



Semantic Enrichment of Short Texts through Conducive Formula

Y Krishna Priya

Department of Computer Science and System
Engineering,
AU College of Engineering,
Visakhapatnam, Andhra Pradesh – 530003, India.

Dr.K.Venkata Ramana

Department of Computer Science and System
Engineering,
AU College of Engineering,
Visakhapatnam, Andhra Pradesh – 530003, India.

Abstract:

Web search queries are regularly short and equivocal. Accurate current categorization of user queries regard increased efficiency, quickness, and return potential in general-purpose web search systems. To characterize these queries into certain target classes is a difficult task. Such categorization becomes critical if the system is to return results not just from a general web collection but from topic-specific databases. Semantic analysis is used for the semantic word collections. For each word co-occurring and conceptualized terms are defined. Random forest algorithm is used for semantic hashing of the queries. The efficiency of retrieving the result is more efficient than the traditional algorithms.

Key words:

Short text, Semantic hashing, Semantic frame ,
Random forest

1. Introduction:

Data mining, the extraction of hidden predictive information from large databases, is a powerful new technology with great potential to help companies focus on the most important information in their data warehouses. Data mining tools predict future trends and behaviours, allowing businesses to make proactive, knowledge-driven decisions [1][2]. The automated, prospective analyses offered by data mining move beyond the analyses of past events provided by retrospective tools typical of decision support systems. Data mining techniques can be implemented rapidly on existing software and hardware platforms to enhance the value of existing information resources, and can be integrated with new products and systems as they are brought on-line.

Semantic metadata enables content intelligence by extracting domain-specific entities and concepts from content and relating them in a meaningful way to identify related content and facilitate intelligent answers to user queries. Semantic enrichment services from scope aim to enhance content by adding contextual information by tagging, categorizing and classifying [7] data in relationship to each other and other base reference sources. For ordinary words such as “cat”, WordNet contains detailed information about its various senses. However, much of the knowledge is of linguistic value, and is rarely evoked in daily usage [3]. For example, the sense of “cat” as gossip or woman is rarely encountered.

2. Related Work:

Web search queries are typically short and ambiguous. To know the effectiveness of the random forest algorithm on query enrichment we follow different steps. We perform an information retrieval experiment on particular data and perform classification based on random forest semantic approach. Based on the idea from existing system Query enrichment [3], which takes a short query and maps it to the intermediate objects. Based on the collected intermediate objects, the query is then mapped to the target categories. To build the necessary mapping functions, we use an ensemble of search engines to produce an enrichment of the queries. The semantic hashing approach encodes the meaning of a text into a compact binary code.

Cite this article as: Y.Krishna Priya & Dr.K.Venkata Ramana, "Semantic Enrichment of Short Texts through Conducive Formula", International Journal of Research in Advanced Computer Science Engineering, Volume 3 Issue 6, 2017, Page 25-29.

Thus, to tell if two texts have similar meanings, we only need to check if they have similar codes. Majorly the meaning of the text that are searched is same but the code is not same for all the text of same meaning. To cluster short texts by their meanings, we propose to add more semantic signals to short texts. The data is derived from different articles such as Wikipedia, open directory projects etc. The existing system uses the deep neural network semantic hashing which is better algorithm with precision and recall. The random forest design for large amount of datasets with best precision and recall than deep neural network [9]. The comparison between the random forest neural network and TFIDF traditional neural network is RF makes computation fast than other neural networks.

3. Preliminaries:

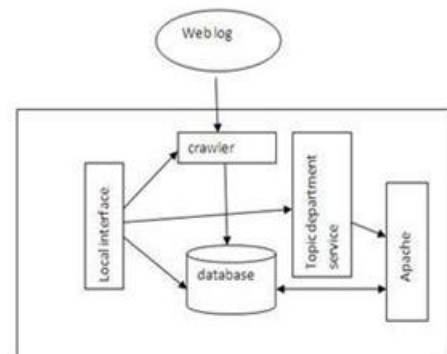
Existing novel approach algorithms are mostly based to classify the queries into certain target categories. The short query classification problem is not as well-formed as other classification problems such as text classification. The difficulties include short and ambiguous queries and the lack of training data. Query enrichment, which takes a short query and maps it to the intermediate objects. Based on the collected intermediate objects, the query is then mapped to the target categories [3]. As short texts do not provide sufficient term co-occurrence information, traditional text representation methods have several limitations when directly applied to short text tasks. Measuring the semantic similarity between two texts has been studied extensively in the IR and NLP communities [2]. However, the problem of assessing the similarity between two short text segments poses new challenges. Text segments commonly found in these tasks range from a single word to a dozen words. Because of the short length, the text segments do not provide enough contexts for surface matching methods such as computing the cosine score of the two text segments to be effective. Salakhutdinov and Hinton [5] proposed a semantic hashing model based on Restricted Boltzmann Machines (RBMs) for long

documents, and the experiments showed that their model achieved comparable accuracy with the traditional methods, including Latent Semantic Analysis (LSA) and TF-IDF. Context is another problem when measuring the similarity between two short segments of text [8]. While a document provides a reasonable amount of text to infer the contextual meaning of a term, a short segment of text only provides a limited context.

4. Proposed system:

Clustering of short texts, such as snippets, presents great challenges in existing aggregated search techniques due to the problem of data sparseness and the complex semantics of natural language. The proposed method employs a hierarchical structure to tackle the data scarcity problems. A novel approach is proposed for understanding short texts. A semantic network based approach for enriching a short text .We present a novel mechanism to semantically enrich short texts with both concepts and co-occurring terms, such external knowledge’s are inferred from a large scale probabilistic knowledge base using our proposed thorough methods [6]. For each autoencoder we design a specific and effective learning strategy to capture useful features from input data. We provide a way to combine knowledge information and random forest network for text analysis, so that it helps machines better understand short texts.

4.1 System Architecture:



4.2 Algorithm:

A random forest works the following way:

1. First, it uses the Bagging (Bootstrap Aggregating) algorithm to create random samples. Given a data set D_1 (n rows and p columns), it creates a new dataset (D_2) by sampling n cases at random with replacement from the original data. About $1/3$ of the rows from D_1 are left out, known as Out of Bag(OOB) samples.

2. Then, the model trains on D_2 . OOB sample is used to determine unbiased estimate of the error.

3. Out of p columns, $P \ll p$ columns are selected at each node in the data set. The P columns are selected at random. Usually, the default choice of P is $p/3$ for regression tree and P is \sqrt{p} for classification tree.

4. Unlike a tree, no pruning takes place in random forest; i.e. each tree is grown fully. In decision trees, pruning is a method to avoid over fitting. Pruning means selecting a subtree that leads to the lowest test error rate. We can use cross validation to determine the test error rate of a subtree. Several trees are grown and the final prediction is obtained by averaging or voting.

5. Enriching short texts:

Given a short text, we first identify the terms that semantic field can recognize, then for each term we perform conceptualization to get its appropriate concepts, and further infer the co-occurring terms. The semantic field text file contain the semantic words for particular word. The semantic words are collected from wordnet or Wikipedia. For example, given "china, India" they can be developing countries or Asian countries same as in electronics the word "TV" specifies different TV sets or different TV companies. We denote this two stage enrichment mechanism as CACT (concepts-and-co-occurring) [1]. After that, a short text can be represented by a vector of term count and fed to our RF model to do semantic hashing.

5.1 Pre-processing:

Pre-processing method plays a very important role in text mining techniques and applications. It is the first step in the text mining process. Extraction is the method used to tokenize the search content into individual words. Stop words are removed from documents because those words are not measured as keywords. Stop words are removed by classical method in pre-processing we treat the semantic field as a dictionary of terms [10].

5.2 Conceptualization:

Consider the true sense of a term is heavily affected by its neighbours, especially for the ambiguous ones, we propose a multiple mechanism to infer. The most appropriate sense for each term in a short text. The most important mechanism we take advantage of is context dependent conceptualization, which uses a probabilistic topic model. Through this method , we first get of a short text s . Let $\sim s$ be the sequence of term indices of s and $\sim z$ be the topic assignment vector of $\sim s$. We then compute the probability of concept c given a term w of s based on the topic distribution where w and c are indices of instance term and concept them respectively, is the typicality of concept c given term w , is the probability of term w given topic k .

Term	Canonical concept	Member concept
phone	device	Mobile device, portable devices, smart device etc.
samsung	company	Corporation, firm, stock, technology

5.3 Co-occurring terms:

We define the co-occurrence score to measure the probability of one term o that co-occurs with a target term t in a short text s . represents the normal co-occurrence probability, which is pre-defined as in

semantic text file, and measures the semantic similarity between o and t under the text s , since we already know the concept c of term t , we aim to make the co-occurring term o have consistent semantic with where c_i represents a concepts of term o and is weight parameter. For example, consider TV in the short text “TV”, video projector, monitors” in different companies would have high normal co-occurring probabilities, but only monitors are appropriate co-occurring terms if we take into consideration the semantic(concept) of “TV” in that text.

$$S(o|t,s) = \alpha \quad (o|t) + (1-\alpha) \quad (o|t,s)$$

Where α is a meta parameter that explores the trade-off between the two component scores. $(o|t)$ is the co-occurrence probability between instance terms o and t . $(o|t,s)$ measures the semantic similarity between o and t .

6. Experiment setup:

To know the effectiveness of the random forest algorithm on query enrichment we follow different steps. The term “electronics” have been chosen for this project. The electronic term has different categories : different company, different products, different prices etc. The data is stored in two formats images and files. Data of some magnitude is derived for each category. In the experiment for every noun term in a query, we enrich it with the top possible concepts and co-occurring terms. As a result, the representation of a query is greatly enriched.

7. Result analysis:

Different experiments are carried out on tasks including information retrieval and classification for short texts. Significant improvements over existing approaches, which confirm that concepts and co-occurring terms effectively enrich short texts, and enable a better understanding of them are determined.

RF based model is able to capture the abstract features and complex correlations from the input text such that the learned compact binary codes can be used to represent the meaning of the text. The duration comparison between the traditional approach TFID and random forest is shown for particular search query.

8. Conclusions:

Query classification is an important as well as a difficult problem in the field of information retrieval. Once the category information for a query is known, a search engine can be more effective and can return more representative Web pages to the users. To solve the query classification problem, we designed an approach based on query enrichment which can map the queries to some intermediate objects. We propose a novel approach for understanding short texts. First, we introduce a mechanism to enrich short texts with concepts and co-occurring terms that are extracted from a probabilistic semantic network. Then we approach random forest algorithm to do semantic hashing. We carry out comprehensive experiments on short text centred tasks including information retrieval and classification. The significant improvements shows that the information retrieval takes less time than the traditional approaches like TFID.

9. Future scope:

There are a number of interesting extensions of this work. As this work is for aggregated search, the efficiency of the whole framework should be optimized for real applications. Moreover, the ranking procedure can be implemented for the frame of work between the random forest algorithm and traditional approach algorithms.

References:

[1]Zheng Yu, Haixun Wang, Xuemin Lin, Senior Member, IEEE, and Min Wang, “Understanding Short Texts through Semantic Enrichment and Hashing”,



iee transactions on knowledge and data engineering,
vol. 28, no. 2, february 2016.

[2] M. Sahami and T. D. Heilman, "A web-based kernel function for measuring the similarity of short text snippets," in Proc. 15th Int. Conf. World Wide Web, 2006, pp. 377–386.

[3] W. tau Yih and C. Meek, "Improving similarity measures for short segments of text," in Proc. 22nd Nat. Conf. Artif. Intell., 2007, pp. 1489–1494.

[4] D. Shen, R. Pan, J.-T. Sun, J. J. Pan, K. Wu, J. Yin, and Q. Yang, "Query enrichment for web-query classification," ACM Trans. Inf. Syst., vol. 24, no. 3, pp. 320–352, 2006.

[5] R. Salakhutdinov and G. E. Hinton, "Semantic hashing," Int.J.Approx. Reasoning, vol. 50, no. 7, pp. 969–978, 2009.

[6] S. Banerjee, K. Ramanathan, and A. Gupta, "Clustering short texts using wikipedia," in Proc. 30th Annu. Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval, 2007, pp. 787–788.

[7] W. Wu, H. Li, H. Wang, and K. Q. Zhu, "Probase: A probabilistic taxonomy for text understanding," in Proc. Int. Conf. Manage. Data, 2012, pp. 481–492.

[8] Y. Song, H. Wang, Z. Wang, H. Li, and W. Chen, "Short text conceptualization using a probabilistic knowledge base," in Proc. 22nd Int. Joint Conf. Artif. Intell., 2011, pp. 2330–2336.

[9] D. Kim, H. Wang, and A. H. Oh, "Context-dependent conceptualization," in Proc. 23rd Int. Joint Conf. Artif. Intell., 2013, pp. 2654–2661.

[10] B. Stein, "Principles of hash-based text retrieval," in Proc. ACM 30th Annu. Int. Conf. Res. Develop. Inf. Retrieval, 2007, pp. 527–53.