

# A Semantic Aspect-Based Vector Space Model to Identify the Event Evolution Relationship within Topics

Yaoyi Xi\*, Bicheng Li, and Yang Liu

Zhengzhou Information Science and Technology Institute, Zhengzhou, China

WIM\_GY@163.com, lbclm@gmail.com, liuyang198610@163.com

## Abstract

Understanding how the topic evolves is an important and challenging task. A topic usually consists of multiple related events, and the accurate identification of event evolution relationship plays an important role in topic evolution analysis. Existing research has used the traditional vector space model to represent the event, which cannot be used to accurately compute the semantic similarity between events. This has led to poor performance in identifying event evolution relationship. This paper suggests constructing a semantic aspect-based vector space model to represent the event: First, use hierarchical Dirichlet process to mine the semantic aspects. Then, construct a semantic aspect-based vector space model according to these aspects. Finally, represent each event as a point and measure the semantic relatedness between events in the space. According to our evaluation experiments, the performance of our proposed technique is promising and significantly outperforms the baseline methods.

**Category:** Smart and intelligent computing

**Keywords:** Topic evolution; Event evolution relationship; Hierarchical Dirichlet process; Semantic similarity

## I. INTRODUCTION

Topic evolution analysis has recently received a great deal of attention [1-4]. Two key issues for topic evolution analysis are event detection and event evolution relationship identification (hereafter EERI). This paper focuses on event evolution relationship identification. Event evolution relationship refers to the relationship between two events appearing in the same topic; that is, one event has an effect on the happening of the following event, such as a causal relationship, progressive relationship, turning relationship, and so on. Identifying event evolution relationship, which is helpful for people to deeply understand the main development trend of the topic and exactly locate the significant events along with the topic's evolu-

tion, plays an important role in topic evolution analysis. However, it is hard to distinguish between the causal relationship, progressive relationship, turning relationship, and so on. Although there has been research [5, 6] that has tried to identify the specific event evolution relationship type, this research is limited to a few domain-specific data and depends on handcraft rules to determine the type of relationship. Most existing research on EERI [7, 8] has only focused on judging whether there is an evolution relationship between two events, and no attempt has been made to identify its type of relationship. This paper still only identifies if an evolution relationship exists between two events.

Generally speaking, for any two events that have some evolution relationship, the similarity in their content is

**Open Access** <http://dx.doi.org/10.5626/JCSE.2015.9.2.73>

<http://jcse.kiise.org>

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/3.0/>) which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

Received 02 Mar 2015; Accepted 15 May 2015

\*Corresponding Author

larger than the one between the events that have no evolution relationship [7, 8]. Therefore, the general process of identifying the event evolution relationship is as follows: firstly, generate each event's representation using a mathematical model; secondly, calculate the similarity between any two events; finally, sort the similarities in order of decreasing similarity and consider that there is an evolution relationship between two events if their similarity is above a threshold. The key to this process is how to represent an event. Different representation models will directly influence the accuracy of the similarity calculation between events. Currently, all the existing research [7-9] used the traditional vector space model (VSM) [10] to represent an event. The advantages of a VSM-based representation method are simplicity and intuitive graphical representation of events and their similarities, but it has the following limitations: 1) the similarity calculation is only comparative in form, not in meaning; 2) it cannot model linguistic phenomena like synonymy and polysemy; 3) it considers each term in isolation and separates them with the assumption that all terms are orthogonal—inter-independent. In fact, terms in the event are related to each other and several of them are combined together to express certain semantics.

Compared with VSM, the topic model thinks that several aspects generate one document and each aspect corresponds to a set of semantically related words. That is, the topic model does not view isolated terms, but groups semantically related words into clusters. Each cluster is called an aspect, and each aspect captures a certain meaning of the text. In addition, the topic model can be less affected by synonymy and polysemy. Latent Dirichlet allocation (LDA) [11] is a typical topic model. Many of the existing topic models are variations of latent Dirichlet allocation [1, 12]. However, LDA assumes a fixed finite number of aspects. Compared with LDA, the hierarchical Dirichlet process (HDP) model does not need to pre-define the number of aspects and can infer the number automatically [13]. This distinct advantage makes it suitable to data that users do not know about.

For the topic-related document set, the HDP can automatically discover the aspects in it. Although each aspect contains semantic information, it does not correspond directly to an event. Therefore, it cannot use the aspect directly to represent an event. To solve this problem, this paper proposes to construct a semantic aspect-based VSM (SAVSM) to represent an event. The idea of SAVSM is to represent each event as a point in space (i.e., a vector in a vector space). Points that are close together in this space are semantically similar and points that are far apart are semantically distant.

The main contributions of this paper are that it: 1) uses the HDP to mine the latent aspects in documents and the HDP has never been used before in the research of EERI; 2) construct SAVSM to represent an event, which bridges the event and semantic aspect and provides the opportu-

nity to measure the semantic similarity between events.

The rest of the paper is organized as follows. In Section II, we give the definition of the problem. Section III presents related work. Section IV gives a formal description of our EERI approach, while Section V describes the experiments and reports the results. Finally, Section VI summarizes our conclusions.

## II. PROBLEM DEFINITION AND NOTATION

For convenience, we first provide some definitions related to EERI. The definition of the event and topic depends on context and granularity [14]. This study adopts the definition of [8] in the topic and event.

### A. Event Evolution Relationship

Event Evolution is defined as the transitional development process of related events within the same topic along the timeline, and we call the development process from one event to another an event evolution relationship. Formally, an event evolution relationship is defined as the directional logical dependencies or relatedness between two events. The word 'directional' emphasizes that an event evolution relationship is not bidirectional and cannot be reversed. If the occurrence of event B depends on the occurrence of event A, then there is an event evolution relationship from event A to event B [8].

### B. Event Evolution Relationship Identification

Given a set of  $N$  distinct documents  $D = \{d_1, d_2, \dots, d_N\}$  on a given topic  $T$  and their time of publication, assume there are a set of  $M$  events  $E = \{e_1, e_2, \dots, e_M\}$  in  $T$ , where EERI aims to detect the evolution relations between these events under the following constraints:

$$\forall i, e_i \cap D \neq \emptyset \quad (1 \leq i \leq M) \quad (1)$$

$$\forall d_i \exists e_j \in E, d_i \in e_j \quad (1 \leq i \leq N, 1 \leq j \leq M) \quad (2)$$

While the first constraint says that each event is an element in the power set of  $D$ , the second constraint ensures that each document can belong to at least one event.

## III. RELATED WORK

In this section, we review the related work on topic detection and tracking, event evolution relationship identification, event causal relation extraction, and event-time temporal relation learning.

### A. Topic Detection and Tracking

Topic detection and tracking (TDT) is directly related

to EERI. TDT is a research program investigating methods for automatically organizing news documents by the topics that they discuss [15]. Topic detection aims to identify significant topics from a document collection, whereas topic tracking aims to follow the evolution of an identified topic [16]. However, TDT processes a document stream in a very sketchy fashion and it only brings topic-related documents together. In fact, topics can be described in different sizes, making it hard to define the ‘correct’ granularity. In TDT-2004, the topic detection task was replaced by a new task called hierarchical topic detection (HTD), which used a hierarchy to capture more possible granularities [17]. HTD represents topics in a hierarchical structure. This hierarchical structure can help people understand how many events are in the topic, but the relationships between these events are still unknown.

### B. Event Evolution Relationship Identification

Nallapati et al. [7] used event threading to capture the rich structure of events and their dependencies in a news topic. Yang et al. [8] added an event-joining pattern on the basis of event threading. When calculating similarities between events, Nallapati et al. [7] thought that the average of the similarities of all pairs of documents between two events led to a better result, whereas Yang et al. [8] used the average of the term vectors of documents as the event term vector and thought the similarity between two events was the cosine similarity of their event term vectors. These studies made use of VSM to represent the event, lacking excavation of semantic information. Luo et al. [9] tried to mine the semantic information in the event by extracting context-keywords. However context-keywords are hardly enough to describe the deep semantic relationship between two events.

### C. Event Causal Relation Extraction

Event causal relation extraction (ECRE) differs with event evolution relationship identification in this paper [18, 19]. The major difference is that they work with different definitions of ‘Event’. ECRE focuses on sentence-level events, whereas event evolution relationship identification concerns document-level events. Consider the following example:

*He died from fire.*

ECRE thinks that the sentence contains the causal relationship pair of ‘fire-death’; that is, there are two events in this sentence. However, EERI takes the sentence as a whole and believes it only describes an event.

### D. Event-Time Temporal Relation Learning

Event-time temporal relationship learning aims to automatically extract temporal relationships between events

[20, 21]. Compared with EERI, event-time temporal relationship learning only focuses on the temporal relationship between events.

## IV. METHODOLOGY

This section clarifies why and how we propose the SAVSM model for EERI.

### A. Semantic Aspect-Based Vector Space Model

In order to measure the semantic similarity between events, this paper proposes SAVSM to represent an event. SAVSM represents events as vectors in a  $K$ -dimensional space  $R$ , which has only positive axis intercepts.

$$R \in \mathbb{R}_{\geq 0}^K \quad \text{with} \quad K \in \mathbb{N}_{\geq 0}$$

Each dimension of  $R$  represents a so-called fundamental semantic aspect. It is defined for fundamental aspects to be orthogonal (i.e., they are assumed to be inter-independent).

An event  $e_i$  is represented by a semantic aspect vector  $\vec{e}_i$  in the  $R$  vector space.

$$\vec{e}_i = (e_{i1}, e_{i2}, \dots, e_{ik}) \in R,$$

where  $e_{ij}$  represents the weight of the  $j$ th element in  $e_i$ .

In Fig. 1, we provide visualization of the SAVSM operational vector space. Here, the multi-dimensional space consists of three dimensions. Dimensions are constructed from fundamental aspects.

By introducing such a vector space, SAVSM is capable of measuring the semantic similarity between events by calculating the angles between event vectors.

The SAVSM formal operational procedure can be divided

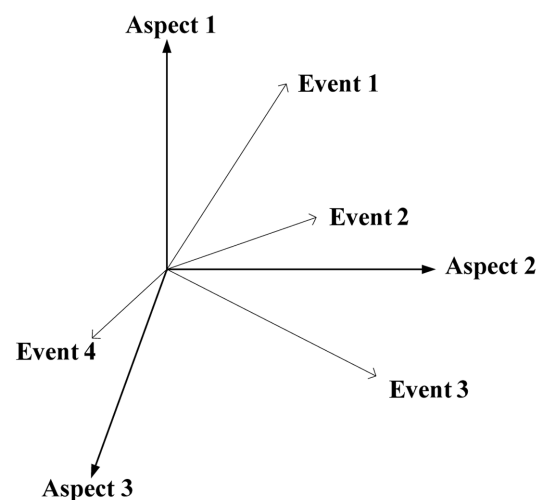


Fig. 1. Semantic aspect-based vector space model visualization.

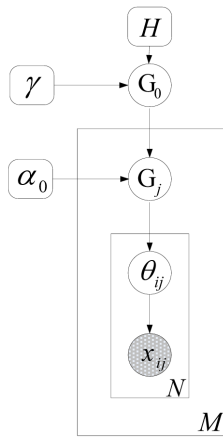


Fig. 2. Hierarchical Dirichlet process.

into three stages. The first stage is document topic modeling. Here semantic bearing aspects are extracted by the HDP. The second stage deals with representing events as vectors. Finally, the third stage is responsible for calculating similarities between events.

1) Document Topic Modeling

The dimensionality of the operational vector space is specified by the number of aspects in the topic. We use HDP to mine the aspects, as shown in Fig. 2.

The generation process of the HDP model is as follows:

1. Draw an overall base measure  $G_0 \sim DP(\gamma, H)$ , which denotes the overall aspect distribution for the topic-related document set  $D$ .
2. For the document  $d_j$  in  $D$ , draw a local measure  $G_j \sim DP(\alpha_0, G_0)$ .
3. For the word  $x_{j,i}$  in  $d_j$ , first draw the aspect assignment  $\theta_{j,i}$  for  $x_{j,i}$ , then sample  $x_{j,i}$  from the aspect corresponding to  $\theta_{j,i}$ .

The posterior inference of the HDP can be found in [13]. Through modeling  $D$  by the HDP, we can get the  $K$  aspects in the topic  $T = \{SE_1, SE_2, \dots, SE_K\}$ . According to the  $K$  aspects, we can construct the  $K$  dimensional semantic vector space  $R$ .

2) Event Model Representation

After modeling  $D$  by the HDP, we can also obtain the aspect distribution of each document. For document  $d_i$ , its aspect distribution is

$$p(\theta_k) = \frac{n_{i,k}}{N_i}, \quad k = 1, 2, \dots, K,$$

where  $N_i$  is the number of words in  $d_i$ ,  $n_{i,k}$  is the number of words generated by aspect  $SE_k$ .

According to its aspect distribution, document  $d_i$  can be represented as a vector in  $R$  and each vector element's weight is the proportion of the corresponding aspect in  $d_i$ .

Formally,

$$\vec{d}_i = (w_{i1}, w_{i2}, \dots, w_{ik}) \in R, \quad w_{ij} = p(\theta_j).$$

Each event may have multiple documents. To obtain the event model, we average the sum of document vectors present in the event and use it as the event model. Assume event  $e_i$  has  $M_i$  documents, then the event model of the event  $e_i \in E$  in SAVSM is represented by event vector  $\vec{e}_i \in R$ .

$$\vec{e}_i = (e_{i1}, e_{i2}, \dots, e_{ik}), \quad \text{where } e_{ij} = \sum_{l=1}^{M_i} w_{lj} / M_i.$$

3) Obtaining Event Similarities

Now, as we have formal event representation (event model) we can obtain event similarities. In most cases the cosine coefficient, which measures the angle between vectors, is used as the similarity measure [10]. Cosine captures the idea that the length of the vectors is irrelevant; the important thing is the angle between the vectors.

In this paper, the similarity between two events  $e_i$  and  $e_j$  is defined as the cosine of the angle between the event vectors.

$$sim(e_i, e_j) = \cos \beta_{e_i, e_j} = \frac{\vec{e}_i \cdot \vec{e}_j}{\|\vec{e}_i\| \|\vec{e}_j\|}.$$

The cosine ranges from 0, when the vectors are orthogonal ( $\theta$  is  $90^\circ$ ), to +1, when they point in the same direction ( $\theta$  is  $0^\circ$ ).

Cosine similarity only considers the event content. However, clearly the time dimension should play an important role in determining event similarity, since events created closer together in time are more likely to discuss the same event than events created at very different moments. This paper introduces the notion of fading similarity to capture both content similarity and time proximity. Formally, we define the fading similarity between a pair of events  $e_i$  and  $e_j$  as

$$sim_f(e_i, e_j) = \frac{sim(e_i, e_j)}{D(\|e_i^T - e_j^T\|)},$$

where  $D(\|e_i^T - e_j^T\|)$  is a distance measure that is monotonically increasing with  $\|e_i^T - e_j^T\|$  and  $D(\|e_i^T - e_j^T\|) \geq 1$ . In this paper, we set  $D(\|e_i^T - e_j^T\|) = e^{\|e_i^T - e_j^T\|}$ .

B. Event Evolution Relationship Identification

We sort the event pairs in order of decreasing fading similarity and consider that there is an evolution relationship between two events if their fading similarity is above a threshold.

## V. EXPERIMENTS AND ANALYSIS

In this paper, we are not aiming at improving existing event detection techniques. Our concern is solely with identification of the event evolution relationship using manually generated and annotated events. Manually generated events can eliminate the biases that may be created by different event detection techniques and, hence, provide a best platform on which we can fairly compare different EERI techniques.

### A. Data Sets

There is no existing standard test set for EERI methods. We randomly choose 4 emergent news topics from four selected news websites, as shown in Table 1. We choose these sites because all of them provide special news about a topic edited by professional editors. Detailed statistics are listed in Table 2.

The major difference between our data set and the corpus adopted in [7] and [8] is that the duration of these topics spans a much longer interval and each topic includes significantly more events and documents in our data set. Therefore, this data set is more close to the actual situation.

After crawling all linked news documents for each topic, we hired two annotators to label the events in each document for each topic independently. The annotators were asked to read the news documents of each topic several times to form a general picture on its development. In the next step, each annotator was asked to identify the events for each topic, and annotate the events in each document independently. The two annotators then met together, reviewed the events and the events in each document annotated individually for each topic, and revised them to arrive at a “consensus” for the topic.

**Table 1.** News sources of four topics

Source	No.
Sina News	380
Tencent News	619
Phoenix News	620
Netease News	143

**Table 2.** Detailed information about four topics

Topic	Size	Events
1. Explosions in the Russian city of Volgograd (December 2013)	137	24
2. Crash of Asiana Airlines Flight 214 in San Francisco (July 2013)	628	128
3. Kunming terror attack (March 2014)	270	49
4. Boston Marathon bombing (April 2013)	727	74

We remove common stop-words and tokens, which are neither verbs, nor nouns, nor adjectives from the news documents by the National Language Processing & Information Retrieval (NLPIR) sharing platform (<http://www.nlpir.org/>).

### B. Evaluation Metrics

This study adopts the measurement of precision and recall for our evaluation, which are also used in [7] and [8].

We denote the event evolution relationship from event  $e_i$  to  $e_j$  (where  $i \neq j$ ), if existing and valid, as  $(e_i, e_j)$ . Suppose there are  $L$  truth event evolution relationships and denote them as  $LTruth$ . A system can generate  $L'$  event evolution relationships and denote them as  $LSystem$ . Since events are predefined in this paper, our evaluation focuses on the differences between the sets of event evolution relationships  $LTruth$  and  $LSystem$ .

Generally speaking, users expect to have as many correct event evolution relationships identified as possible in  $LSystem$ , but usually it is at the cost of a fast increasing volume of returned results. In that case, users need to spend much time on distinguishing between correct event evolution relationships and incorrect ones. Thus, we need to measure both the number of correct event evolution relationships identified and the total size of returned results and allow users to make a balance between them [8].

The identified correct set  $LRelevant$  is the overlapping part of  $LTruth$  and  $LSystem$ .

$$LRelevant = LTruth \cap LSystem .$$

We define precision and recall as follows:

$$Precision = \frac{|LRelevant|}{|LSystem|} ,$$

$$Recall = \frac{|LRelevant|}{|LTruth|} .$$

Obviously, we can tune the input parameters of our models and then generate different sets of precision and recall rates. Theoretically, if the tuning of model parameters is continuous, we can plot a smooth precision and recall curve. Here it is enough to know that perfect performance would be a curve close to the upper and right boundaries of the graph.

To measure global performance, two averaging methods are used. The micro-average method assigns equal weight to each decision for each relationship and accumulates the precision/recall probabilities over all topics. The macro average accumulates the precision/recall probabilities separately for each topic and then averages the probabilities over topics, with equal weight assigned to each topic. Existing research [22] has shown that the macro average method can provide better estimates of performance, so this paper adopts the topic-weighted method to estimate the EERI performance.

### C. Algorithms for Comparison

We implement the following widely used EERI algorithms as the baseline systems. To these methods we take the average score as their performance. For fairness we conduct the same preprocessing for all algorithms.

- Event threading [7]: The method uses the average of the similarities of all pairs of documents between two events as the event similarity.
- Event evolution graph [8]: The method uses the average of the term vectors of documents as the event term vector and thinks the event similarity is the cosine similarity of their event term vectors.
- Event river [9]: The method extracts the context-keywords of an event as its representation, and calculates the event similarity based on its context-keywords using the Jaccard coefficient.

Our proposed algorithm considers two components: event content similarity and temporal proximity. Thus, our proposed algorithm is tested as SAVSM-TD for the method considering time proximity between events and SAVSM that does not consider temporal proximity.

Fig. 3 shows the evaluation results of these methods.

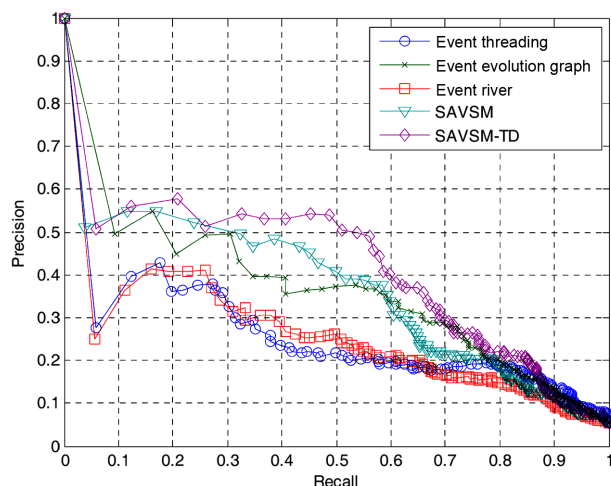


Fig. 3. Comparison of different EERI methods. EERI: event evolution relationship identification, SAVSM: semantic aspect-based vector space model, TD: topic detection.

We observe that both the SAVSM and SAVSM-TD methods substantially outperform the other methods. Event river extracts context-keywords to represent the event and uses the Jaccard coefficient to measure the event similarity. The results indicate that it does not work well in EERI. This is mainly because the diversity of event context-words within a topic is not obvious. Besides, computing event similarity based on the number of context-words in the intersection set will only lose some useful information, such as term frequency, and so on. The result of the event evolution graph is better than event threading and event river. The reason in this case may be that the event evolution graph excludes the IDF factor when representing events, which makes the event representation more representative. In addition, it uses event term vectors for measuring event similarity, which highlights the role of some representative terms and so outperforms the average pair-wise similarity measures that event threading adopts. Event threading has the worst performance. Besides the poor performance of the event representation and similarity measure it adopts, another main reason is that event threading assumes that the evolution relationship only exists between two time-adjacent events. As a result, many evolution relationships are missed.

Both SAVSM and SAVSM-TD outperform baselines, indicating that the event representation model and similarity measure we propose for EERI are beneficial. Some event pairs that have an evolution relationship may have little lexical congruence. Take the following two events about San Francisco air disaster as an example:

*e<sub>1</sub>: The Chinese consulate in San Francisco verifies the Chinese casualties.*

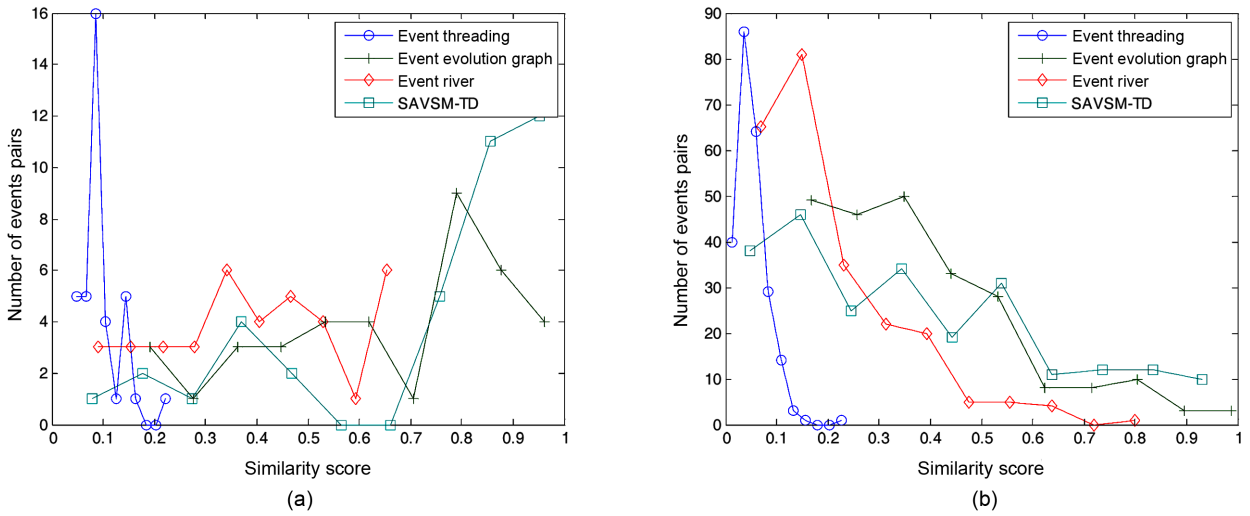
*e<sub>2</sub>: The Chinese embassy in South Korea requires Asiana Airlines to report details about the people on board as soon as possible.*

The two events have a progressive relationship in practice, indicating the Chinese government’s effort to verify the casualty information. However, the traditional methods cannot correctly identify their relationship in general due to the small number of same words, whereas the SAVSM can mine the latent semantic information in the events and help calculate the event similarity accurately. As a result it improves the performance of EERI.

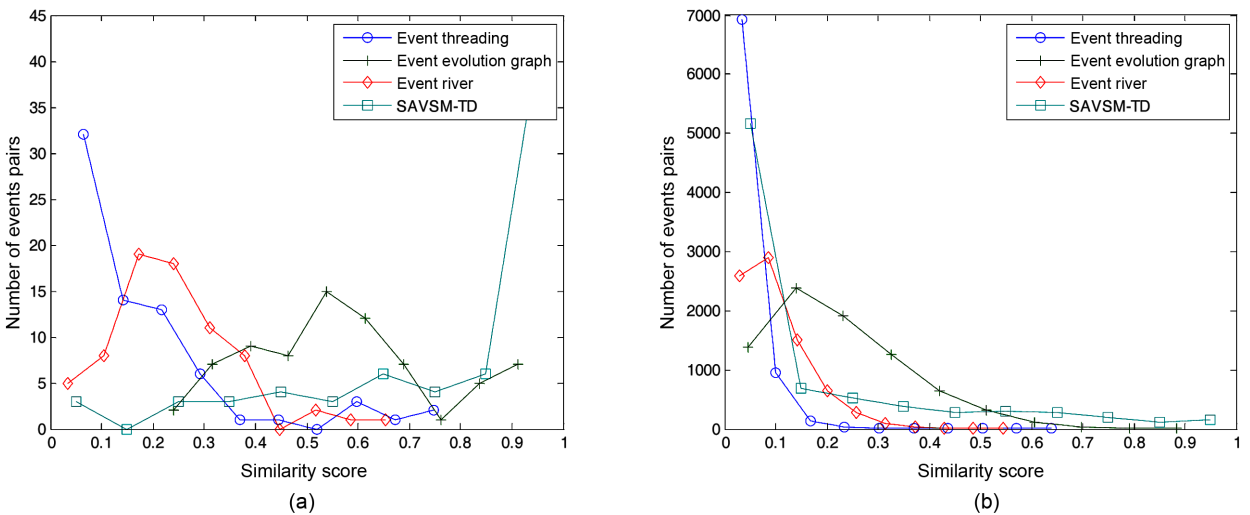
SAVSM-TD performs better than SAVSM. This suggests that temporal proximity is helpful for evaluating event evolution relationships.

For the purposes of convenience for description, we denote the event pairs that have evolution relationships as real event pairs and the ones that do not have evolution relationships as false event pairs. Generally speaking, the similarities between real event pairs are larger than the false ones. As such, a good similarity measure can separate the real event pairs from the false ones. The objective of our second experiment is to evaluate the effectiveness of our proposed event similarity measure.

The main goal of our effort was to come up with a way



**Fig. 4.** Event pairs similarity score distribution in topic 1: (a) real event pairs, (b) false event pairs. SAVSM: semantic aspect-based vector space model, TD: topic detection.



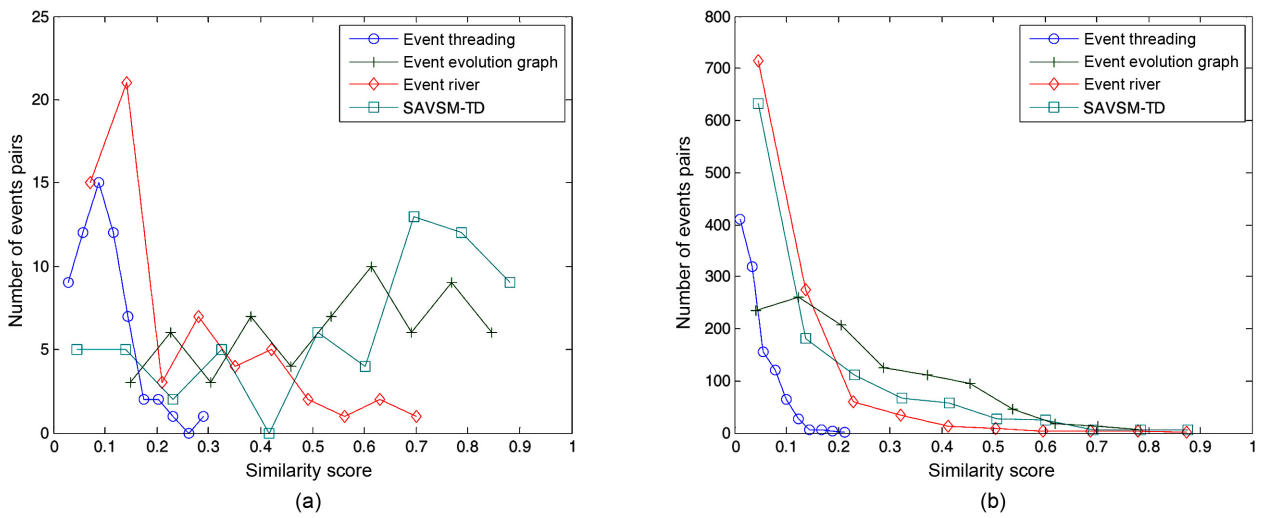
**Fig. 5.** Event pairs similarity score distribution in topic 2: (a) real event pairs, (b) false event pairs. SAVSM: semantic aspect-based vector space model, TD: topic detection.

to correctly measure the semantic similarity between events. To understand what we had actually achieved by using this method, we studied the distribution of the similarity scores assigned to real and false event pairs for the baselines and SAVSM-TD (Fig. 4–7).

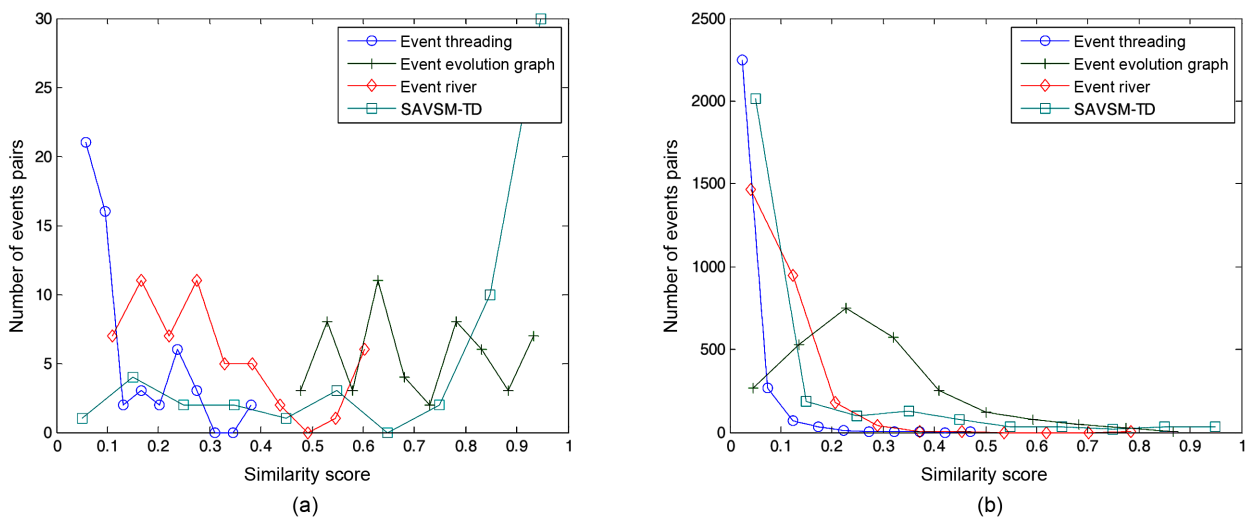
Figs. 4, 5, 6, and 7 correspond to topic 1, 2, 3, and 4, respectively. The left of each figure is the similarity score distribution of real event pairs and the right is the one of false event pairs. From these figures, we observe that event threading performs worst on every topic. The similarities between real event pairs do not differ much from the false ones. As shown in Fig. 4, the big fraction of similarities computed by event threading ranges from 0.1 to 0.2, which has poor distinctiveness. In contrast, results

showed that the SAVSM-TD method has high performance to differentiate the similarities between event pairs and the similarities of a significant number of real event pairs have been increased to be closer to one. The differentiation ability of the other two methods is somewhere between SAVSM-TD and event threading.

From these figures, we can see that using SAVSM-TD to measure the semantic similarities between events is effective. Existing research depends too much on the lexical similarity, while SAVSM-TD tries to mine the latent semantic information in events and compare the semantic similarity between them, which can effectively distinguish the event pairs that are close in meaning from the semantically different ones.



**Fig. 6.** Event pairs similarity score distribution in topic 3: (a) real event pairs, (b) false event pairs. SAVSM: semantic aspect-based vector space model, TD: topic detection.



**Fig. 7.** Event pairs similarity score distribution in topic 4: (a) real event pairs, (b) false event pairs. SAVSM: semantic aspect-based vector space model, TD: topic detection.

## VI. CONCLUSION AND FUTURE WORK

Identifying event evolution relationships helps to understand how the topics evolve along the timeline. Considering that traditional methods cannot accurately compute the semantic similarity between events and, thus, reduce performance of EERI, this paper proposes constructing a SAVSM to represent the event with the help of a topic model. The semantic similarity between events are effectively measured based on this representation, which improves the performance of EERI.

The paper develops the event evolution relationship identification technique, but it is carried in an offline condition and under the premise that all the events are cor-

rectly detected. In the future, we plan to combine the event detection and evolution relationship identification together. We also hope to devise an online algorithm for identifying an event evolution relationship.

## REFERENCES

1. L. Huang and L. Huang, "Optimized event storyline generation based on mixture-event-aspect model," in *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Seattle, WA, 2013, pp. 726-735.
2. A. Ahmed and E. P. Xing, "Timeline: a dynamic hierarchi-



- cal Dirichlet process model for recovering birth/death and evolution of topics in text stream,” in *Proceedings of the 26th Conference on Uncertainty in Artificial Intelligence (UAI2010)*, Catalina Island, CA, 2010.
3. J. H. Lau, N. Collier, & T. Baldwin, “on-line trend analysis with topic models: twitter trends detection topic model online,” in *Proceedings of COLING 2012: Technical Papers*, Mumbai, India, 2012, pp. 1519-1534.
  4. P. Lee, L. V. Lakshmanan, & E. E. Milios, “Event evolution tracking from streaming social posts,” <http://arxiv.org/pdf/1311.5978v1.pdf>.
  5. A. Feng and J. Allan, “Finding and linking incidents in news,” in *Proceedings of the 16th ACM Conference on Information and Knowledge Management (CIKM'07)*, Lisbon, Portugal, 2007, pp. 821-830.
  6. A. Feng and J. Allan, “Incident threading for news passages,” in *Proceedings of the 18th ACM Conference on Information and Knowledge Management (CIKM'09)*, Hong Kong, China, 2009, pp. 1307-1316.
  7. R. Nallapati, A. Feng, F. Peng, and J. Allan, “Event threading within news topics,” in *Proceedings of the 13th ACM International Conference on Information and Knowledge Management (CIKM'04)*, Washington, DC, 2004, pp. 446-453.
  8. C. C. Yang, X. Shi, and C. P. Wei, “Discovering event evolution graphs from news corpora,” *IEEE Transactions on Systems, Man and Cybernetics, Part A: Systems and Humans*, vol. 39, no. 4, pp. 850-863, 2009.
  9. D. Luo, J. Yang, M. Krstajic, W. Ribarsky, and D. Keim, “EventRiver: visually exploring text collections with temporal references,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 18, no. 1, pp. 93-105, 2012.
  10. P. D. Turney and P. Pantel, “From frequency to meaning: vector space models of semantics,” *Journal of Artificial Intelligence Research*, vol. 37, no. 1, pp. 141-188, 2010.
  11. D. M. Blei, A. Y. Ng, & M. I. Jordan, “Latent Dirichlet allocation,” *Journal of Machine Learning Research*, vol. 3, pp. 993-1022, 2003.
  12. J. Y. Delort and E. Alfonseca, “DualSum: a topic-model based approach for update summarization,” in *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, Avignon, France, 2012, pp. 214-223.
  13. Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei, “Hierarchical Dirichlet processes,” *Journal of the American Statistical Association*, vol. 101, no. 476, pp. 1566-1581, 2006.
  14. K. N. Vavliakis, A. L. Symeonidis, and P. A. Mitkas, “Event identification in web social media through named entity recognition and topic modeling,” *Data & Knowledge Engineering*, vol. 88, pp. 1-24, 2013.
  15. J. Allan, *Topic Detection and Tracking: Event-Based Information Organization*. Boston, MA: Kluwer Academic Publishers, 2002.
  16. W. Ding and C. Chen, “Dynamic topic detection and tracking: a comparison of HDP, C-word, and cocitation methods,” *Journal of the Association for Information Science and Technology*, vol. 65, no. 10, pp. 2084-2097, 2014.
  17. A. Feng and J. Allan, “Hierarchical topic detection in TDT-2004,” Center for Intelligent Information Retrieval Technical Report, 2005.
  18. Z. Kozareva, “Cause-effect relation learning,” in *Workshop Proceedings of TextGraphs-7 on Graph-Based Methods for Natural Language Processing*, Jeju Island, Korea, 2012, pp. 39-43.
  19. J. Fu, Z. Liu, W. Liu, and Q. Guo, “Using dual-layer CRFs for event causal relation extraction,” *IEICE Electronics Express*, vol. 8, no. 5, pp. 306-310, 2011.
  20. S. A. Mirroshandel, M. Khayyamian, and G. Ghassem-Sani, “Syntactic tree kernels for event-time temporal relation learning,” in *Human Language Technology: Challenges for Computer Science and Linguistics*. Heidelberg: Springer, pp. 213-223, 2011.
  21. S. A. Mirroshandel and G. Ghassem-Sani, “Towards unsupervised learning of temporal relations between events,” *Journal of Artificial Intelligence Research*, vol. 45, pp. 125-163, 2012.
  22. National Institute of Standards and Technology, “The 2004 Topic Detection and Tracking (TDT2004) task definition and evaluation plan,” <http://www.itl.nist.gov/iad/mig/tests/tdt/2004/TDT04.Eval.Plan.v1.2.compare.1.1c.pdf>.



### Yaoyi Xi

Yaoyi Xi received his M.S. degree in signal analysis and processing from the Zhengzhou Information Science and Technology Institute in 2011. He has been working on theoretical research and the experimental application of text processing technologies. The main research in his Ph.D. program is around topic evolution analysis, topic detection and tracking, and multi-document summarization.



### **Bicheng Li**

---

Bicheng Li teaches at Zhengzhou Information Science and Technology Institute as a professor and Ph.D. supervisor. His teaching subjects include pattern recognition and artificial intelligence, and wavelet transformation. He has published many professional books on various topics including information fusion and application, and pattern recognition principles and their application. His research fields include text analysis and understanding, speech/image/video processing and recognition, and information fusion.



### **Yang Liu**

---

Yang Liu received his M.S. degree in communications and information systems from Zhengzhou Information Science and Technology Institute in 2011. He has been working on theoretical research and experimental application of text processing technologies. The main research in his Ph.D. program is around social network analysis and data mining.