (1985). *Organizing information: Principles of database and retrieval systems*. Orlando, FL: Academic Press.
- [9]. Blair, D. C. (1990). *Language and representation in information retrieval*. Amsterdam: Elsevier Science.
- [10]. Songbo Tan (2005). *Neighbor- Weighted K-nearest neighbor for unbalanced text corpus*: Elsevier Science.
- [11]. Sigrid Roehling. *Cross-Language information Retrieval*.
- [12]. Gina-Anne Levow, Douglas W Oard, Philip Resnik " Dictionary-based techniques for cross-language information retrieval" *Elsevier, information processing and management* 41 (2005) 523-547
- [13]. *Cluster Analysis: Basic Concepts and Algorithms*. Page no-491.
- [14]. *Information Retrieval the scope of IR. HCIEncyclopediaIRShortEForDS*. Page no-1.

Authors

Ms.Archana.M, is presently working as Assistant Professor in Department of Master of Computer Application, R.V.College of Engineering, Bangalore, India, She has completed her MCA from Gulbarga University, MPhil from Alagappa University and has 5 years of teaching experiences. She has published papers in national and international Conferences. She is also a member of IAENG.



Dr.SumithraDevi K.A, is serving as a Director of Master of Computer Applications Department at R.V.Collage of Engineering, Bangalore, India, she earned PhD in Computer Science and Engineering from the Avinashilingam University for Women, Coimbatore, INDIA and she draws a strong back ground in VLSI Partitioning CAD Tool. She has been awarded by many national and international awards. She is the lifetime member of IEEE, WIE, ISTE and CSI. She is a registered PhD guide under VTU and is guiding two research students and one research candidate under the funded project.

