



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirce.com

Vol. 5, Issue 3, March 2017

Text Summarization Formed by Ranking Sentence from Key Phrase Extraction

Vidya Panchal, Asra Anjum

M.Tech Student, Dept. of CSE, Maharashtra Institute of Technology, , BAMU, Maharashtra, India

Professor, Dept. of CSE, Maharashtra Institute of Technology, BAMU, Maharashtra, India

ABSTRACT: Text Summarization involves reducing large amount of text into smaller in size. So that we can get information about what kind of the data are present in that document. We can get more correct amount of information of the large text or document. Now days, various techniques are available to obtain the summary from the document or large text. In this, we represent the technique that help you to get summary of the document. This technique based on the ranking the sentence from key phrase extraction algorithm.

KEYWORDS: Ranking, Key phrase, text summary, Data mining algorithm.

I. INTRODUCTION

Now a day's tremendous amount of the data available everywhere in the form of wiki website, large text, documents and internet. So that to read that much amount of the data is the time consuming. Text summarization helps us to summarize that data into in the form of abstraction so that we can get exact and generic and more specific amount of information of that document or internet website or large text. According to MichelangeloCecia ,CorradoLogliscia , LucreziaMacchia[1] , they stated thatsummary of the document image is obtain using keyphrase extraction algorithm. According to Jackson, P., they stated that summary can be obtained in the form of extracts, they are formed by extracting the key phrases of the original text or document and abstracts of the document. When they are formed by concluding the meaning of the original document or by re-generating the content of original document.[2] According to Sparck Jones K., They stated that summary can be obtained by more extractive way like selecting more salient sentences through Machine Learning or Data Mining algorithms. So that we can obtained the summary of the original document by more effective way. According to PaiceC.D[3]. , they stated that sentences are mentioned in the form of lexical and structural features (e.g., keywords frequency, title keywords, sentence location, indicator phrases, etc.) and described or represented by as vectors of quantitative and categorical measures of those features (attribute-value representation) in the original document. And Zhuli, X., Li X., Barbara, D.E., Peter, N., Weimin, X., Thomas T[4] , they mentioned that the summary of the document is also obtained by Several supervised learning techniques which are the part of the Data Mining algorithms. To rank new sentences of the original document final ranking sentence function is used. Svore et al.[9] stated that resort to a neural network pair to generate a ranking composed of three sentences is also used to developed summary of the original document.

II. RELATED WORK

In based paper they use WISDOM++ tool for the transform the document image into XML format of that image. This Process complete into 6 steps. The document image is then divide into several part which shows coherent components or imp part of the document.

Ranking the sentence is based on the extracted based technique which leads to supervised data mining algorithm in data mining. To learn the ranking the model data mining algorithms is used. Several supervised learning techniques is found in the data mining algorithm and these are used to generate the ranking functions. We can generate the ranking functions using neural network pair [10].

III. PROPOSED ALGORITHM

1. Sentence Ranking for Summary:

Learning a preference model that leads to identify preference relations



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirccce.com

Vol. 5, Issue 3, March 2017

to be applied to new documents is used to solved the problem of calculating summary of the Document Text or URLs.

1.1 Mining Preference relations:

The problem of mining preference relations between sentences can be obtained as follows.

Given:

- A database schema S with h relational tables $S = \{T_1, T_2, \dots, T_h\}$
- Two sets PK and FK of primary and foreign key constrains on tables in S
- A target relation $T \in S$ representing sentences that play the role of reference objects
- A precedence relation $PT \in S$ with two attributes. Each tuple in this table represents an ordered pair of reference objects where the first reference object precedes the second one.

Find: A probability estimation $P(a < b | a, b)$ for any couple of sentences a and b which are present to a new document represented according to the schema $S - PT$. Objects in $S - \{T, PT\}$ play the role of task relevant objects, while the

Precedence relation implicitly defines a partial ordering between two sentences.

In our approach, it is also possible to avoid to consider some sentences belonging to parts of the document that are not considered relevant for the task at hand (e.g. sentences in tables or sentences in references of a original paper). This means that the preference relation in training data does not necessarily express total ordering of sentences in training documents.

By applying the Bayes theorem, $P(a < b | a, b)$ can be computed as:

$$P(a < b | a, b) = P(a < b)P(a, b | a < b) / P(a, b) \dots\dots\dots(1)$$

where:

- $P(a < b)$ in (1) denotes the prior probability that a sentence precedes another. This probability might be different from 0.5.
- $P(a, b) = P(a < b)P(a, b | a < b) + P(b < a)P(a, b | b < a)$.

1.2. Patterns construction:

The relational pattern discovery is obtained by exploring level-by-level the lattice of relational patterns ordered according to a generality relation (\geq) between patterns. Given two patterns P_1 and P_2 , $P_1(\geq)P_2$ shows that P_1 (P_2) is more general (specific) than P_2 (P_1). Hence, the search gives us the most general pattern for the summary and iteratively alternates the candidate generation and candidate evaluation phases to obtain the summary. In [14], the authors mention an enhanced version of the level-wise method [15] to discover patterns from data in multiple tables of a relational database for the calculating the summary of the document. The space of linked relational patterns are searched by candidate patterns, which is constructed according to the θ -subsumption generality order [16].

This shows possible to obtain a levelwise exploration of the lattice of relational patterns ordered by θ -subsumption. Specifically, patterns are obtained by forming the pattern space one level at a time starting from the most specific Pattern (the pattern have only the *preference/2* predicate) and then by applying a breadth-first evaluation in the relational patterns ordered according to $\geq \theta$.

1.3 Ranking Reconstruction

The main goal of this step is to produce ranking the sentences the sentences for the construction of the summary of the original document.

For the ranking the sentence for construction of the summary following algorithm are used:

Algorithm 1 ranking identification algorithm:

- 1: find ranking ($G = [V, E]$): Ranking L
- 2: $L \leftarrow \Theta$;
- 3: while ($\#L < \#V$) do
- 4: $L.add(\arg \max_{(bi \in V/L)} SUMPREFG(bi))$;
- 5: end while



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirce.com

Vol. 5, Issue 3, March 2017

Where,

- $G = \{V, E\}$ be a labeled directed graph where $V = \{b \in T\}$ and $E = \{(a, b, w(a,b) \in V^*V \times [0, 1]) | w(a,b) = P(a < b | a, b)\}$ is the set of weighted edges of the graph where weights are the probabilities $P(a < b | a, b)$ computed according to equation (1),
- $SUMPREFG: V \rightarrow [0, \#V]$.

This Algorithm works on the method for the ranking identification. If the $SUMPREFG()$ of the sentence is high then that sentence is added to the summary of the document text. Once this implementation of the algorithm is done then that 'n' sentences is obtained from the algorithm is added to the Summary of the original document.

1.4. Data extraction and representation

Extraction of the key phrases from the document involves tokenization, sentence splitting, part-of-speech (POS) tagging, stop-word removing and stemming.

The representation of the sentences form by the ranking the key phrases extracted by extraction algorithm is done by means of a phase of natural language processing.

IV. SIMULATION RESULTS

Comparison based Confusion Matrix

	System Yes	System No
Manual yes		
Manual No		

Accuracy: $(TP + TN)/Total$

Misclassification Rate: $(FP + FN)/ Total$

True Positive Rate: $TP/ actual\ yes$

False positive Rate: $FP/actual\ no$

Specificity : $TN/actual\ no$

Precision : $TP/ predicated\ yes$

Prevalence : $actual\ yes/ total$

Where:

TP = True Positive

TN = True Negatives

FP = False Positives

FN = False Negatives

For Link1

	System Yes	System No
Manual Yes	2	3
Manual No	3	2

Table 4.1 Link 1



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijircce.com

Vol. 5, Issue 3, March 2017

For Link2

	System Yes	System No
Manual Yes	4	0
Manual No	4	2

Table 4.2 Link 2

For Link3

	System Yes	System No
Manual Yes	4	1
Manual No	1	4

Table 4.3 Link 3

For Link4

	System Yes	System No
Manual Yes	3	1
Manual No	2	4

Table 4.4 Link 4

For link6

	System Yes	System No
Manual Yes	3	0
Manual No	2	3

Table 4.5 Link 6

For Link7

	System Yes	System No
Manual Yes	0	5
Manual No	5	0

Table 4.6 Link 7

For link8

	System Yes	System No
Manual Yes	1	4
Manual No	3	2

Table 4.7 Link 8



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijircce.com

Vol. 5, Issue 3, March 2017

For link 9

	System Yes	System No
Manual Yes	3	2
Manual No	3	2

Table 4.8 Link 9

For link10

	System Yes	System No
Manual Yes	2	2
Manual No	4	2

Table 4.9 Link 10

For Link11

	System Yes	System No
Manual Yes	2	3
Manual No	3	2

Table 4.10 Link 11

For WordFile 1

	System Yes	System No
Manual Yes	3	1
Manual No	2	4

Table 4.11 WordFile1

For WordFile 2

	System Yes	System No
Manual Yes	3	2
Manual No	2	3

Table 4.12 WordFile2

For WordFile 3

	System Yes	System No
Manual Yes	4	0
Manual No	4	2

Table 4.13 WordFile3



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirce.com

Vol. 5, Issue 3, March 2017

For WordFile 4

	System Yes	System No
Manual Yes	4	1
Manual No	1	4

Table 4.14 WordFile4

For WordFile 5

	System Yes	System No
Manual Yes	3	1
Manual No	2	4

Table 4.15 WordFile5

4.3 Performance Evaluation Table:

	Accuracy	M.Rate	True positive rate	False positive rate	Specificity	Precision	Prevalence
Link1	0.4	0.6	0.4	0.6	0.4	0.4	0.5
Link2	0.6	0.4	1	0	0.33	1	0.4
Link3	0.8	0.2	0.8	0.2	0.8	0.8	0.5
Link4	0.7	0.3	0.75	0.66	0.33	0.6	0.4
Link6	0.8	0.2	1	0.4	0.6	0.6	0.5
Link7	0.5	0.5	0	1	0	0	0
Link8	0.3	0.7	0.25	0.6	0.4	0.25	0.5
Link9	0.5	0.5	0.6	0.6	0.4	0.5	0.5
Link10	0.4	0.6	1	0.6	0.3	0.5	0.4
Link11	0.4	0.6	0.6	0.6	0.4	0.6	0.5
WordFile1	0.7	0.3	0.75	0.3	0.7	0.6	0.4
WordFile2	0.6	0.4	0.75	0.3	0.7	0.6	0.5
WordFile3	0.6	0.4	1	0.3	0.7	0.3	0.4
WordFile4	0.8	0.2	0.8	0.2	0.8	0.8	0.5
WordFile5	0.7	0.3	0.75	0.3	0.7	0.6	0.5

Accuracy of all 15 Documents:

$((\text{Sum of Accuracy of Docs})/15)*100$

$= ((8.8)/15)*100$



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirccce.com

Vol. 5, Issue 3, March 2017

= 59%

Misclassification Rate of all Documents:

$((\text{Sum of all M. Rate of Docs})/15)*100$

$= ((6.2)/15)*100$

= 41%

V. CONCLUSION

In this paper, we purpose the how to obtained the summary from large amount data or large document or URLs. For this we use sentence ranking algorithm to form summary of the original document using the key phrases. Key Phrases for ranking sentences is obtained by key phrases extraction algorithm using relational data mining algorithm or tools.

REFERENCES

- 1.MichelangeloCecia,*, CorradoLogliscia, LucreziaMacchia.” Ranking Sentences for Keyphrase Extraction: A Relational Data Mining Approach”
Procedia Computer Science 38 (2014) 52 – 59
- 2.Jackson, P.. The oxford handbook of computational linguistics edited by ruslanmitkov. *Computational Linguistics* 2004;30(1):103–106.
- 3.Sparck Jones, K.. Automatic summarising: The state of the art. *Inf Process Manage* 2007;43(6):1449–1481.
doi:<http://dx.doi.org/10.1016/j.ipm.2007.03.009>
- 4.Paice, C.D.. Constructing literature abstracts by computer: Techniques and prospects. *Inf Process Manage* 1990;26(1):171–186.
- 5.D`zeroski, S., Lavra`c, N.. *Relational Data Mining*. Springer-Verlag; 2001
- 6.Ceci, M., Loglisci, C., Ferilli, S., Malerba, D.. Project d.a.m.a.: Document acquisition, management and archiving. In: Agosti, M., Esposito, F., Meghini, C., Orio, N., editors. *IRCDL*; vol. 249 of *Communications in Computer and Information Science*. Springer. ISBN 978-3-642-27301-8; 2011, p. 115–118.
- 7.Ceci, M., Berardi, M., Malerba, D.. Relational data mining and ILP for document image understanding. *Applied Artificial Intelligence* 2007;21(4&5):317–342.
- 8.Ceci, M., Berardi, M., Porcelli, G., Malerba, D.. A data mining approach to reading order detection. In: *ICDAR*. IEEE Computer Society.ISBN 978-0-7695-2822-9; 2007, p. 924–928.
- 9.Svore, K., Vanderwende, L., Burges, C.. Enhancing single-document summarization by combining ranknet and third-party sources. In: *Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*. Association for Computational Linguistics; 2007, .
- 10.Helft, N.. *Progress in Machine Learning*; chap. Inductive generalization: a logical framework. Sigma Press; 1987, p. 149–157.
- 11.Ceci, M., Appice, A., Malerba, D.. Emerging pattern based classification in relational data mining. In: Bhowmick, S.S., K`ung, J.,Wagner,R., editors. *DEXA*; vol. 5181 of *Lecture Notes in Computer Science*. Springer. ISBN 978-3-540-85653-5; 2008, p. 283–296.
12. Robinson, J.A.. A machine oriented logic based on the resolution principle. *Journal of the ACM* 1965;12:23–41
- 13.Ceci, M., Appice, A., Malerba, D.. Discovering emerging patterns in spatial databases: A multi-relational approach. In: Kok, J.N., Koronacki, J., de M`antaras, R.L., Matwin, S., Mladenic, D., Skowron, A., editors. *PKDD*; vol. 4702 of *Lecture Notes in Computer Science*. Springer. ISBN 978-3-540-74975-2; 2007, p. 390–397.
- 14.Mannila, H., Toivonen, H.. Levelwise search and borders of theories in knowledge discovery. *Data Mining and Knowledge Discovery* 1997;1(3):241–258.
- 15.Kamishima, T., Akaho, S.. Learning from order examples. In: *ICDM*. IEEE Computer Society. ISBN 0-7695-1754-4; 2002, p. 645–648.
- 16.Liu, K., Terzi, E., Grandison, T.. Manyaspects: a system for highlighting diverse concepts in documents. *PVLDB* 2008;1(2):1444–1447.