

# Analysis of different techniques for information retrieval from high dimensional Unstructured Text Data

**Ankit J. Faldu\* & Dr. Kishor H. Atkotiya\*\***

\*Author, Assistant Professor-Faculty of Science-MCA, Atmiya University.

\*\*Corresponding Author, Professor-Department of statistic, Saurashtra University.

Received: January 28, 2019

Accepted: March 02, 2019

**ABSTRACT:** *Text analysis (TA) is a trending research area that used to identify helpful information that could derived from large amount of high dimensional unstructured text data by using difference method from very trending technology that is machine learning, data mining, NLP means natural language processing and information retrieval (IR). Text analysis includes the preliminary processing method to collecting data and early aware of the different patterns and attributes that present in the collecting data from different resources. The different Methods like "Information Retrieval (IR)", "Information Extraction (IE)", "Categorization of data", "Clustering and Summarization" of data that are used to analyze these in-between representations like "distribution analysis (DA)", "association rules (DR)" and "visualization". Text analysis helps to find valuable information from customer email, documentation, comments, social network etc.*

**Key Words:** *Big Data, Unstructured Data, Unstructured Text Data, Algorithms, Information retrieval.*

## 1. Introduction:

Nowadays a most of business organizations collect and accumulate massive amounts of new information in data warehouses and cloud platforms and this unstructured data remains to raise exponentially by the minute as new information comes driving in from several sources. As outcome, it becomes a great issues and challenge for corporate organizations to store, process, and analyze massive amounts of unstructured textual data with currently available traditional tools. This is the main reason that Text analysis comes in the picture. The text analysis obtains the helpful information from different data sources like emails, html, comments, documents etc. through the elaboration and identifications of exciting patterns and unique design in data. So Natural processing language is used to translate unstructured data into structured data. Text analysis is used to discovery useful information for this data. Text analysis combines and assimilates the tools of statistics, data mining, information retrieval (IR), machine learning, and computational linguistics. Text analysis also pacts with Natural language processing (NLP) texts that is stored in semi-structured or unstructured formats. Text is an important intermediate for exchange different type of information. The aim of text analysis applications is to trace data, retrieve data and operate data of applicable information proficiently from the large amount of bulk of text which

continues to grow expeditiously and exponentially.

## 2. Text-Analysis Pre-Processing method:

Pre-processing method are dynamic in TA and its different applications. It is also known as cleaning the text, which is used to identify bug and eliminate bugs in the given large amount of unstructured text. Different pre-processing methods are available to make the text to be prepared for more analysis process to accomplish notable performance [1]. All Most of the functionalities of pre-processing methods work with Natural language processing (NLP). The succeeding are some key pre-processing techniques: Word Sense Disambiguation (WSD), Parts-of-Speech (POS) Tagging, Tokenization, Grammatical Parsing and Chunking, Stop Word Filtering and Chunking, Text Summarization, Term Frequency and Inverse Document Frequency, Lemmatization and Stemming.

## 3. Information Retrieval (IR)

IR is a procedure of mining appropriate and related patterns base on a given large set of different words or expressions. It's very strong relationship in text analysis and information retrieval form large amount of unstructured text data. Information retrieval systems used various algorithms that are helpful to analysis and find pattern the user's comportment and find appropriate data [2]. Most of search engines are more or less using information retrieval (IR)

system regularly to find appropriate information and documents according to input text. That's way search engines used query related different algorithms to find the more relevant output. [3]

#### **4. Information Extraction (IE)**

Information extraction (IE) is used of repeatedly taking out structured data from verity of unstructured machine-readable documents. This activity focus on processing human language under stable texts by help of natural language processing (NLP). Current activities in hypermedia document processing like programmed comment and content mining out of different unstructured documents could be understood as IE.

The purpose of IE techniques is the to obtaining helpful information from unstructured text. It's used finds the entities, extraction or obtaining of events and relationships between unstructured text data. IE recognizes important phrases and relationships available within given unstructured text data. Information extraction (IE) is anxious with obtained of related data and information from the unstructured text data. [4].

#### **5. Natural Language Processing (NLP)**

NLP focus on automatic analysis and processing of high volumes unstructured text information. NLP does use various types of data analysis techniques like NER means "Named Entity Recognition" for recorded and its alternative word finding and to identified the relationships between words [5]. NER recognize the occurrences of definite entity from different large cluster of various documents. These all the entities and their object permit the finds the association and other info to achieve their important concept. Main disadvantage of NER technique is to absence of whole dictionary for completely named objects used for recognition [5][6]. Query related to the algorithms essential to be used to achieves suitable output. In reality data, a single unit or entity has abundant relations such as freeze and refrigerator. Many times, a collection of successive words have a multiple word names to recognize the boundaries value and solve overlying issues to use classification technique. NER systems have accomplished significance upto accurate result between 75 % to 85% [7].

To abstract alternative expression and abbreviation from unstructured text data, coreferencing techniques are regularly taking for used. Natural Languages Processing (NLP) its self to much of difficulties as an unstructured text mined from different resources do not have similar words or abbreviation. So, it is an essential

to distinguish such problems and make guidelines or protocol for equals identification [8].

#### **6. TEXT CATEGORIZATION**

Aim of text categorization is to categorize large set of verities documents in persistent number of advance established categories. Every text file, document or article relating to one or multiple class. Responsibility of categorization is organized a given unstructured data object to predefine set of categories. TC is one kind of "supervised" learning approach. When categorizing any document an algorithm will regularly consider as a "Bag of Words". It is a group of unstructured text file, article and documents, the procedure of discovery the perfect result or result for every document. Automated text categorization techniques used mix of contexts. Further some important methods of textual analysis were presented with ordinary text analysis process to expand the significance and correctness of results from unstructured data [9].

To categorize the various text documents one of method is weighted heuristics which extracts useful information and features by given listed some specific rules. fixed phrase, Sentence length, related word, paragraph. Those methods are implementing and analyzing for text categorization. Text categorization has been used on many file, documents and article at the same time. Quality of classifiers & types of classifiers is totally depending on the theme and nature of the unstructured text data [10].

#### **7. CLUSTERING**

Clustering is most attentive and most important methods in text analysis which is helped to grouping similar type of documents in one group. Its purpose is to identify basic structures available in information, and organize information into relevant sub-groups for additional study of data and analysis of data. It is a non-supervised procedure. With these different instances are arranged into different groups known as a cluster. The problem with clustering is that creating group without label gathering into relevant clusters without any prior information. Any labels related to instances are found from the given unstructured data. More advantage of clustering is that documents can materialize in multiple topic and sub-topics, so certifying that a valuable document must be present to search results. A fundamental clustering algorithm makes vector topics for every valuable data and procedures the weights of how the data fix into every create cluster. Clustering techniques is useful in many real time applications like pattern

recognition, biology, data analysis, image segmentation, document retrieval, pattern classification, web search, security and BI. Frequently used clustering algorithms are the "K-means", "the EM-based mixture resolving", and the HAC. [11]

## 8. VISUALIZATION

Visual content analysis or data visualization places huge unstructured text data sources in a hierarchy order. The data given by form of chart, picture, graph etc. is best, complete and fast processing than plain text-based explanation so it is better for analysis the huge set of documents. Visualization is applicable in broad variety of documents and find related information. Visualization can expand to finding of related useful patterns or related information for text analysis. Information or data that permit a pictorial representation includes ontology, result set, relations of keyword are considered the important elements of the search activities. [12] construction of information visualization conducted into following step.

1. Data assembling
2. Data examination or analysis and extraction
3. Visualization

## 9. ISSUES RELATED TO TEXT ANALYSIS

Various issues face in the text analysis process. Difficulties have rise at in-between of text analysis techniques. Pattern analysis is used to translate unstructured text data into structured but analysis process has generated its own difficulties [13]. Another main problem is a text modification of group is dependent that make problems. Existing some tools that support multiple languages. Numerous algorithms and different techniques do useful for autonomously and support multi language text. Important documents continue external the text analysis process. Several tools and technologies don't support it. [14]. Requirements of the domain, experts' people have been required to work together from precise results [5], [6].

Text's synonyms and antonyms available in the given data generate issues for the process of text analysis tools. It's tough to classify the documents where documents are huge. Abbreviations result gives different meaning in different conditions is large problem [15]. So, it's the necessity to define rules depending to the domain that rules will be used as a standard rule in the required field and shall be integrated in text analysis as a requirement [14], [15]. Natural language processing (NLP) have many problems in text enhancement techniques and the kind of

relationship between entity. Many times, word's spelling is same but meaning different, for example, bank & bank. TA tools considered both words same. Grammatical regulation is still time an open-ended research issue in text analysis [16].

## 10. CONCLUSION

Present of large text data required to be analyzed to find helpful information. Text analysis are used to analyze required information from very huge volumes of unstructured text information. In this research paper we explain text analysis techniques that is helpful to upgrade and expand the text analysis mechanism. Choosing accurate techniques depending to the domain knowledge support to text analysis to easy, understand efficient. expertise of various domain and its integration, multiple text enhancement, varying concepts granularity and uncertainty in NLP are important problems and challenges that stand up during text analysis process. Many efforts are doing on the document analysis using text analysis methods. But text analysis is till the challenging open-ended issue on present world situation. In future, we will try to design an algorithm and that algorithms are help to solve the different issues mention in this work.

## References:

1. Saira, Gillani Andleeb, "From text mining to knowledge mining: An integrated framework of concept extraction and categorization for domain ontology", PhD Dissertation, Budapesti Corvinus Egyetem, 2015.
2. R. Steinberger, "A survey of methods to ease the development of highly multilingual text mining applications," Language Resources and Evaluation, vol. 46, no. 2, pp. 155–176, 2012.
3. Andreas Hotho, Andreas Nürnberger, and Gerhard Paaf, "A Brief Survey of Text Mining. In Ldv Forum", Vol. 20.19–62. 2005.
4. D. Jasmine Guna Sundari, D. Sundar, "A Study of Various Text Mining Techniques", vol 08 issue 02 pg.no 82-85 (2017)
5. B. Laxman and D. Sujatha, "Improved method for pattern discovery in text mining," International Journal of Research in Engineering and Technology, vol. 2, no. 1, pp. 2321–2328, 2013.
6. A. Henriksson, H. Moen, M. Skeppstedt, V. Daudaravičius, and M. Duneld, "Synonym extraction and abbreviation expansion with ensembles of semantic spaces," Journal of biomedical semantics, vol. 5, no. 1, p. 1, 2014.
7. A. M. Cohen and W. R. Hersh, "A survey of current work in biomedical text mining," Briefings in bioinformatics, vol. 6, no. 1, pp. 57–71, 2005.
8. E. A. Calvillo, A. Padilla, J. Muñoz, J. Ponce, and

- J. T. Fernandez, "Searching research papers using clustering and text mining," in Electronics, Communications and Computing (CONIELECOMP), 2013 International Conference on. IEEE, 2013, pp. 78–81.
9. C. P. Chen and C.-Y. Zhang, "Data-intensive applications, challenges, techniques and technologies: A survey on big data," *Information Sciences*, vol. 275, pp. 314–347, 2014.
10. R. Al-Hashemi, "Text summarization extraction system (tses) using extracted keywords." *Int. Arab J. e-Technol.*, vol. 1, no. 4, pp. 164– 168, 2010.
11. B. L. Narayana and S. P. Kumar, "A new clustering technique on text in sentence for text mining," *IJSCEAT*, vol. 3, no. 3, pp. 69–71, 2015.
12. Rashmi Agrawal, Mridula Batra, "A Detailed Study on Text Mining Techniques", IJSCE, ISSN: 2231-2307, Vol. 2, Issue-6, January 2013.
13. A. Henriksson, J. Zhao, H. Dalianis, and H. Boström, "Ensembles of randomized trees using diverse distributed representations of clinical events," *BMC Medical Informatics and Decision Making*, vol. 16, no. 2, p. 69, 2016.
14. A. Kumaran, R. Makin, V. Pattisapu, and S. E. Sharif, "Automatic extraction of synonymy information:-extended abstract," *OTT06*, vol. 1, p. 55, 2007.
15. A. Kaklauskas, M. Seniut, D. Amaratunga, I. Lill, A. Safonov, N. Vatin, J. Cerkauskas, I. Jackute, A. Kuzminske, and L. Peciure, "Text analytics for android project," *Procedia Economics and Finance*, vol. 18, pp. 610–617, 2014.
16. N. Samsudin, M. Puteh, A. R. Hamdan, and M. Z. A. Nazri, "Immune based feature selection for opinion mining," in Proceedings of the World Congress on Engineering, vol. 3, 2013, pp. 3–5.