

# An Experiential Learning Of Ontology-Based Multi-Document Summarization By Removal Summarization Techniques

**Pranjali Yadav Deshmukh**  
Sinhagad Institute of Technology,  
Pune ,Maharashtra, India  
pranu313@gmail.com

**Prof. Rahul Ambekar**  
Sinhagad Institute of Technology,  
PuneMaharashtra, India  
rahulambekar@yahoo.com

**Abstract**—Outstanding development of the Internet alongside the New technologies, for example, high-speed networks and cheap monstrous storage, alongside the, have prompted a gigantic increment in the price and accessibility of digital document. For any person, perusing of this information is huge tedious so need to get to multi-document summarization (MDS) frameworks, which can successfully consolidate data found in a few documents into a short, readable synopsis, or summary. For semantic representation of textual data in ontology region, as a hypothetical model, our framework gives a significant structure. The feasibility of utilizing the ontology as a part of taking care of multi-document summarization issues in the area of disaster administration is finding in propose model. Saliency score is normally assigned out to every sentence and sentences are positioned by score, then the top positioned sentences are chosen as the summary. Regarding the summary magnificence, wide tests on an accumulation of press discharges applicable to "Jammu Kashmir Flood in 2014" documents. Ontology based multi-document summarization systems utilizing "NLP based extraction" outflank different baselines. Our commitment in propose module is to utilize the data extraction strategies (NLP) to improve and enhance summary result.

**Keywords**—Multi - document summarization, Disaster Management, Ontology, Extraction technique, k-means, NLP, Sentence Extraction, and Generic.

## I. INTRODUCTION:

India has been by and large weak against natural pain due to its striking geo-climatic conditions. Surges, dry spells, tornados, tremors and landslide have been irregular phenomena. Around 60% of the landmass is slanted to tremors of diverse intensities; in abundance of 40 million hectares is slanted to surges; around 8% of the total extent is slanted to tornados and 68% of the range is helpless to dry

season to dry season. In the area of disaster administration, over a great many several reports are often discharged by the neighborhood government or local emergency offices amid the disaster, which wrap most occasions material to the disaster also, the time separation will be days to months, contingent upon how brutal the disaster is. The information will be introduced in a arrangement of newswire, containing a considerable measure of routine providing details regarding various parts of the calamity. In such a case, it is greatly troublesome for domain specialists to rapidly discover either the most critical data generally overall

(generic summarization) or the most significant data to a predetermined (query / topic focused summarization).

The appearance of WWW has made a lot of information. The motivational part is the way to summarize considerable information which is show on Internet furthermore to check whether summarization can be utilized as a part of domain of disaster management. Furthermore, additionally to pass on the quintessence of the documents, helps in finding important data rapidly. In this paper, attainability of utilizing ontology in understanding multi-document summarization issue in calamity administration area is investigated. There exists a rich collection of writing on multi-document summarization. Some analysts use the unequivocal ideas inside sentences to address multi-document synopsis. Processing sentence significance in view of the idea of eigenvector centrality (prestige) that we call Lex page rank. In this model, a sentence integration matrix is built in view of cosine likeness. In the event that the cosine similitude between two sentences surpasses a specific predefined threshold, a comparing edge is added to the network. New multi-document summarization system in light of sentence-level semantic investigation and symmetric non-negative matrix factorization. Most exceedingly terrible compute sentence likenesses utilizing semantic investigation and build the comparability matrix. At that point symmetric matrix factorization, which has been indicated to be proportional to standardized unearthly clustering, is utilized to gathering sentences into clusters. At last, the most useful sentences are chosen from every gathering to

frame the summary. To help clients of the Internet discover the data they are searching for rapidly, a productive algorithm for building the summarize of an accumulation of document found by a search engine because of a client query, called DISCO (Distribution Scoring) is proposed. Utilized both connections among sentences and connections between the given queries what's more, the sentences by complex ranking. Likelihood models have likewise been proposed with diverse presumptions on the era procedure of the reports and the query.

Summarization is merging the source content into a shorter variation shielding its information substance and general criticalness. It is incredibly troublesome for people to physically abridge broad documents of substance. Summarization techniques can be approved into extractive and abstractive summary. An extractive Summarization framework includes selecting discriminating sentences, areas et cetera from the first documents and connecting them into shorter structure. The significance of sentences is picked concentrated around numerical furthermore, linguistic quirks of sentences. An abstractive Summarization strategy includes appreciation the first substance and re-telling it know in less words. It uses linguistic systems to examine and decipher the substance and after that to find the new ideas and understandings to best depict it by delivering another shorter message that goes on the most basic information from the first substance report.

The nonexclusive summarization gives a general summary of all the data contained in a report. Answers the question what is this report about? A report is normal to contain a few themes. Fundamental subjects are talked about widely by numerous sentences. Minor themes have less sentence support further more; exist to backing the fundamental subjects. The particular objective of bland synopsis is as per the following: 1) for a given archive with  $n$  sentences, pick  $k$  sentences (as determined by the client) from the record that best portray the principle points of the record. 2) Keep excess of the picked sentences to a least.

The query-based summarization errand is to produce an summary of single/various documents which is centered towards the clients query. All in all, an query score was computed for every sentence in light of the conveyance of query terms and added to its general score got by sentence extraction strategies. The top scoring sentences were utilized as a summary for each of the recovered record. Ventures in inquiry based summarization are distinguishing proof. In query-focused summarization, the data identified with a given subject or query ought to be consolidated into summary. The sentences suiting the clients query ought to be extricated. Utilizing ontology, an experiential investigation of distinctive methodologies for summarization occupation is introduced in this paper.

## II. RELATED WORK:

In generic summarization, every sentence is connected with a saliency score. As showed by situating, sentences are top situated and picked the summary concentrated around the situating result. To analyze the information contained in a report set what's more, to think particularly exceptional sentences into the summary concentrated around syntactic measurable gimmicks [2] [3] [4] [6] [5], unsupervised furthermore directed procedures are proposed. Regardless, most of the current frameworks, the sensible information in the sentence level is ignored. Client's all the more clear results for once-overs can be given by the connected information. To address multi-document rundown [5] [4], a couple of researchers use the unequivocal idea inside sentences like using Wikipedia. After that the sentences are positioned by saliency score. As indicated by positioning, sentences are top positioned furthermore, chose the summary in view of the positioning result. To investigate the data contained in an document set and to concentrate exceedingly notable sentences into the summary based on syntactic or statistical features [8] [9] [10] [11] [12], unsupervised and supervised methods are proposed. In any case, the vast majority of the current methods, the calculated data in the sentence level are overlooked. Clients more coherent results for summaries can be given by the reasonable data. to address multi-document summarization [13] [14], some specialists utilize the explicit concepts inside sentences like utilizing Wikipedia. For domain-specific document summarization undertakings such systems a not be straightforwardly connected as Wikipedia contains such a large number of ideas not identified with a particular domain. Be that as it may, such procedures can't be specifically connected to domain particular document summarization tasks, since Wikipedia contains an excess of ideas not pertinent to a particular domain. In query-based summarization, the data identified with a given subject or query ought to be fused into summaries. The sentences suiting the clients a query ought to be removed. To fuse the query data, different strategies for non specific summarization can be reached out to contain query data. A strong summarization framework grew inside the GATE architecture is proposed in Saggion et al. [7]. It utilizes the powerful parts for semantic labeling and co-reference determination given by GATE. Weiet et al. [15] expressed the query impact into the shared fortification fasten to manage the need for query arranged multi-document summarization. Wan et al. [16] made utilization of both connections among sentences and connections between the given query and the sentences by complex positioning. Likelihood models have likewise been proposed with diverse suppositions on the era methodology of the documents and the query [17][18]. The methodology of enlarging the clients query with extra terms so as to move forward search results are Query development. For example, when client is prepared to pursuit panther by some search engine, we can

extend such query by adding equivalent words of panther to the query, such as jaguar, cougar, and so forth. In the field of document summarization query expansion development is extremely well known. Since here the nature of the created summary can be move forward. Case in point, Daume what's more, Marcu [19] state a defended query extension strategy in the language modeling for IR system. Notwithstanding, it does not consider the semantic relatedness between the sentences what's more, the question string.

### III. IMPLEMENTATION DETAILS:

The fundamental focus of our framework is investigating the practicality of utilizing the ontology as a part of unraveling multi- document summarization issues in the area of fiasco administration of (Jammu Kashmir Flood in 2014). In this paper j and K surge news [21] [22] data are gathered physically and produce text reports for alluding as dataset in this framework. Likewise information extraction systems to further enhance summarization results. In proposed framework information extraction system is utilized that is NLP system. This technique for term acknowledgment utilizing linguistic and statistical strategy, making utilization of logical data to bootstrap learning.

#### A. System Architecture:

Firstly, the ontology into multi-document summarization issues in calamity administration area is utilized. The likelihood of using the ontology is investigated to attain to the objective of diminishing data redundancy. The NLP based information extraction strategies to extra enhance outline results is utilized.

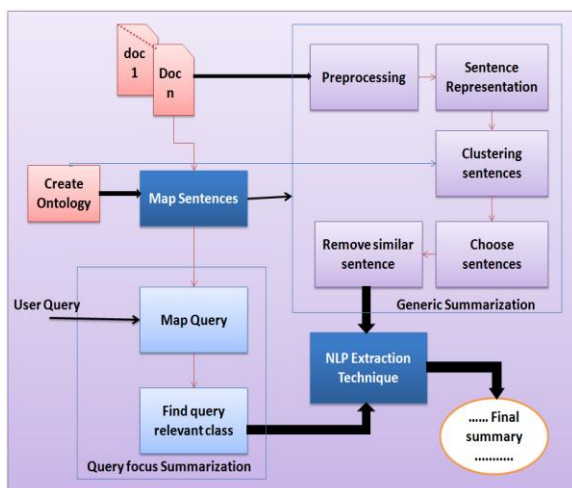


Fig. 1. System Architecture

System Architecture is defined as following way:

- 1) Input: Multiple record of debacle related reports, news i.e. advanced content are use as input.
- 2) Preprocessing: Before handling on printed information we need do some preprocess on input data. Process incorporates tokenization, sentence division, stop words removing, and stemming.

- 3) Create Ontology: Before making of our framework first need to make ontology of our chose space by expert.
- 4) Map sentences: First venture of our framework to utilization of made ontology by mapping sentences in ontology chain of importance. Sentence appoint to its connected hub of ontology.
- 5) Apply clustering sentence algorithm: use standard k mean algorithm on all documents sentences.
- 6) Choose sentence: Clustered sentences are select by utilizing centroid based system and select L sentences from every group.
- 7) Remove redundancy: Similar significance sentences are discovered at keep one and only of them.
- 8) Query map: client inquiry is mapped in made ontology.
- 9) Find inquiry relevant class: discover hub of ontology where inquiry greatest coordinated. Also, select that hub sentences as query relevant sentences.
- 10) NLP based sentence extraction Technique: Apply six features on every sentence and figure out last score of every sentence. Select top k sentence as last summary.
  - a) Term Feature: Word which are all the more regularly happens in sentences are vital so we consider this words for weighting sentence.
  - b) Sentence Position: Sentence position is a sentence area in a passage. We accepted that the first sentence of every passage is the most essential sentence. Hence, we sort the sentence based on its position.
  - c) Sentence centrality: Sentence centrality is the vocabulary cover between this sentence and different sentences in the report.
  - d) Sentence incorporation of name entity: Generally the sentence that contains more proper noun, places or things is essential and it is most likely included in the document summary.
  - e) Sentence inclusion of numerical data: The sentence that contains numerical information is likewise a vital and generally included in last summary so we consider this highlight for scoring sentence.
  - f) Sentence Length : This highlight is utilized to penalize sentences that are too short, following these sentences are most certainly not anticipated that would have a place with the summary.

#### B. Algorithm:

Algorithm 1 K-mean Algorithm:

- X: a set of N data vectors Data set  
 CI: initialized k cluster centroids Number of clusters,  
 C: the cluster centroids of k-clustering random initial centroids
- 1:  $P = p(i) \text{ — } i = 1, \dots, N$
  - 2: is the cluster label of X
  - 3:  $KMEANS(X, CI) \rightarrow (C, P)$
  - 4: DO

- 5:  $C_{previous} \leftarrow C_i$ ;
- 6: FOR all  $i \in [1, N]$  DO
- 7: Generate new optimal partitions
- 8:  $p(i) \leftarrow \arg \min d(x_i, c_j); 1 \leq j \leq k$
- 9: FOR all  $j \in [1, k]$  DO
- 10: Generate optimal centroids
- 11:  $c_j \leftarrow \text{Average}, e \text{ of } x_i$
- 12: whose  $p(i) = j$ ;
- 13: While  $C = \text{Previous}$

IV. RESULT AND DISCUSSION:

Our framework performs better than IC, for a large portion of the metric. This is because of the way that for building the NLP based situated summaries; we depend on the sentences scored which is imperative in the ranking stage. Subsequently clearly the coming about summaries is shorter than the Input sentences, and no additional data is included. Taking after tables and graph foresee the execution based last result of our framework. Table 1 shows anticipated execution of proposed framework in view of F measure which consider both precision and recall value esteem. Furthermore, general result demonstrates that proposed framework is superior to IC based sentence.

Table 1:-F Measures value of systems according to different no of documents.

No of Doc	IC based F measure	NLP based F measure
5	0.87	0.91
10	0.86	0.92
20	0.75	0.86
40	0.79	0.80

F measure is really test the exactness, f measure consider both precision and recall estimation of framework for figuring f score. F score is harmonic mean of precision and recall. F score esteem is between 0 for most exceedingly terrible and 1 for best score.

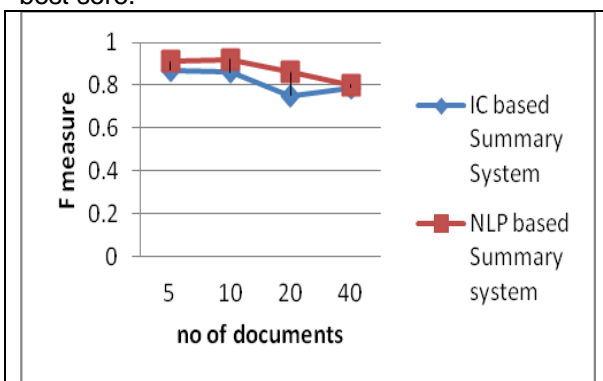


Fig 2.shows F-measure base comparison between Existing and proposed System

Table 2:-Comparison of preprocessing effect on output performance base on time

No of Doc	Time require before preprocessing(ms)	Time require after preprocessing(ms)
5	72	46
10	106	66
20	245	145
40	420	320

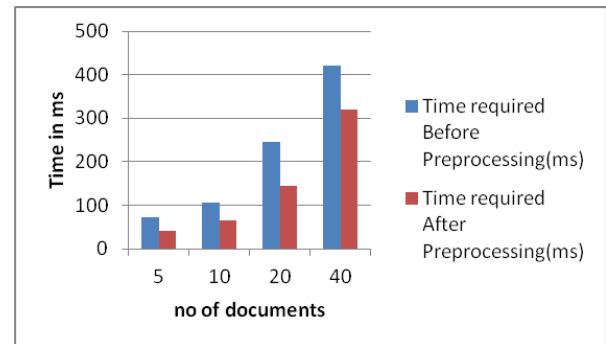


Fig.3. No of documents vs. time in sec

Preprocessing is vital venture in generic summary generation. Preprocessing take a shot at information does some critical procedure like stemming and un required words evacuating methodology. With doing this procedure table1 and diagram shows results. Without preprocessing framework Oblige additional time and gives low precision.

Table 3. Comparison of Preprocessing Effect on Output Performance Base on Time

No of Doc	Accuracy before preprocessing(%)	Accuracy after preprocessing(%)
5	75	85
10	68	86
20	61	90
40	60	88

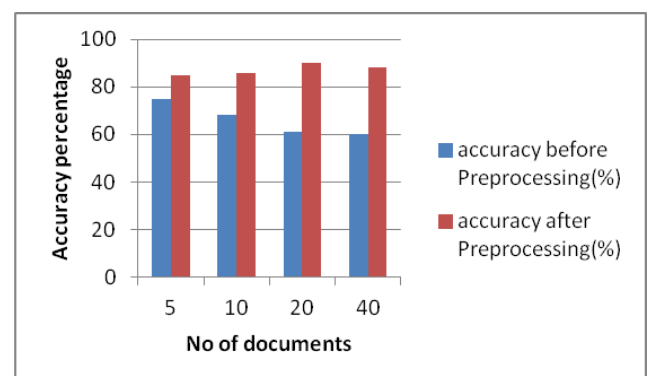


Fig.4: No of documents vs. Accuracy percentage

At last, as per the results got when surveying the user will be fulfill, summary created by our framework is give preferable result over past framework. Final generic summary shows general data of J and K flood in short. Were query based summary will be returns imperative information identified with necessity.

#### V. CONCLUSION:

In the proposed system, an experimental study on a few approaches that use the ontology is given to fathom diverse multi-document summarization issues in disaster management area i.e. J and K surge 2014. For generic summarization, distinctive vector space models are utilized to speak to sentences in the report gathering, and the practicality of distinctive blends of the VSMs is investigated. At that point clustering algorithm is connected on all sentences used to cluster the sentence set and the vital sentences are extricated. For query focused summarization, dove into the impact of query development in summarization assignments. The last summary was therefore created by decreasing data excess and ranking sentences utilizing NLP based sentence extraction strategy which is extremely viable.

#### ACKNOWLEDGMENT:

We are thankful to the teachers for their valuable guidance. We are thankful to the authorities of Savitribai Phule University of Pune and reviewer for their valuable suggestions. We also thank the college authorities for providing the required infrastructure and support. Finally, we would like to extend a heartfelt gratitude to friends and family members.

#### REFERENCES:

[1] Lie li and tao li, "An Empirical study of ontology based Multi-Document Summarization in Disaster Management," IEEE Transaction, vol 44, NO.2. FEBRUARY 2014

[2] V. Nastase, "Topic driven multi-document summarization with encyclopedicKnowledge and spreading activation," in Proc. EMNLP, 2008, pp. 763-772.

[3] C. Lee, Z. Jian, and L. Huang, "A fuzzy ontology and its application to news summarization," IEEE Trans. Syst., Man, Cybern., B Cybern., vol. 35, no. 5, pp. 859-880, Oct. 2005.

[4] H. Saggion, K. Bontcheva, and H. Cunningham, "Robust generic and query-based summarization," in Proc. ECAL, 2003, pp. 235-238.

[5] F. Wei, W.Li, Q. Lu, and Y. He, "Query-sensitive mutual reinforce ent chain and its application in query-oriented multi-document summarization," in Proc. SIGIR, 2008, pp. 283-290.

[6] X. Wan, J. Yang, and J. Xiao, "Manifold ranking based topic focused multi-document summarization," in Proc. IJCAI, 2007, pp. 2903-2908.

[7] J. Tang, L. Yao, and D. Chen, "Multi-topic based query oriented summarization," in Proc. SDM, 2009.

[8] A. Haghighi and L. Vanderwende, "Exploring content models for multi document summarization," in Proc. HLT-NAACL, 2009, pp. 362-370.

[9] E. Klien, M. Lutz, and W. Kuhn, "Ontology based discovery of geographicInformation services an application in disaster management, Compute" Environ. Urban Syst., vol. 30, no. 1, pp. 102-123, 2006.

[10] H. Hsu, C. Tsai, M. Chiang, and C. Yang, "Topic generation for web document summarization," in Proc. IEEE SMC, 2008, pp. 3702-3707.

[11] X. Yong dong, W. Xiao long, L. Tao, and X. Zhi-ming, "Multi-document summarization based on rhetorical structure: Sentence extraction and evaluation," in Proc. IEEE SMC, 2008, pp. 3034-3039.

[12] D. Wang, T. Li, S. Zhu, and C. Ding, "Multi-document summarization via sentence level semantic analysis and symmetric matrix factorization," in Proc. SIGIR, 2008, pp. 307-314.

[13] G. Erkan and D. Radev, "Lex page rank: Prestige in multi-document text summarization," in Proc. EMNLP, vol. 4, 2004, pp. 365-371.

[14] X.Wan and J. Yang, "Multi document summarization using cluster based link analysis," in Proc. SIGIR, 2008, pp. 299-306.

[15] D. Radev, H. Jing, M. Sty, and D. Tam, "Centroid based summarization of multiple documents" Inf. Process. Manage. vol. 40, no. 6, pp.919-938, 2004.

[16] H. Daume and D. Marcu, "Bayesian query focused summarization," in Proc. ACL, vol. 44, no. 1. 2006, p. 305.

[17] D. Jurafsky, J. Martin, A. Kehler, K. Vander Linden, and N. Ward, Speech And Language Processing. Englewood Cliffs, NJ, USA: PrenticeHall, 2000.

[18] S.T.Yuan and J. Sun, "Ontology based structured cosine similarity in speech document summarization," in Proc. WI, 2004, pp. 508-513.

[19] S. Yuan and J. Sun, "Ontology-based structured cosine similarity in document summarization: With applications to mobile audio-based knowledge management," IEEE Trans. Syst., Man, Cybern., B Cybern., vol. 35, no. 5, pp. 1028-1040, Oct. 2005.

[20] P. Tan, M. Steinbach, and V. Kumar, Introduction to Data Mining. Boston, MA, USA: Pearson Addison Wesley, 2006.