

## A STUDY ON INVESTIGATING CONTENT MODELS FOR MULTI-DOCUMENT SUMMARIZATION

1 Joshva Premkumar.P. M.C.A.,M.Phil.,  
2 Boaz Gladson.R. M.C.A.,  
3. Martin Paul Rufus Kumar.K. M.Sc.,B.Ed.,M.Phil.,

1Asst Prof, PG Department of Computer Science,  
Voorhees College,Vellore.

2 Asst Prof, PG Department of Computer Science,  
Voorhees College,Vellore.

3 Asst Prof, Department of Computer Application ( BCA),  
Voorhees College,Vellore.

### Abstract

Robotized text rundown has drawn a ton of interest among the Natural Language Processing and Information Retrieval people group lately. The underlying interest for computerized text outline began during the last part of the 1960s in American examination libraries, where an enormous number of logical papers and books were to be carefully put away and made accessible. Prior to the development of PCs and the rise of the World Wide Web (WWW) as a worldwide computerized library, finding text materials of importance was a demanding assignment. After the approach of WWW the structure and capacity has been modified, where in individuals, academicians, scientists or lay end clients get colossal advantages by perusing the substance on the web. Despite the fact that this has decreased the weight of data assembling, the assignment of obtaining the applicable data in a succinct way is as yet a test.

**Keywords :** clustering, summarization, data mining, documents, language processing, Sentence Rank,Model.

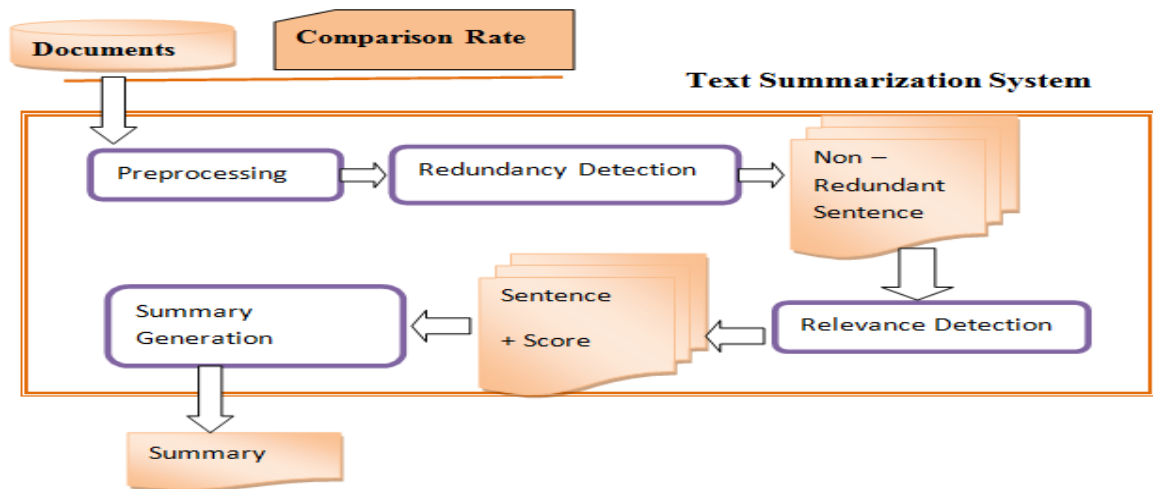
### 1. INTRODUCTION

Text outline is the answer for address this issue. Outline is a method where a PC consequently makes a theoretical or rundown of at least one records. Robotized text outline is the cycle of consequently building synopses for a book contingent upon the client's requirements. A synopsis is an exact portrayal of data relying upon the predefined target pressure proportion. Frameworks summing up single archives are called single report rundown frameworks, while frameworks which play out similar assignment with different related arrangements of records are called multi-record synopsis frameworks.

#### Evolution

Examination on rundown errands took its underlying foundations a very long while back and has kept on being a consistent subject of exploration. Frameworks created in the mid 1950s were named as surface level methodologies misusing topical highlights, for example, Term recurrence, Term Occurrence, result of Term Frequency and Inverse Document Frequency, area based highlights, and presence of foundation terms like title, sign words and expressions Followed by this, substance level methodologies dependent

on syntactic relations, comparability connections, co-event and co-reference were created during the 1960s. Later on, augmentations to the substance level methodologies named as talk based methodologies were created during the 1970s utilizing the expository construction of text and configuration of the record.



## Applications

The application areas for automated text summarization are extensive. Due to the rapid growth of online information, it becomes hard and indeed very difficult to retrieve relevant information in an efficient way. Information is published simultaneously in many media channels in different versions, for instance, a paper newspaper, web newspaper, SMS message, radio newscast, and a spoken newspaper for the visually impaired.

Moreover that is all familiar with summaries such as:

- Headlines (around the world)
- Outlines (notes for students)
- Minutes (of a meeting)
- Previews (of movies)
- Synopses
- Reviews (of a book, CD movie, etc.)
- Digests (TV guide)
- Biography (resumes)
- Abridgements (Shakespeare for children)

- Bulletins (weather forecasts / stock market reports)
- Sound bites (politicians on a current issue)
- Histories (chronologies of salient events).

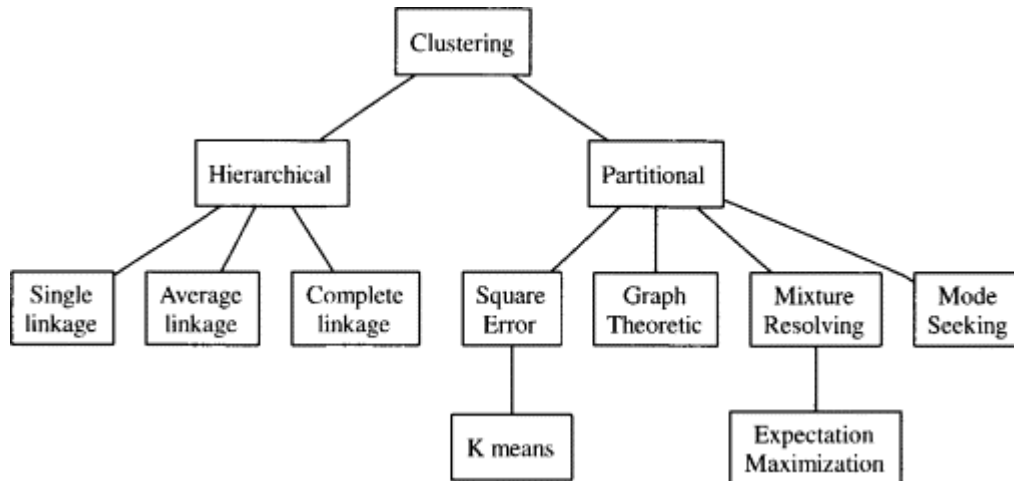
Summarization helps us to save computer resources and bandwidth. For example, if the reader does not understand the source language and if he / she intends to translate such a large document, this translation becomes totally wasted. Instead, translation could be made of only a summary, so that the user can assess if the whole document is worth the effort of translating. Similarly, text-to-speech for the visually impaired can profitably be applied to a summary, before the pronunciation of the whole text is attempted. Summarization is highly beneficial in several information acquisition tasks such as:

- Promoting current awareness
- Saving reading time
- Facilitating effective selection
- Facilitating literature searches
- Improving indexing efficiency
- Aiding the preparation of reviews

As a result, even for a simple topic it becomes difficult to read through all the documents that are related to it. Therefore, the demand to condense documents has increased. The most significant characteristic of multiple document summarization compared with the single one, is that there is a great redundancy of information in a document set. Summarization is an important approach to manage the large amount of text that people must read. Summarization can reduce the amount of text people have to read to let them decide if a document is relevant to their information need. Since the inception of using computers to process a written text, one of the first tasks undertaken was that of summarizing the text by shortening a long document to present the document's content briefly, while preserving the underlying meaning.

#### Existing approach

Summarization of text documents involves investigations in several connected aspects, like methodologies adopted for picking the top ranked sentences, removal of redundancies and evaluation of extracted summaries. This chapter provides a detailed literature survey of similarity measures, evaluation of summaries and summarization methodologies.



Document structure similarity algorithms like the Optimal Tree Edit distance algorithm, Tag similarity, Fourier transforms and Path similarity algorithms. These approximation algorithms include the simple weighted tag similarity algorithm, the Fourier transform and a new application of the Shingle technique for structural similarity of documents. Structural similarity as applied by the Shingle technique measures the similarity among the structures of documents. This technique is found to reduce the set of words or tokens in a document into a list of hashes that can be directly compared with the document using set difference, union and intersection to determine the similarity. Investigations show that path similarity algorithms combined with the Shingle technique are found to be most effective as compared to the basic Fourier Transform and Tree Edit Distance algorithms.

Classes of similarity measures to evaluate the similarity among the publications provided by the ACM digital libraries. These classes are categorized as text-based and citation-based similarity measures. The evaluation of these metrics was carried out in terms of accuracy, separability and independence. These measures evaluate the similarity in the title, abstract, index terms and body of the document. The quality of publication similarity measures used for locating related or similar publications were also discussed. Using the Vector Space Model and the TF/IDF weighting scheme, the similarity between two publications is measured by using cosine, Jaccard, dice or other document metrics.

Proposed approach

The proposed two enhancements to the already existing graph based approaches. These enhancements are applicable to all the existing methods. This approach describes the two enhancements of the discounting technique and the incorporation of the position weight in the next two subsections.

### Discounting Technique

Input: Symmetric adjacency Matrix – A;

Compression ratio r;

Method Chosen method;

Output : Sorted list of sentences s\_list<>;

```
begin
s_list    empty;
n' = n * r ;
do while n' > 0
    begin
call chosen_method( );
method
SWmax = 0.0 ;
for i = 1 to n do
if SW(i) > SWmax
{
SWmax = SWi ; nn = i ; }
s_list ← s_list + nn ;
for i = 1 to n do
{ a[i,nn] = 0;
a[nn,i] = 0;
}
n' = n - 1;
end;
sort s_list <>
end;
```

### SentenceRank (Continuous)

This is a modification of Method **Discounted LexRank (Continuous)** with the incorporation of the position weight. This is rather fortuitous and the actual performance comparison has to be based on the average values obtained from a collection of the document set

**Input:** Symmetric adjacency Matrix - A;

compression ratio r;

d : damping factor

n : number of sentences

**Output :** Sorted list of sentences s\_list<>;

Array : SR(thres) [n] , SR(Cont)[n]

$\xi$  : error tolerance value;

begin

s\_list  $\leftarrow$  empty;

n' = n \* r ;

w[i] = 1/n;

call discounting( );

do

{

for i 1 to n do

for j 1 to n do

a[i][j] = idf-modified-cosine(S[i],S[j]);

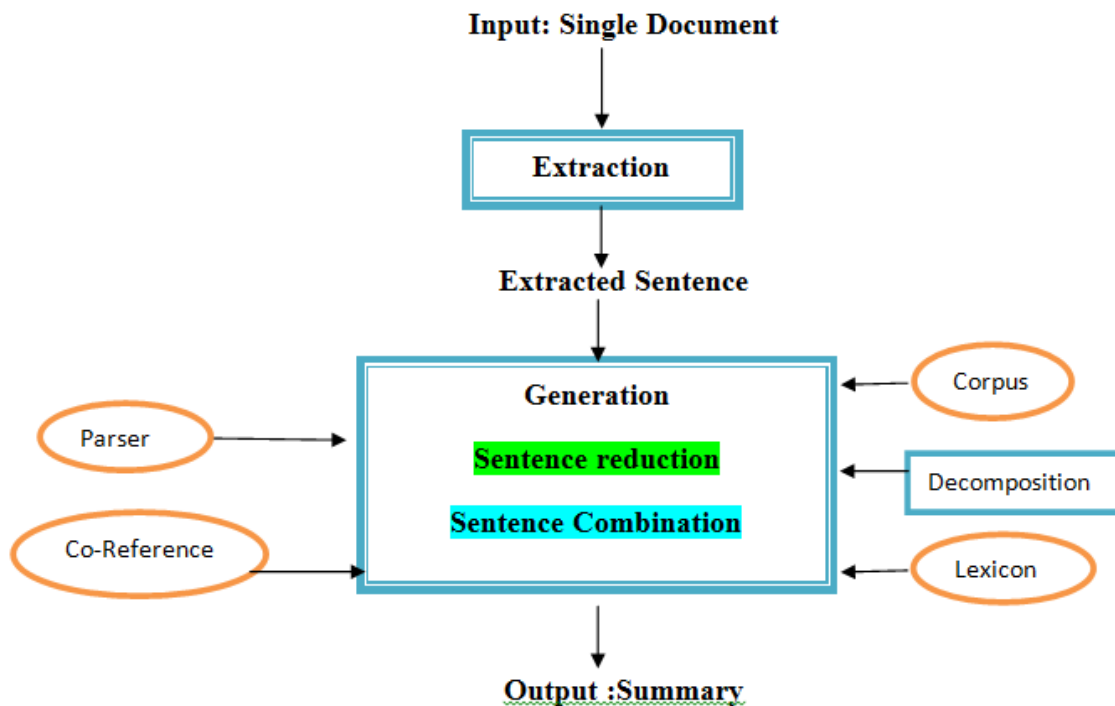
if a[i][j] > t then

a[i][j] = 1; Degree[i] ++;

else

a[i][j] = 0;

end

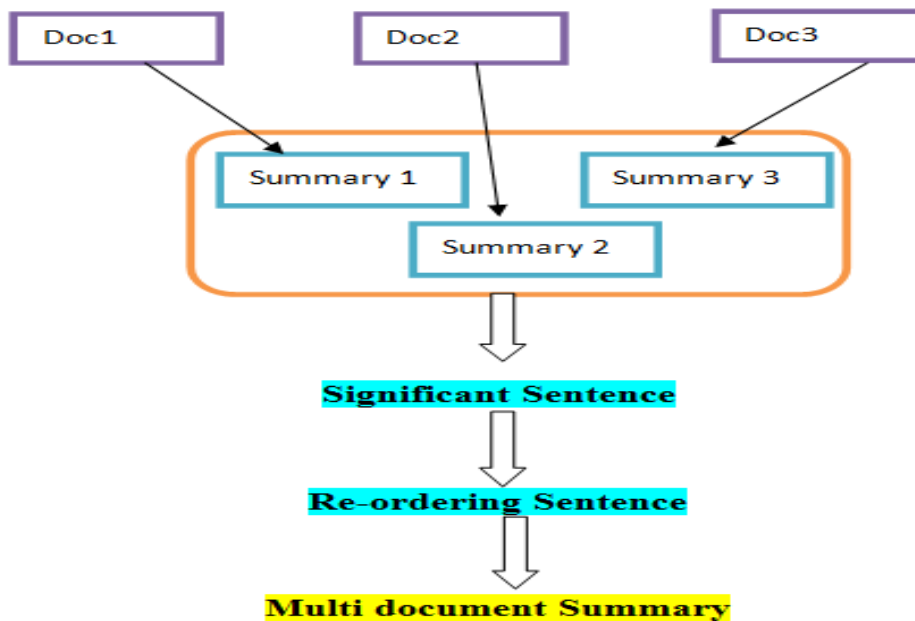


This part researched top to bottom, two classes of graphical techniques for text outline. The top of the line compares to strategies for the non-PageRank type, while the subsequent gathering depends on the PageRank type calculations. In each class, the limiting strategies proposed in this section are better than the essential techniques and the proposed limiting in addition to situate weight approach tolls the best. Every one of the twelve strategies are promising in that they yield better outcomes as looked at than irregular choice, in view of the ordinary exactness metric just as by the proposed measurements Effectiveness1 (E1) and Effectiveness2 (E2). It is brought out from the examinations introduced, that dependent on the normal execution of over a 30-report set, techniques Sentence Rank (Threshold) and Sentence Rank (Continuous) – the proposed Sentence Rank (Threshold) and Sentence Rank (Continuous), yields the best consequences of the multitude of 12 strategies considered. The following section presents the examinations did for multi-record processed.

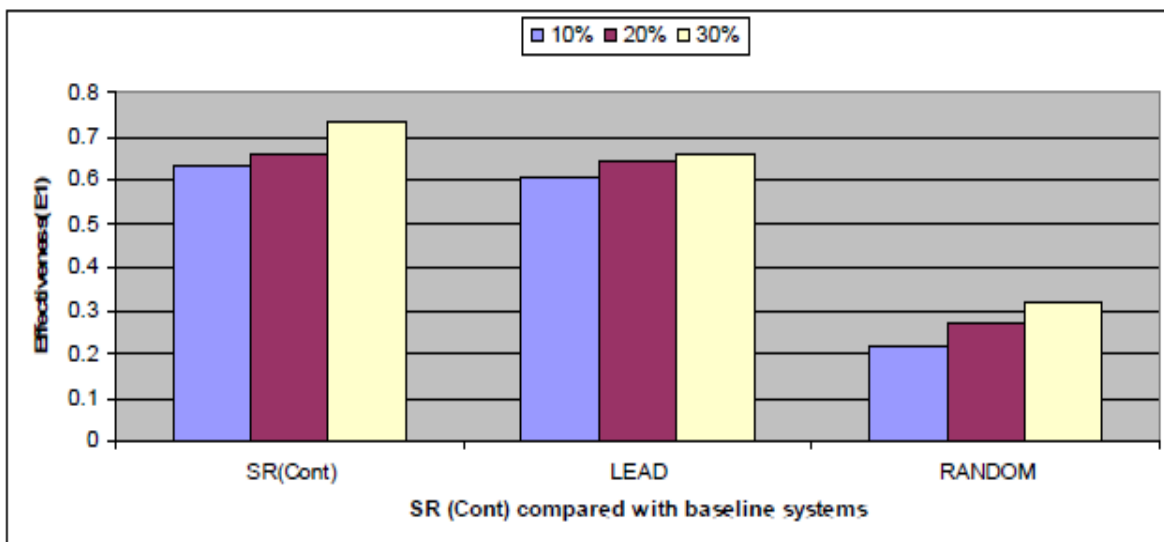
#### Graph Approach

Diagram based methodologies for outlines are very famous. These strategies are displayed under two sorts of informal organizations. Allow us to think about a true circumstance to characterize these two kinds to understand their significance. An individual having broad contacts with individuals in an association is viewed as more significant than an individual with less contacts. Subsequently, the

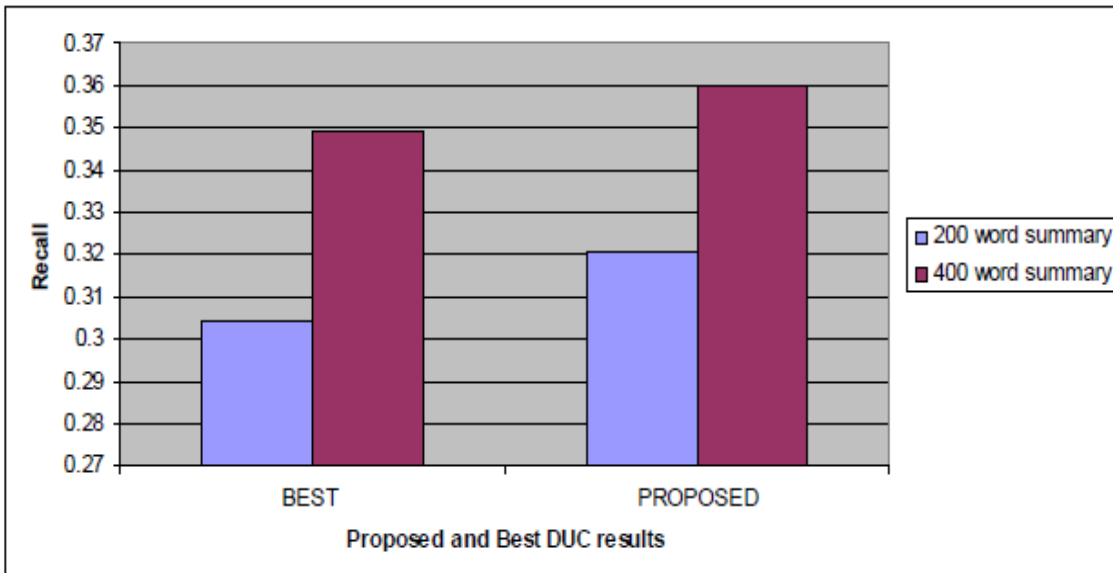
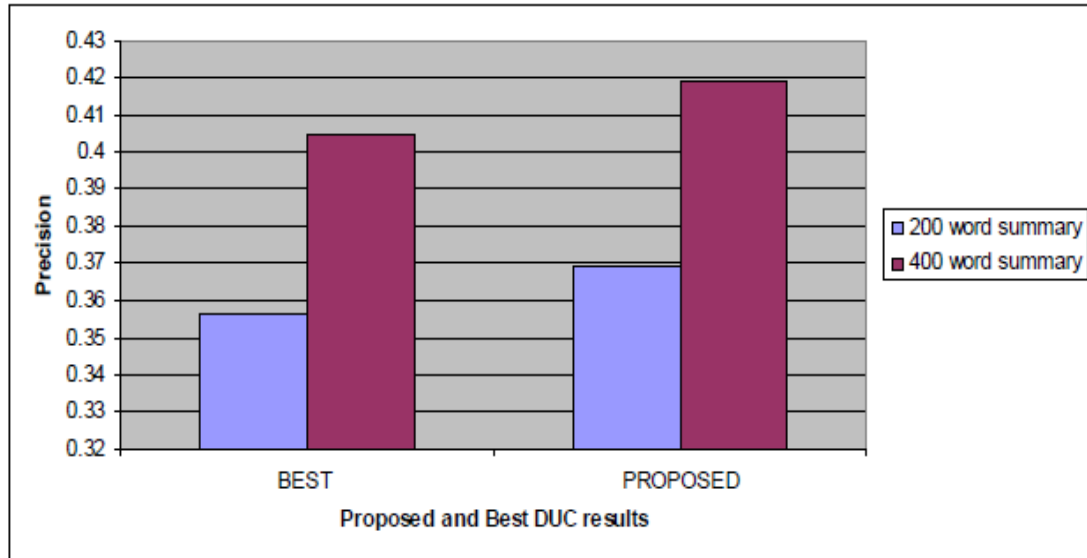
individual's conspicuousness can be essentially decided in a vote based way, by the quantity of contacts he has. Then again, let us consider the instance of a second individual who has less contacts, however the entirety of his contacts are exceptionally positioned and persuasive people. Obviously, in the present circumstance, the subsequent individual may have significant impact and eminence contrasted with the previous.



### Analysis







## CONCLUSION

The rundown of text reports has been an intensely explored zone. This proposal has researched two classes of graphical strategies for text outline. The five star relates to essential strategies for non PageRank type, while second gathering depends on PageRank type calculations. It is shown that in each class limiting strategies proposed in this theory is better than fundamental techniques and the proposed limiting procedure in addition to situate weight strategy admissions the best. The Sentence Rank (Continuous) technique is found to yield better outcomes as thought about than the best distributed outcomes no of informational index. The postulation has investigated elective techniques for the natural assessment of synopses and has proposed another measurement called 'Adequacy'. Further advances have been formalized for the arrangement of the 'best quality level' reference synopsis. Studies done utilizing the corpus of archives gathered from business and reserch locales and DUC 2002 informational collection set up the prevalence of the techniques proposed. The Sentence Rank (Threshold) and Sentence Rank (Continuous) approaches proposed, yield better outcomes for both the informational indexes, independent of the assessment measures.

## REFERENCE

- [1].C. Wang and D. M. Blei, "Collaborative topic modeling for recommending scientific articles," in Proc. 17th ACM SIGKDD Int. Conf. Knowl. Discov. Data Min., 2011, pp. 448–456.
- [2].Y. Gao, Y. Xu, Y. Li, and B. Liu, "A two-stage approach for generating topic models," in Advances in Knowledge Discovery and Data Mining, PADKDD'13. New York, NY, USA: Springer, 2013, pp. 221–232.
- [3].N. Zhong, Y. Li, and S.-T. Wu, "Effective pattern discovery for text mining," IEEE Trans. Knowl. Data Eng., vol. 24, no. 1, pp. 30–44, Jan. 2012.
- [4].A. Tagarelli and G. Karypis, "A segment-based approach to clustering multi-topic documents," Knowl. Inform. Syst., vol. 34, no. 3, pp. 563–595, 2013.
- [5].Hongyan Liu, Ping'an Liu, Wei Heng and Lei Li, 2011.*The CIST Summarization System at TAC2011. TAC 2011 Proceedings.*
- [6].Alguliev R.M. and Aliguliyev R.M. (2005), 'Effective Summarization Method of Text Documents', Proceedings of the 2005 IEEE/WIC/ACM International Conference on Web Intelligence, Compiegne University of Technology France, pp. 264-271.

- [7].Ali M., Ghosh M.K. and Abdullah Al Mamun (2009), 'Multidocument Text Summarization: SimWithFirst Based Features and Sentence Co-selection Based Evaluation', Proceedings of the International Conference on Future Computer and Communication, Kuala Lumpur, Malaysia, pp. 93-96.
- [8].Antiqueira L., Oliveira O.N., Costa L.D.F. and Nunes M.D.G.V. (2009), 'A Complex Network Approach to Text Summarization', Information Sciences, Vol. 179, No. 5, pp. 584-599.
- [9].Banu M., Karthika C., Sudarmani P. and Geetha T.V. (2007), 'Tamil Document Summarization Using Semantic Graph Method', Proceedings of International Conference on Computational Intelligence and Multimedia Applications, Sivakasi, Tamilnadu, India, Vol. 2, pp. 128-134.