# Generic Multi-Document Summarization Using Topic-Oriented Information

Yulong Pei, Wenpeng Yin, and Lian'en Huang

Shenzhen Key Lab for Cloud Computing Technology and Applications
Peking University Shenzhen Graduate School
Shenzhen, Guangdong 518055, P.R. China
`peiyulong@sz.pku.edu.cn, mr.yinwenpeng@gmail.com, hle@net.pku.edu.cn`

**Abstract.** The graph-based ranking models have been widely used for multi-document summarization recently. By utilizing the correlations between sentences, the salient sentences can be extracted according to the ranking scores. However, sentences are treated in a uniform way without considering the topic-level information in traditional methods. This paper proposes the topic-oriented PageRank (ToPageRank) model, in which topic information is fully incorporated, and the topic-oriented HITS (ToHITS) model is designed to compare the influence of different graph-based algorithms. We choose the DUC2004 data set to examine the models. Experimental results demonstrate the effectiveness of ToPageRank. And the results also show that ToPageRank is more effective and robust than other models including ToHIST under different evaluation metrics.

**Keywords:** Multi-Document Summarization, PageRank, HITS, LDA.

## 1 Introduction

As a fundamental tool for understanding document data, text summarization has attracted considerable attention in recent years. Generally speaking, text summarization can be defined as the process of automatically creating a compressed version of a given text set that provides useful information for the user [3]. In the past years, two types of summarization have been explored: extractive summarization and abstractive summarization. Extractive summarization aims to generate summary directly by choosing sentences from original documents while abstractive summarization requires formulating new sentences according to the documents. Although abstractive summarization could be more concise and understandable, it usually involves heavy machinery from *natural language processing*. In this paper, we focus on extractive multi-document summarization.

In order to generate the highly comprehensive summary of a given document set, multi-document summarization can be divided into 3 steps: 1) computing the scores of sentences in the document set; 2) ranking the sentences based on the scores (or combined with some other rules); 3) choosing the proper sentences as the summary. Among the three steps, computing the scores plays the most

important role. Most recently, some graph-based ranking algorithms have been successfully applied for computing the sentence scores by deciding on the importance of a vertex in the graph according to the global information of the document set. Normally, a set of documents are represented as a directed or undirected graph [3] based on the relationship between sentences and then the graph-based ranking algorithm such as PageRank [12] or HITS [4] is used. According to the previous study, a concise and effective summarization should meet three requirements: diversity, coverage and balance [5]. However, these methods often employed the sentences or terms in a uniform usage [16] without considering the topic-level information, which could lead to less satisfaction of these requirements because of ignoring the information of hidden diverse topics.

To deal with the problem in previous graph-based ranking models, we propose the topic-oriented PageRank (ToPageRank) model, which is to divide traditional PageRank into multiple PageRanks with regard to various topics and then extract summaries specific to different topics. Afterwards, based on the topic distribution of all the documents the entire summaries will be obtained. Aim to explore the influence of different graph-based ranking algorithms, we compare the ToPageRank model with a modified HITS model, named topic-orientd HITS (ToHITS) accordingly. In ToHITS, topics and sentences are regarded as hubs and authorities respectively and the hub scores and authority scores are computed iteratively in a reinforcement way. Experiments on DUC2004[1] dataset have been performed and the results demonstrate the good effectiveness of the proposed ToPageRank model which outperforms all other models under various evaluation metrics. The parameters in experiments have also been investigated and the results show the robustness of the proposed model.

The rest of the paper is organized as follows. First we introduce the related work in Section 2. The basic models are discussed in Section 3. We discribe the ToHITS and ToPageRank models in Section 4. The experiments and results are represented in Section 5 and finally we conclude this paper in Section 6.

## 2   Related Work

Multi-document summarization aims to generate a summary by reducing documents in size while retaining the main characteristics of the original documents [18]. According to the differences of provided information, multi-document summarization can be classified into generic summarization and query-oriented summarization. In particular, generic summarization would generate the summary only based on the given documents and query-oriented summarization would form the summary on a certain question or topic. Both generic summarization and query-oriented summarization have been explored recently.

Traditional feature-based ranking methods exploited different features of sentences and terms to compute the scores and rank the sentences. One of the

---

[1] `http://www-nlpir.nist.gov/projects/duc/guidelines/2004.html`

most popular feature-based methods is centroid-based method [13]. MEAD as an implementation of cetroid-based summarizer applied a number of predefined features such as TF*IDF, cluster centroid and position to score the sentences. Lin and Hovy [6] used term frequency, sentence position, stigma words and simplified Maximal Marginal Relevance (MMR) to build the NeATS multi-document summarization system. You Ouyang et al. [11] studied the influence of different word positions in summarization.

Cluster information has also been explored in generating the summaries. Wang et al. [18] integrated both document-term and sentence-term matrices in a language model and utilized the mutual influence of document clusters and summaries to cluster and summary the documents simultaneously. Wan and Yang [16] clustered the documents and combined the cluster information with the Condition Markov Random model and HITS model to rank the sentences. Cai et al. [2] studied three different ranking functions and proposed a reinforcement method to integrate sentences ranking and clustering by making use of term rank distributions over clusters.

Graph-based ranking algorithms such as PageRank and HITS nowadays are applied in text processing i.e. TextRank [9] and LexPageRank [3]. Graph-based ranking algorithms take global information into consideration rather than rely only on vertex-specific information, therefore have been proved successful in multi-document summarization. Some works [17] [14] have extended the traditional graph-based models recently. In [17], a two-layer link graph was designed to represent the document set and the ranking algorithm took into account three types of relationship, i.e. relationship between sentences, relationship between words and relationship between sentences and words.

However, in the previous methods the topic-level information has seldom been studied while in our work we incorporate the topic distribution in graph-based ranking algorithms so the topic-level information can be well utilized.

## 3   Basic Models

### 3.1   Overview

The graph-based ranking algorithms stem from web link analysis. Hyperlink-Induced Topic Search (HITS) [4] algorithm and PageRank [12] algorithm are the most widely used algorithms in recent years and a number of improvements of both algorithms have been exploited recently as well. Based on PageRank, some graph-based ranking algorithms, i.e. LexPageRank and TextRank, have been introduced to text summarization. In these methods documents are represented as graph structure, more specifically, the nodes in the graph stand for the sentences and the edges stand for the relationship between a pair of sentences. The graph-based ranking algorithms are used to compute the scores of sentences and the sentences with high scores would be chosen as the summaries.

## 3.2   HITS

In HITS [4], Kleinberg proposed that a node had two properties, named hub and authority. A good hub node is one that points to many good authorities and a good authority node is one that is pointed to by many good hubs. In the basic HITS model, the sentences are considered as both hubs and authorities.

Given a document set $D$, we denote the graph as $G_H=(S, E_{SS})$ to represent the documents. $S = \{s_i|0 \leq i \leq n\}$ is the set of vertices in the graph and stands for the set of sentences, and $E_{SS} = \{e_{ij}|s_i, s_j \in S, i \neq j\}$ corresponds to the relationship between each pair of sentences. Each $e_{ij}$ is associated with a weight $w_{ij}$ which indicates the similarity of the pair of sentences. The weight is computed by using the standard cosine measure between two sentences as follows.

$$w_{ij} = sim_{cosine}(s_i, s_j) = \frac{\vec{s_i} \cdot \vec{s_j}}{|\vec{s_i}| \times |\vec{s_j}|} \quad , \tag{1}$$

where $\vec{s_i}$ and $\vec{s_j}$ are the term vectors corresponding to the sentence $s_i$ and $s_j$, respectively. The graph we propose to build is undirected so we have $w_{ij} = w_{ji}$ here and we follow [15] to define $w_{ii} = 0$ to avoid self transition. TF*ISF (inverse sentence frequency) value of each term is applied to describe the elements in the sentence vector. Then the weight value of sentences in the documents can be denoted as a symmetric matrix $W$.

The authority score $Auth^{(t+1)}(s_i)$ of sentence $s_i$ and the hub score $Hub^{(t+1)}$ $(s_j)$ of sentence $s_j$ in the $(t + 1)^{th}$ iteration are computed based on the corresponding scores in the $t^{th}$ iteration as follows.

$$\begin{aligned} Auth^{(t+1)}(s_i) &= \sum_{s_j \in S, i \neq j} w_{ij} \cdot Hub^{(t)}(s_j) \\ Hub^{(t+1)}(s_j) &= \sum_{s_j \in S, i \neq j} w_{ij} \cdot Auth^{(t)}(s_i) \end{aligned} \quad , \tag{2}$$

where $Hub^{(t+1)}(s_i)$ is equal to $Auth^{(t+1)}(s_i)$ in the $(t+1)^{th}$ iteration because the sentences are considered as both hubs and authorities; besides, the matrix $W$ is a symmetric matrix, therefore $w_{ij} = w_{ji}$ here. To guarantee the convergence of the iterative process, $Auth(\cdot)$ (or $Hub(\cdot)$) would be normalized after each iteration as $Auth^{(t)} = Auth^{(t)}/\|Auth^{(t)}\|$. Both $Auth(s_i)$ and $Hub(s_i)$ can be initialize as $1/N$ and $N$ is the number of sentences.

Usually the iterative process is stopped when the difference between the scores of sentences in two sequential iterations falls below a defined threshold or the number of iteration exceeds a given value. After iteration, the authority scores are used as the scores of sentences.

## 3.3   PageRank

PageRank is one of the most important factors for Google to rank the search results. This algorithm computes the scores of nodes by making use of the voting or recommendations between nodes.

Using the similar denotations in previous section, $G_{PR}=(S, E_{SS})$ is an undirected graph to represent the sentences and sentence correlations in a given document set $D$. $S$ and $E_{SS}$ correspond the set of sentences and the relationship between pairs of sentences respectively. The weight $w_{ij}$ affiliated with the edge $e_{ij} \in E_{SS}$ is also computed by using the standard cosine measure represented in Equation (1). Furthermore, the definition of weight matrix $W$ is the same as the introduction in Section 3.2.

After that, $W$ is normalized to $\tilde{W}$ to make the sum of each row to be 1:

$$\tilde{W}_{ij} = \begin{cases} W_{ij}/\sum_{j=1}^{|S|} W_{ij}, & if \sum_{j=1}^{|S|} W_{ij} \neq 0 \\ 0, & otherwise \end{cases} . \tag{3}$$

Therefore, the PageRank score $Score_{PR}(s_i)$ for sentence $s_i$ is defined based on the normalized matrix $\tilde{W}$:

$$Score_{PR}(s_i) = \lambda \cdot \sum_{j:j\neq i} Score_{PR}(s_j) \cdot \tilde{W}_{ji} + \frac{(1-\lambda)}{|S|} . \tag{4}$$

For convenience, Equation (4) can be denoted in the matrix form:

$$\overrightarrow{\omega} = \lambda \tilde{W}^T \overrightarrow{\omega} + \frac{1-\lambda}{|S|} \overrightarrow{e} , \tag{5}$$

where $\overrightarrow{\omega}$ is a $|S| \times 1$ vector made up of the scores of sentences in set $S$ and $\overrightarrow{e}$ is a column vector with all the elements equal to 1. $\lambda$ is a damping factor which ranges from 0 to 1 and $(1 - \lambda)$ indicates the probability for node $s_i$ to jump to a random node in the graph.

## 4   Proposed Models

### 4.1   Overview

As mentioned in the Introduction, a concise summary should meet three requirements: diversity, coverage and balance. Traditional graph-based ranking algorithms vote the prestigious vertices based on the global information recursively extracted from the entire graph [9]. However, the document set normally contains a number of different topics with different importance. Computing the sentences in a uniform method would violate the coverage principle because it ignores the differences among topics.

Under the assumption that the sentences containing more important topics should be ranked higher than the sentences contain less important topics, we leverage the topic-level information. In this study, we propose the model to make use of the correlation between topics and sentences. The proposed topic-oriented PageRank (ToPageRank) model follows [8] to divide the random walk into multiple walks and compute the PageRank scores specific to individual random walk, and then combines the entire scores according to the topics distribution. Besides, we modify the basic HITS model, named topic-oriented HITS (ToHITS), to compare the influence of different graph-based algorithms on summarization incorporating the topic-level information.

## 4.2   Topic Detection

The goal of the topic detection is to identify the topics of a given document set automatically. In the experiments, Latent Dirichlet Allocation [1] (LDA) is utilized to detect the topics. LDA is an unsupervised machine learning technique which can identify latent topic information from a document collection. In LDA, each word $w$ in a document $d$ is regarded to be generated by first sampling a topic from $d$'s topic distribution $\theta^{(d)}$, and then sampling a word from the distribution over words $\phi^{(z)}$ which characterizes the topic $z$. Both $\theta^{(d)}$ and $\phi^{(z)}$ have Dirichlet priors with hyper-parameters $\alpha$ and $\beta$, separately. Therefore, the probability of word $w$ in the given document $d$ and prior is represented as follows.

$$p(w|d, \alpha, \beta) = \sum_{z=1}^{T} p(w|z, \beta) \cdot p(z|d, \alpha) \ , \tag{6}$$

where $T$ is the number of topics.

We use GibbsLDA++ toolkit[2] to detect the topics in this study. After iteration, we can get the word-topic distributions, i.e. the probability of each word $w$ on a topic $t$. Since the unit used in multi-document summarization is the sentence, the probability of each sentence $s$ on a topic $t$ should be computed. In the experiments, we choose a heuristic method to compute the contribution degree $degree(s|t)$ of each sentence $s$ on a topic $t$ instead of the probability, i.e. the sum of the probability of every word occurring in the sentence on a topic. Furthermore, to keep the longer sentence from getting the larger value, we normalize the sum by dividing the length of the sentence as shown in Equation (7).

$$degree(s|t) = \frac{\sum_{w \in s} p(w|t)}{|s|} \ , \tag{7}$$

where $|s|$ is the length of the sentence $s$.

## 4.3   Topic-Oriented HITS

In ToHITS model, we build a topic-sentence bipartite graph shown in Figure 1, in which topic nodes represent the hubs and sentence nodes correspond to the authorities.

Formally, the new graph is denoted as $G_{ToH} = (A_{Sent}, H_{Topic}, E_{TS})$, where $A_{Sent} = \{s_i\}$ is the set of sentences and $H_{Topic} = \{t_i\}$ is the set of topics detected by LDA introduced in Section 4.2. $E_{TS} = \{e_{ij} | t_i \in H_{Topic}, s_j \in A_{Sent}\}$ corresponds to the correlations between a sentence and a topic. Each $e_{ij}$ is associated with a weight $w_{ij}$ which indicates the quantitative relationship of the topic $t_i$ and the sentence $s_j$. We compute the contribution degree of the sentence $s$ on the topic $t$ through Equation (7) to denote the relationship. The authority score $Auth(s_i)$ of sentence $s_i$ and the hub score $Hub(t_j)$ of topic $t_j$ are computed by the same iteration process showed in Equation (2). After the iteration converges, the authority scores are used as the scores of sentences.
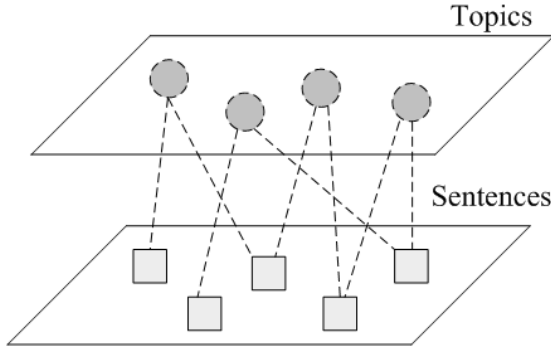
---

[2] http://gibbslda.sourceforge.net/

**Fig. 1.** The topic-sentence bipartite graph in ToHITS model

### 4.4 Topic-Oriented PageRank

In ToPageRank model, we leverage the topic distribution of the documents. Each node in the graph is represented by different topics rather than one [10]. In order to incorporate the topic-level information, we divide traditional PageRank into multiple PageRanks [8] based on different topics, and then the topic-specific PageRank scores is computed in each subgraph. Whereafter, the topic distribution of the entire document set is taken into account, and we can further obtain the final scores of sentences to obtain summaries that are relevant to the documents and at the same time cover the major topics in the documents set. This process is shown in Figure 2.

In PageRank model shown in Equation (4), the probability for a node to jump to a random node is set to a fixed value based on the number of sentences, however the idea in ToPageRank is to run PageRank on each individual topic. Therefore we use the Biased PageRank [8] on each topic separately. The degree value $degree(s|t)$, which is computed according to Equation (7), is assigned to each sentence $s$ in the Biased PageRank on the specific topic. For topic t, the ToPageRank score $Score_{ToPR}^{t}(s_i)$ of sentence $s_i$ is defined as follow:

$$Score_{ToPR}^{t}(s_i) = (1 - \lambda)degree(s|t) + \lambda \cdot \sum_{j:j\neq i} Score_{ToPR}^{t}(s_j) \cdot \tilde{W}_{ji} , \qquad (8)$$

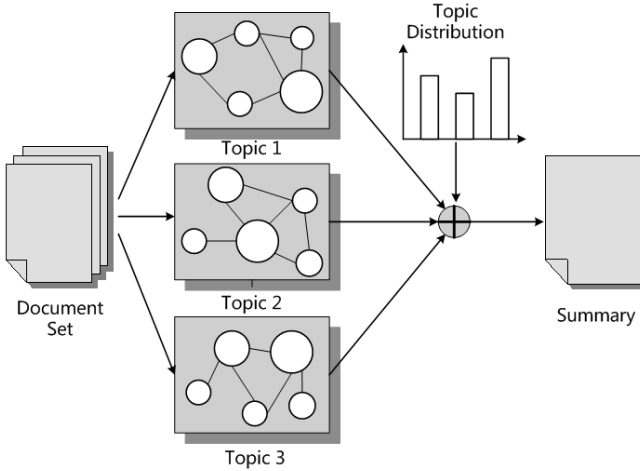where $\tilde{W}$ is the normalized similarity matrix introduced in Equation (3).

The final score $Score_{ToPR}(s)$ of sentence $s$ is computed as follows:

$$Score_{ToPR}(s) = \sum_{t\in T} Score_{ToPR}^{t}(s) \cdot avg_p(t) , \qquad (9)$$

where $T$ stands for all topics and $avg_p(t)$ represents the average value of the sum of probabilities for topic $t$ in the document set. Instinctively every document in the dataset should be treated equally, therefore we obtain $avg_p(t)$ by computing the average value of the sum of probabilities for a topic $t$ in every document:

$$avg_p(t) = \frac{\sum_{d \in D} p(t|d)}{|D|} \quad,$$ (10)

where $p(t|d)$ is the probability of topic $t$ in the document $d$. $D$ is the document set and $|D|$ is the number of documents in the set.



**Fig. 2.** The process of ToPageRank Model (reference to [8])

## 5     Experiments

In order to choose more informative but less redundancy sentences as the final summary, in the experiments we apply the variant version of MMR algorithm proposed in [14]. This method is to penalize the sentences that highly overlap with the sentences that have been chosen as the summary.

### 5.1     Data Set

To evaluate the summarization results empirically, we use DUC2004 dataset since generic multi-document summarization is one of the fundamental tasks in DUC2004. DUC2004 provided 50 document sets and every generated summary is limited to 665 bytes. Table 1 gives a brief description of the dataset.

**Table 1.** Description of the Data Set

|  | DUC2004 |
|---|---|
| **Data source** | TDT[*](Topic Detection and Tracking) |
| **Number of collections** | 50 |
| **Number of documents** | 500 |
| **Summary length** | 665 bytes |

[*] http://www.itl.nist.gov/iad/mig/tests/tdt/

## 5.2   Evaluation Methods

We use the ROUGE [7] toolkit[3] to evaluate these models, which has been widely applied for summarization evaluation by DUC. It evaluates the quality of a summary by counting the overlapping units between the candidate summary and model summaries. ROUGE implements multiple evaluation metrics to measure the system-generated summarization such as ROUGE-N, ROUGE-W and ROUGE-SU.

The ROUGE toolkit reported separate scores for 1, 2, 3 and 4-gram and among these different metrics, unigram-based ROUGE score (ROUGE-1) has been shown to correlate well with human judgments [7]. Besides, longest common subsequence (LCS), weighted LCS and skip-bigram co-occurrences statistics are also used in ROUGE. In this experimental results we show three of the ROUGE metrics: ROUGE-1 (unigram-based), ROUGE-2 (bigram-based), and ROUGE-SU4 (extension of ROUGE-S, which is the skip-bigram co-occurrences statistics) metrics.

## 5.3   Evaluation Results

The proposed ToPageRank model are compared with the ToHITS model, the basic HITS, basic PageRank, the DUC Best performing system and the lead baseline system on DUC2004. The DUC Best performing system is the system with highest ROUGE scores among all the systems submitted in DUC2004. The lead baseline takes the first sentences one by one in the last document in a document set, where documents are assumed to be ordered chronologically. Table 2 shows the comparison results on DUC2004.
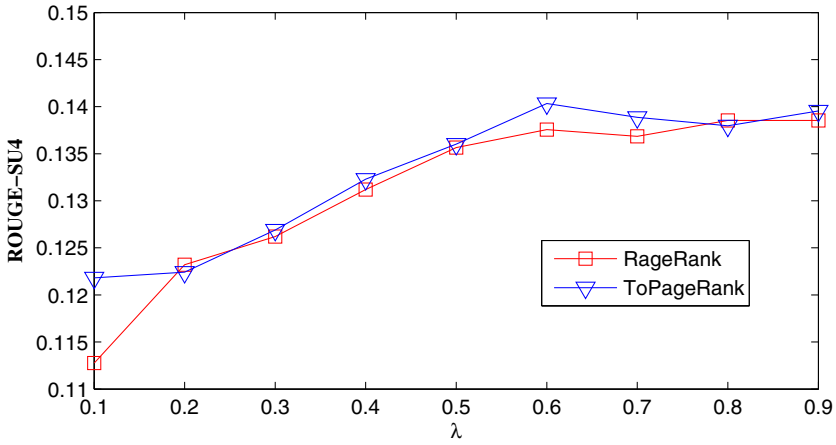
It can be seen that ToPageRank model can outperform the ToHITS model, two basic models and the DUC best system. The results indicate the effectiveness of the proposed ToPageRank model. However, it is worth mentioning that the performance of ToHITS model is better than its corresponding basic HITS model but worse than basic PageRank model. It might result from that in the ToHITS model only the topic-sentence information is applied but sentence relationships are ignored, while in the PageRank model the relationships between sentences play a vital role. This comparison indicates that the sentence-level information is quite important in summarization to some extent.

To investigate how the damping factor and topic number influence the performance of these models, we compare the different combinations with $\lambda$ ranges from 0 to 1 and Figure 3 shows the ROUGE-SU4 curve on the DUC2004. Correspondingly, we further compare the influence of topic number in the models and the results are shown in Figure 4 which indicate that when the topic number is relatively small (around 15) ToPageRank performs well and meanwhile ToHITS prefers relatively larger number of topics (around 20). The trends of other metrics such as ROUGE-1 and ROUGE-2 are similar.
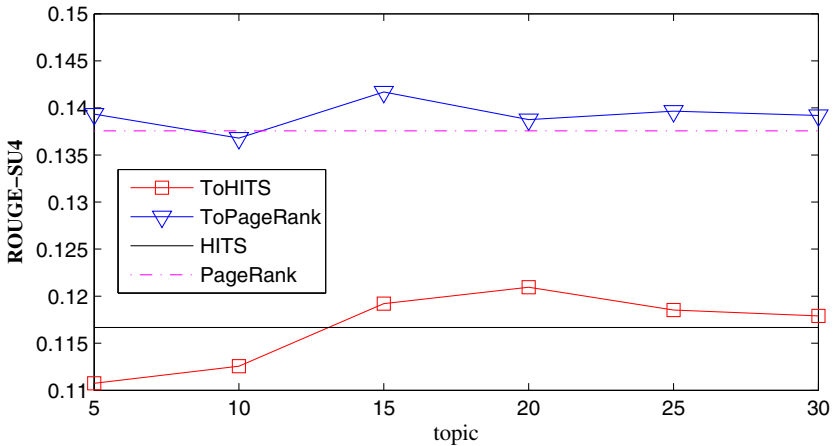
---

[3] ROUGE version 1.5.5 is used in this study which can be found on the website
`http://www.isi.edu/licensed-sw/see/rouge/`

**Table 2.** Comparison results on DUC2004

| Systems | ROUGE-1 | ROUGE-2 | ROUGE-SU4 |
|---------|---------|---------|-----------|
| Lead | 0.33182 | 0.06348 | 0.10582 |
| DUC Best | 0.38279 | 0.09217 | 0.13349 |
| HITS | 0.36305 | 0.06892 | 0.11668 |
| PageRank | 0.38043 | 0.07815 | 0.12473 |
| ToHITS | 0.37264 | 0.07526 | 0.12095 |
| ToPageRank | **0.40501** | **0.09555** | **0.14034** |



**Fig. 3.** ROUGE-SU4 vs. $\lambda$



**Fig. 4.** ROUGE-SU4 vs. *topic*

# 6 Conclusion and Future Work

In this paper, we propose a novel summarization model to incorporate the topic-oriented information in the document set, named ToPageRank. To compare with different algorithms, the ToHITS model is introduced. Experimental results on DUC2004 dataset demonstrate that ToPageRank could outperform the corresponding basic models, ToHITS and the top performing systems in various evaluation metrics.

In this study, the probabilities of sentences on certain topic is computed by accumulating the probabilities of words occurring in the sentence. In future we will design some new model to compute sentence probabilities in a more meaningful method. Moreover, we will explore other correlations between sentences in graph-based ranking methods.

# References

1. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. The Journal of Machine Learning Research 3, 993–1022 (2003)
2. Cai, X., Li, W., Ouyang, Y., Yan, H.: Simultaneous ranking and clustering of sentences: a reinforcement approach to multi-document summarization. In: Proceedings of the 23rd International Conference on Computational Linguistics, pp. 134–142. Association for Computational Linguistics (2010)
3. Erkan, G., Radev, D.R.: Lexpagerank: Prestige in multi-document text summarization. In: Proceedings of EMNLP, vol. 2004, pp. 365–371 (2004)
4. Kleinberg, J.M.: Authoritative sources in a hyperlinked environment. Journal of the ACM (JACM) 46(5), 604–632 (1999)
5. Li, L., Zhou, K., Xue, G.R., Zha, H., Yu, Y.: Enhancing diversity, coverage and balance for summarization through structure learning. In: Proceedings of the 18th International Conference on World Wide Web, pp. 71–80. ACM (2009)
6. Lin, C.Y., Hovy, E.: From single to multi-document summarization: A prototype system and its evaluation. In: Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, pp. 457–464. Association for Computational Linguistics (2002)
7. Lin, C.Y., Hovy, E.: Automatic evaluation of summaries using n-gram co-occurrence statistics. In: Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology, vol. 1, pp. 71–78. Association for Computational Linguistics (2003)
8. Liu, Z., Huang, W., Zheng, Y., Sun, M.: Automatic keyphrase extraction via topic decomposition. In: Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, pp. 366–376. Association for Computational Linguistics (2010)
9. Mihalcea, R., Tarau, P.: Textrank: Bringing order into texts. In: Proceedings of EMNLP, vol. 2004, pp. 404–411. ACL, Barcelona (2004)

10. Nie, L., Davison, B., Qi, X.: Topical link analysis for web search. In: Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 91–98. ACM (2006)
11. Ouyang, Y., Li, W., Lu, Q., Zhang, R.: A study on position information in document summarization. In: Proceedings of the 23rd International Conference on Computational Linguistics: Posters, pp. 919–927. Association for Computational Linguistics (2010)
12. Page, L., Brin, S., Motwani, R., Winograd, T.: The pagerank citation ranking: Bringing order to the web (1999)
13. Radev, D., Jing, H., Stys, M., Tam, D.: Centroid-based summarization of multiple documents. Information Processing & Management 40(6), 919–938 (2004)
14. Wan, X.: Document-Based HITS Model for Multi-document Summarization. In: Ho, T.-B., Zhou, Z.-H. (eds.) PRICAI 2008. LNCS (LNAI), vol. 5351, pp. 454–465. Springer, Heidelberg (2008)
15. Wan, X.: An exploration of document impact on graph-based multi-document summarization. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing, pp. 755–762. Association for Computational Linguistics (2008)
16. Wan, X., Yang, J.: Multi-document summarization using cluster-based link analysis. In: Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 299–306. ACM (2008)
17. Wan, X., Yang, J., Xiao, J.: Towards an iterative reinforcement approach for simultaneous document summarization and keyword extraction. In: Annual Meeting-Association for Computational Linguistics, vol. 45, p. 552 (2007)
18. Wang, D., Zhu, S., Li, T., Chi, Y., Gong, Y.: Integrating clustering and multi-document summarization to improve document understanding. In: Proceeding of the 17th ACM Conference on Information and Knowledge Management, pp. 1435–1436. ACM (2008)