

# Personalized news recommendation based on articles chain building

Wanrong Gu<sup>1</sup> · Shoubin Dong<sup>1</sup> · Mingquan Chen<sup>1</sup>

Received: 22 September 2014 / Accepted: 3 June 2015 / Published online: 16 June 2015  
© The Natural Computing Applications Forum 2015

**Abstract** News personalized recommendation has long been a favorite research in recommender. Previous methods strive to satisfy the users by constructing the users' preference profiles. Traditionally, most of recent researches use users' reading history (content based) or access pattern (collaborative filtering based) to recommend newly published news to them. In this way, they only considered the relationship between news articles and the users and ignored the context of news report background. In other words, they fail to provide more useful information with considering the progression of the news story chain. In this paper, we propose to define the quality of a news story chain. Besides, we propose a method to construct a news story chain on a news corpus with date information. At last, we use a greedy selection method for filtering the final recommended news articles with considering accuracy and diversity. In this way, we can provide the news articles for users and meet their requirement: after reading the recommended news, the user gains a better understanding of the progression of the news story they read before. Finally, we designed several experiments compared to the state-of-the-art approaches, and the experimental results show that our proposed method significantly improves the accuracy, diversity and NDCG metrics.

**Keywords** News story chain · News recommendation · Hybrid recommender

---

✉ Shoubin Dong  
sbdong@scut.edu.cn

<sup>1</sup> School of Computer Science and Engineering, South China University of Technology, Guangzhou 510641, People's Republic of China

## 1 Introduction

As the World Wide Web becomes the source and distribution channels of news, more and more users read news online rather than buying newspaper or watching TV. However, massive online news leads to *information overload* problem and has long been recognized in the computing research and industry. People are constantly faced to the larger and larger amounts of internet data every day. With the massive data, users are often missed the whole picture of the news articles.

Traditional approaches strive to satisfy the users by constructing the users' preference profiles and recommend the news articles based on the reading history news content or the collaborative relationships between users. In this way, the preference of user and the news content are seem to be static. For this reason, personalized news recommendation needs to be researched to present data in a meaningful and effective manner. In this paper, we investigate personalized news recommendation via news articles chain—providing a useful personalized news recommendation utilized the news story chain building. We believe that the user's reading preference is affected by the news stories tracking.

We focus on the personalized news recommendation domain: given two sets of news articles from a user, in which one is the begin set of the user's reading news, and the other is the end set of this user's reading news, our method automatically finds a coherent recommendation list linking them together. In this way, we can recommend the middle news articles, which are likely to meet the user's preference to the user. In our work, we propose a method for producing a news articles list. The list should be useful for the user: After reading this list, the user may gain a better understanding of the whole picture of the news story.

To the best of our knowledge, the news recommendation via constructing the chain of news is novel. Previous studies (e.g., [1–4]) focused on recommending news articles to the user based on the user's profile or the user access pattern, but did not address the notion of the development of news story.

Our main contributions are as follows:

- *We define an equation for measuring the quality of a news story chain.* The quality of a news story chain can be effected by its measurement. We explored a measurement of the news story chain quality with considering multiple factors.
- *We proposed to model the chain in news articles instead of traditional recommendation.* In this way, there was no information loss in news personalization recommendation. Besides, the recommended reason based on news story chain was better than the previous studies. In this way, the news recommendation could recommend news articles with news story information to the users.
- *We proposed and implemented a method of news story chain building.* This is an important part of this paper. We explored and simulated the suitable and feasible approach of news story building. This method can also be used in the other applications and researches.

The remaining of this paper is organized as follows. Section 2 covers related work relevant to our study, including content-based news recommender, collaborative filtering-based method and hybrid approach. In Sect. 3, we describe the problem statement in this paper. And then, we introduce the framework we proposed in Sect. 4. Section 5 shows the experiments and results compared to the state-of-the-art approaches. Finally, we conclude this paper and discuss the future work in Sect. 6.

## 2 Related work

In this section, we provide a brief review of related work about news recommender. The previous studies can be roughly divided into three categories: content-based, collaborative filtering and hybrid approaches.

### 2.1 Content-based methods

Content-based news recommender is a classical method, which construct user profile based on the news content and recommend news articles to meet the preference of user [5]. In practice, a piece of news is often represented as vector space model. Each key word has a weight in this news article vector. The user profile is constructed by the history reading news articles. At last, the recommender

calculates the similarity between user profile and newly published news. Some previous works belong to this category. For example, [6] utilized K-nearest neighbor method to recommend news articles for the given user. [7] employed the Naïve Bayesian method to classify web pages and construct user profile, and then recommend web pages to the user based on the similarity between web page and user profile. Liu et al. [2] (called ClickB in experimental section) proposed a recommender using news article content tracking on click behavior. However, the construction of user profile based on news content would be adhesive in a special center in higher-dimensional space. In this way, the recommended list would not provide a good diversity, and the user would be recommended to read many news articles very similar to his/her read before.

### 2.2 Collaborative filtering-based methods

Collaborative filtering-based recommender utilizes the collaborative access pattern among users and items. By this way, this method seems to be content-free and can be roughly divided into two categories: Heuristic-based and Model-based. The former is inspired by the real-world phenomena [8], and the latter trains a model for predicting the utility of the current user  $u$  on item  $i$ , such as [9] and [3] (called Goo in experimental section). In news recommendation practice, the rating of user on the news can be seen as binary, where a click or a comment on a piece of news can be rated as 1, and 0 otherwise [3]. However, there are many newly published news articles enter the recommender system. It suffers from the well-known *cold-start* problem and needs to utilize the news article content for recommending.

### 2.3 Hybrid approaches

Hybrid recommender combines content-based method and collaborative filtering-based method, considering other factors [10] to recommend. Many recent studies are hybrid, such as Bilinear [11], Bandit [12], SCENE [1] and TwoHy [4], which will be discussed and analyzed in experimental section. From the perspective of news recommender, our work is similar to EMM News Explorer [13], Newsjunike [14] and TwoHy [4]. However, EMM News system did not provide personalized recommender; Newsjunike and TwoHy did not address the news story tracking in recommendation.

Although the previous approaches have achieved great success in personalized news recommender applications, they failed to make full use of the news story chain in model news recommendation. In our work, we propose to use to construct the news story chain and apply it into our recommendation model. In this way, we can provide a reasonable recommended list to user.

There are many related works about *narrative generation* community [15, 16], *event tracking* [17] and *event detection* [18, 19]. For the narrative generation researches, researchers strived to explore the ways to model the narratives. For event tracking studies, they classified news articles into categories by machine learning. In this way, they needed to insure that the labels are known beforehand. At last, for the event detection-related works, they dealt with discovering the new news events. In contrast, our contribution lies in supplying the appropriate news articles through the news story chain building. Nevertheless, some studies also provide us some useful ideas for our research [20].

### 3 Problem statement

The problem of this paper is personalized news recommendation, and we can give the definition as follows:

**Problem 1** (*Personalized news recommendation*) Given a collection of news  $\mathcal{N}$  and a given user  $u$ , recommend a result of news articles  $\mathcal{C}$  that maximally match  $u$ 's satisfy, including preference, novelty and freshness.

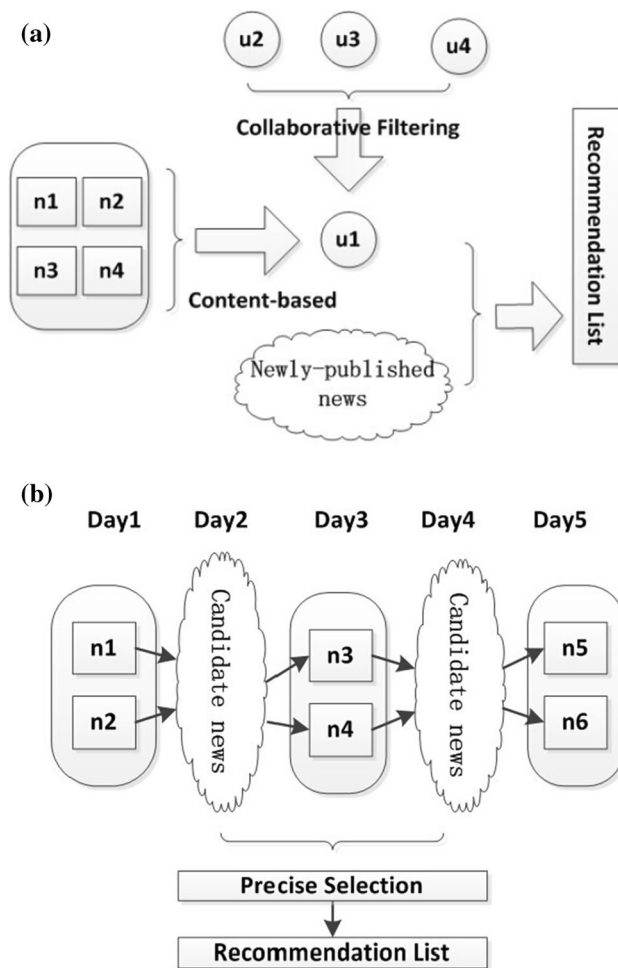
In some previous studies, they solved this problem using the similarity of user profile and the newly published news articles, see [2]. Some previous researches proposed to use the collaborative filtering method to find the similar users or news topics, and then, recommended news articles to the user refer to considering the similar users or news topics [3]. In contrast to the previous studies, from Fig. 1, we can see that it has not any temporal information in traditional news recommendation. In our proposed model, the recommended news articles are closely linked with the user read before in temporal dimension. Besides, it can provide a coherent news story for user. Let  $\mathcal{D}$  be a set of news articles, given a chain  $(d_1, \dots, d_n)$  from  $\mathcal{D}$ , similar to [21], we can estimate its coherence from the following definition:

$$\text{Coherence}(d_1, \dots, d_n) = \min_{i=1, \dots, n-1} \sum_{w \in (\text{active in } d_i, d_{i+1})} \text{Influence}(d_i, d_{i+1} | w), \tag{1}$$

where  $\text{Influence}(d_i, d_{i+1} | w)$  denotes the influence of key word  $w$  in the link between  $d_i$  and  $d_{i+1}$ , and we will discuss in Sect. 4.

**Subproblem 1** (*News story chain building*) Let  $s$  and  $t$  be two news articles, and  $\mathcal{D}$  be a set of news articles with temporal information between  $s$  and  $t$ . Select a chain  $(d_1, \dots, d_n)$  from  $\mathcal{D}$ .

**Subproblem 2** (*News selection*) Given a collection of news articles  $\mathcal{C}$  and a user  $u$ , filter a subset  $\mathcal{C}^*$  from the set  $\mathcal{C}$



**Fig. 1** An illustrative example of personalized news recommendation problem focuses on traditional approach and our proposed method. Remark  $n$  denotes a piece of news,  $u$  denotes a user. **a** Traditional recommendation. **b** Our recommendation

where  $\mathcal{C}^*$  match the  $u$ 's reading preference optimally and the diversity of this subset is maximized.

In the previous researches, the two subproblems were not considered synthetically. In this paper, these two subproblems were regarded as the necessary parts of the personalized news recommendation.

## 4 Proposed framework

### 4.1 News story chain building

In this section, we discuss how to find a good chain between two news articles. And then, we will analyze how to measuring the influence of this chain. Finally, we introduce how to construct the chain in our proposed recommendation.

(a) s: The microblog about Zhang Wen falls in love with Di Yao has gone viral.

A1: Di Yao endorsed South Korean tablet computer and showed the advertisement on her microblog.

A2: Microsoft new Win RT tablet Metro will cancel the desktop function.

A3: Promotions are confirmed by Microsoft as a "blunder".

A4: Two girls fought due to the tiny blunder of their rumored boyfriend.

t: Rumored boyfriend of Di Yao says: we should not slander her due to the tiny blunder.

	s	A1	A2	A3	A4	t
Di Yao	✓	✓				✓
Microblog	✓	✓				
Microsoft			✓	✓		
Blunder				✓	✓	✓
Boyfriend					✓	✓
Zhang Wen	✓					
Tablet		✓	✓			

(b) s: The microblog about Zhang Wen falls in love with Di Yao has gone viral.

B1: Di Yao and Zhang Wen are traveling arm in arm.

B2: Proof of Zhang Wen's infidelity: hug Di Yao in the street.

B3: Zhang Wen and Ma Yili appeared and claimed not to divorce.

B4: The status of Di Yao: did not affect the filming.

t: Rumored boyfriend of Di Yao says: we should not slander her due to the tiny blunder.

	s	B1	B2	B3	B4	t
Di Yao	✓	✓	✓		✓	✓
Zhang Wen	✓	✓	✓	✓		
Boyfriend						✓
Ma Yili				✓		
Divorce				✓		
Microblog	✓					
Infidelity			✓			

◀Fig. 2 Two examples of news story chains connecting the same endpoints. Left chain created by shortest path using similarity weight. Right a good story chain. a News story chain 1#. b News story chain 2#

### 4.1.1 How to build a good story chain?

In this section, the goal is to build a good story chain between two given news articles,  $d_s$  and  $d_t$ . A natural way is to construct an undirected weighted complete graph over the news articles and find a shortest  $d_s - d_t$  path, in which each edge is linked between each pair of two articles with the weight of similarity. However, this approach does not necessarily build a good story chain. For example, we try to construct a news story chain between  $d_s$ : *The microblog about Zhang Wen falls in love with Di Yao has gone viral* (2014.03.28) and  $d_t$ : *Rumored boyfriend of Di Yao says: we should not slander her due to the tiny blunder* (2014.07.25), and then find two story chains between them. The result is shown in Fig. 2.

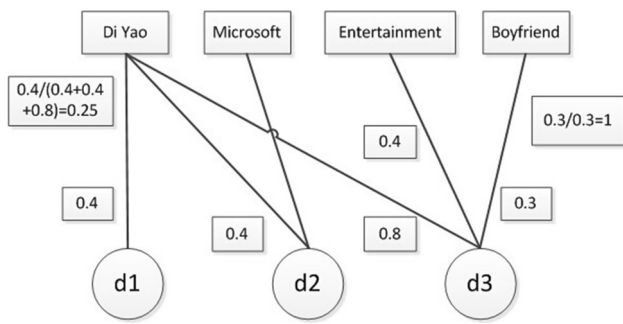
The straightforward reason may lie with the locality of approach. The method ensures that every two consecutive news articles are related, but there is no the best news story chain in global environment. We can take a closer analysis at the two chains. In the good chain in Fig. 2, for example, the named entity word *Di Yao* and *Zhang Wen* appeared almost throughout the whole good chain. But in the bad chain, some key words appear and then disappear for a long period. Besides, the key words appeared more evenly in the bad chain. This observation motivates that we need to compute the weight of a key word, measure the influence of a chain and propose a definition of coherence in a chain.

### 4.1.2 Measuring influence of a chain

Before constructing a news story chain, it requires evaluation of  $Influence(d_i, d_j|w)$ —the influence of  $d_i$  on  $d_j$  w.r.t. key word  $w$ . Several approaches for measuring news story chain influence have been proposed, such as random graph simulations [22], authority computation [23] and random works method [24].

In this work, we explore a different method for measuring the influence. First of all, we define a bipartite undirected graph,  $G = (V, E)$ . The vertices set  $V = V_N \cup V_W$  correspond to the news articles and key words. For each key word  $w$  in a news article  $n$ , we add an edge  $(w, n)$ . Refer to Fig. 3 for a simple bipartite graph: there are three (circle) news articles and four (square) key words.

The edge weights represent the value of the correlation between a news article and a key word in this document. We used the tool GATE [25] for named entity extraction and used the extracted named entities as the key words. We need to do some preprocesses for measuring the influence



**Fig. 3** A bipartite graph used to measure influence

as the following two stages: first, we assign TF-IDF weight for each candidate key word on a article. And then, we normalize them over all key words in the article. For instance, the  $d_3$  document is connected to *Di Yao* (0.8), and somewhat about *Entertainment* (0.4) and *Boyfriend* (0.3). The word *Boyfriend* can only be reached through a single news article. Therefore, the normalized weight of this edge is  $\frac{0.3}{0.3} = 1$ . Similarly, the normalized edge weight of the word *Di Yao* to document  $d_1$  is  $\frac{0.4}{0.4+0.4+0.8} = 0.25$ . We now utilize this normalized weighted graph to calculate influence between news articles.

As analysis before,  $\text{Influence}(d_i, d_j|w)$  should be high if the two articles are highly connected and the word  $w$  plays an important role in this connection. Normally, if two documents are tightly connected, a random walk starting from  $d_i$  should reach  $d_j$  frequently. We first calculate the initial distribution for random walks starting from  $d_i$ . And then, we control the expected random walks restart probability,  $\epsilon$ . The distribution can be defined as follows:

$$\Pi(v) = \begin{cases} \epsilon \sum_{(u,v) \in E} \frac{\Pi(u)}{P(v|u)} & (v \neq d_i) \\ (1 - \epsilon) \sum_{(u,v) \in E} \frac{\Pi(u)}{P(v|u)} & (v = d_i) \end{cases}, \quad (2)$$

where  $P(v|u)$  denotes the probability of reaching  $v$  from  $u$ , and we use the normalized edge weight instead.

### 4.1.3 Finding a good story chain

In the above sections of this paper, we discussed an approach to score a news development chain in a corpus set. In this section, we propose the method of finding a good chain in a news reading environment.

Similar to Sect. 3, we formulate the problem as finding the chain which maximizes our scoring function. We use the bisearch optimization method to build the chain as follows:

1. Find document  $m$  in the middle of the day in  $s$  and  $t$  and insure that the  $\text{Coherence}(s, m, t)$  is max.

2. In the starting point  $s$  and ending point  $m$ , repeat the step (1) and find the middle point  $m'$ . Similarly, we can find the middle point  $m''$  in  $m$  and  $t$ .
3. Repeat the bisearch process, until each date is filled with a piece of news.

If a user has read more than a piece of news, for example, a user read 2 articles on day  $s$  and 3 articles on day  $t$ . We need to construct  $2 \times 3 = 6$  news story chain for the pair of  $s$  and  $t$ . Besides, a user may read a lot of news articles on one day, i.e., more than 3 articles. In this condition, we use k-means method to cluster the news articles and set the  $k$  to 3. Each cluster represents a piece of news described above.

## 4.2 News selection

For a given user, if he/she read a lot of news articles every day, it would product many candidate news articles for him/her. Besides, the candidate news would appear many similar news articles. Based on this intuition, our news article selection method can be described as:

Assuming that  $\mathcal{C}$  denotes the candidate news article set,  $\mathcal{S}$  denotes the selected news article set, and  $\zeta$  denotes the article being selected. After selecting a piece of news  $\zeta$ , our selection strategy must insures that:

- The diversity should not deviate much in  $\mathcal{S}$ . This strategy insure that the recommended list would provide more different news articles from different topics.
- $\mathcal{S}$  should give more satisfaction to the user based on his/her reading history. In other words, the selected news set should be more similar to the user’s reading history.
- $\mathcal{S}$  should be similar to the general topic in  $\mathcal{C} \setminus \mathcal{S}$ . In this way, the selected news set would be more taking into account the global candidate set.

Similar to [26], we define a quality function  $q(\mathcal{S})$  to represent the value of current selected news set  $\mathcal{S}$  for each corresponding strategy above as follows:

$$q(\mathcal{S}) = \frac{1}{\binom{|\mathcal{S}|}{2}} \sum_{\substack{n_1, n_2 \in \mathcal{S} \\ n_1 \neq n_2}} -\text{sim}(n_1, n_2) + \frac{1}{|\mathcal{S}|} \sum_{n_1 \in \mathcal{S}} \text{sim}(u, n_1) + \frac{1}{|\mathcal{C} \setminus \mathcal{S}| \cdot |\mathcal{S}|} \sum_{n_1 \in \mathcal{C} \setminus \mathcal{S}} \sum_{n_2 \in \mathcal{S}} \text{sim}(n_1, n_2) \quad (3)$$

where  $n_1$  and  $n_2$  denote two news items,  $u$  denotes the user profile built by his/her reading history, and  $\text{sim}(\cdot, \cdot)$  function returns the similarity of its two parameters. Equation (3) contains three components corresponding to the selection strategy we list above.  $q(\mathcal{S})$  balances the contribution of different strategies. Suppose  $\zeta$  is the selecting candidate news article, the quality increase can be represented as follows:

$$I(\zeta) = q(\mathcal{S} \cup \zeta) - q(\mathcal{S}) \quad (4)$$

The goal of our selection strategy is to select a list of news articles, which provides the largest possible value within the budget (i.e., the budget can be seen as the size of the recommended list). After implementing the selection strategies, we can obtain a list of recommended news articles for each user.

## 5 Experiments

To evaluate our proposed method, we collect the news data from SINA<sup>1</sup>, where the data set ranges from Aug 1, 2012, to Aug 31, 2012. We preprocess the corpus by removing the unpopular news articles (i.e., the news articles read by <5 users) and the non-active users (i.e., the users who read <5 pieces of news) to verify our proposed approach performance. After preprocessing, 1127 users are stored with 28,737 news articles. In this data set, we randomly select 10 % as testing data in each user and the other 90 % reading history data for training.

### 5.1 Baseline methods

In our work, we mainly compared several state-of-the-art news recommendations as baseline methods: ClickB [2], Goo [3], SCENE [1] and TwoHy [4] for comparison. The core idea of these approaches is described as follows:

- ClickB: This is a content-based recommender, in which each user's profile is constructed by his/her reading history based on the click behavior. A Bayesian framework is utilized for predicting users' preferences.
- Goo: It is a collaborative filtering-based recommender. In this method, researchers combined MinHash clustering, PLSI topic model and co-visitation counts into a personalized news recommendation framework.
- SCENE: This approach is a hybrid method, and it proposes a two-stage recommendation with considering the named entity, user access pattern and news content together. For the first stage, this recommender selects the news topic clusters which the user may be interested.
- TwoHy: Similar to the above approach, this method is also a hybrid recommender. This method clusters the users, named entities, news articles and news topics in a hypergraph model at the first stage. After clustering processing, this method merges newly published news into this model and implements a hypergraph-based ranking and gets the recommended lists for the users.

Try to make sure the fair comparison, the parameters of the above approaches are optimally tuned in our experiments.

<sup>1</sup> <http://news.sina.com.cn>.

### 5.2 Evaluation

To evaluate the performance of our algorithm and other baseline approaches, we use F1-score, Diversity and Normalized Discount Cumulative Gain (NDCG) as our evaluation metrics. Precision is the number of correctly recommended news articles divided by the total number of articles in recommended list. Recall is the number of correctly recommended news articles divided by the total number of articles which should be recommended. F1-score is the harmonic mean of Precision and Recall and shown as follows.

$$F1 = \frac{2PR}{P + R} \quad (5)$$

where  $P$  denotes Precision and  $R$  denotes Recall.

Diversity can be defined as: Given a list of recommended news  $\mathcal{L}$  to a given user  $u$ ,  $\mathcal{L}$  can be said to be diverse to  $u$  if it matches as many as possible of the given user's different reading interests. Let  $R(u)$  be a news recommended list of a user  $u$ , and the diversity of  $u$  can be computed as follows:

$$\text{Diversity} = 1 - \frac{\sum_{i,j \in R(u), i \neq j} \text{sim}(i,j)}{\frac{1}{2}|R(u)|(|R(u)| - 1)}, \quad (6)$$

where  $i$  and  $j$  are two different news articles in recommendation list for user  $u$ , and  $\text{sim}(i, j)$  denotes the news profile similarity between the news item  $i$  and  $j$ . We use NDCG metric to measure the quality of ranking. NDCG value at position  $n$  can be defined as follows:

$$\text{NDCG}@n = \frac{1}{\text{IDCG}} \times \sum_{i=1}^n \frac{2^{r_i} - 1}{\log_2(i + 1)}, \quad (7)$$

where  $r_i$  is the relevance rating of item at rank  $i$ . In news recommendation case,  $r_i$  is set to be 1 if the user has read this recommended news article and 0 otherwise. IDCG is chosen to ensure that the perfect ranking has a NDCG value of 1.

### 5.3 Parameter tuning

In our method, there is only one parameter for tuning,  $\varepsilon$  in Eq. (2). We utilize F1-score value for tuning this parameter, and Fig. 4 shows the results.

From the result, we set the parameters  $\varepsilon$  as 0.2 for our experiments. If the  $\varepsilon$  is set as too small value, the relationship between the key word and the news document would become the direct relation. In this way, the news story chain would loss many useful news article points. Conversely, if  $\varepsilon$  is too large, the relationship between key word and news document would become obscure. Using large  $\varepsilon$  might obtain many non-related news articles in a news story chain.

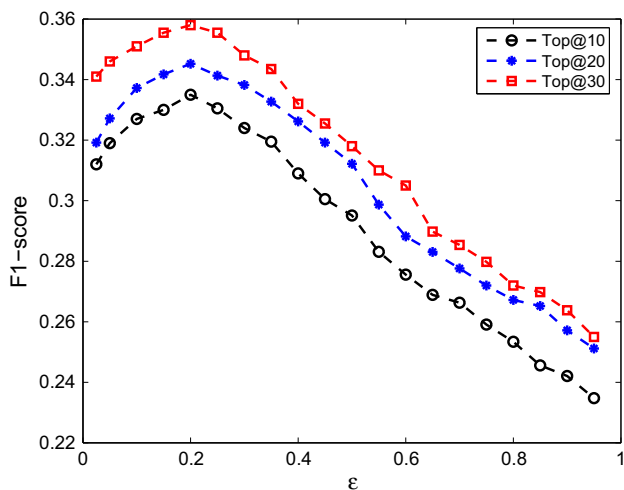


Fig. 4 Parameter tuning

5.4 Accuracy evaluation

In this evaluation, we utilize NDCG and F1-score for comparison. The experimental results are shown in Fig. 5. It is can be seen that our approach significantly outperforms the other methods on NDCG and F1-score metrics. The results of SCENE and TwoHy are comparable to ours. In the study of SCENE, they use named entity and news topic as recommended factors. Therefore, SCENE considers more topic relation and named entity and consequently obtains reasonable result. TwoHy explicitly takes into account the relationship between named entity, news topic, news article and user, and encapsulates them in hypergraph ranking. As the observation, the performance of ClickB and Goo is relatively poor. The straightforward reason may be that the recommended results of this two approaches are heavily related with users’ co-visiting histories; however, the corpus set used in our experiment contains a lot of cold users, i.e., users who read <10 news articles.

Discussion NDCG and F1-score are accuracy evaluation metrics. Normally, the NDCG would be good when the F1-score is good. We proposed method considered the local

reading profile of the given user. In this way, the user profile building process is starting in the news chain building. Compared to the user profile built from all training corpus of the baselines, our proposed method could perform a better accuracy.

5.5 Handing cold-start problem

In traditional recommenders, it needs lots of history news content or behaviors to construct the user’s preference profile. Therefore, the cold-start problem cannot be well solved due to the data sparsity. Our proposed method can resolve it. The straightforward explanation is that: Our news story chain building method not only considers the items are preferred by this user or not, but also takes the temporal factor into account.

In this section, we randomly select 100 cold users who read <10 (but more than 5) news articles for experiment. Figure 6a, b shows the results. As shown in experimental results, the cold user start problem can be significantly alleviated via our approach.

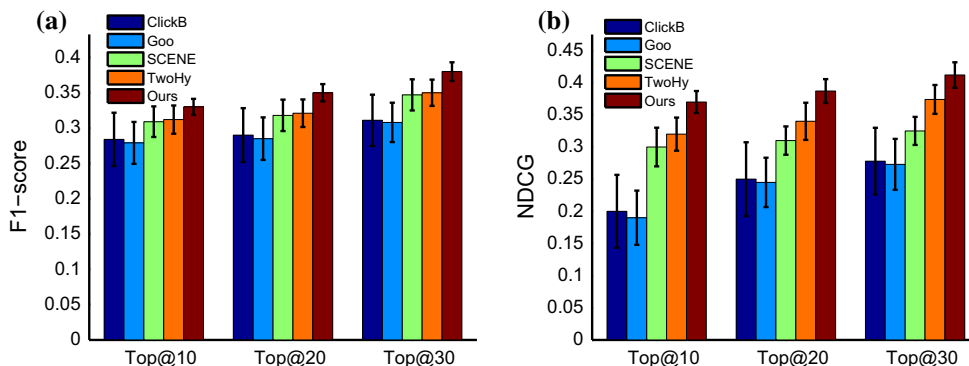
From the result of Figs. 5 and 6, we can observe that the performances of the baseline approaches SCENE and TwoHy are similar to our proposed method. However, in the cold users data set, our proposed method is significantly outperformed the other baselines. The straightforward reason may be that our proposed method selects the news articles along with the users’ reading time series. Therefore, our method can mine the accurate interest of the cold users with considering the given reading history news articles pair.

5.6 Stability evaluation

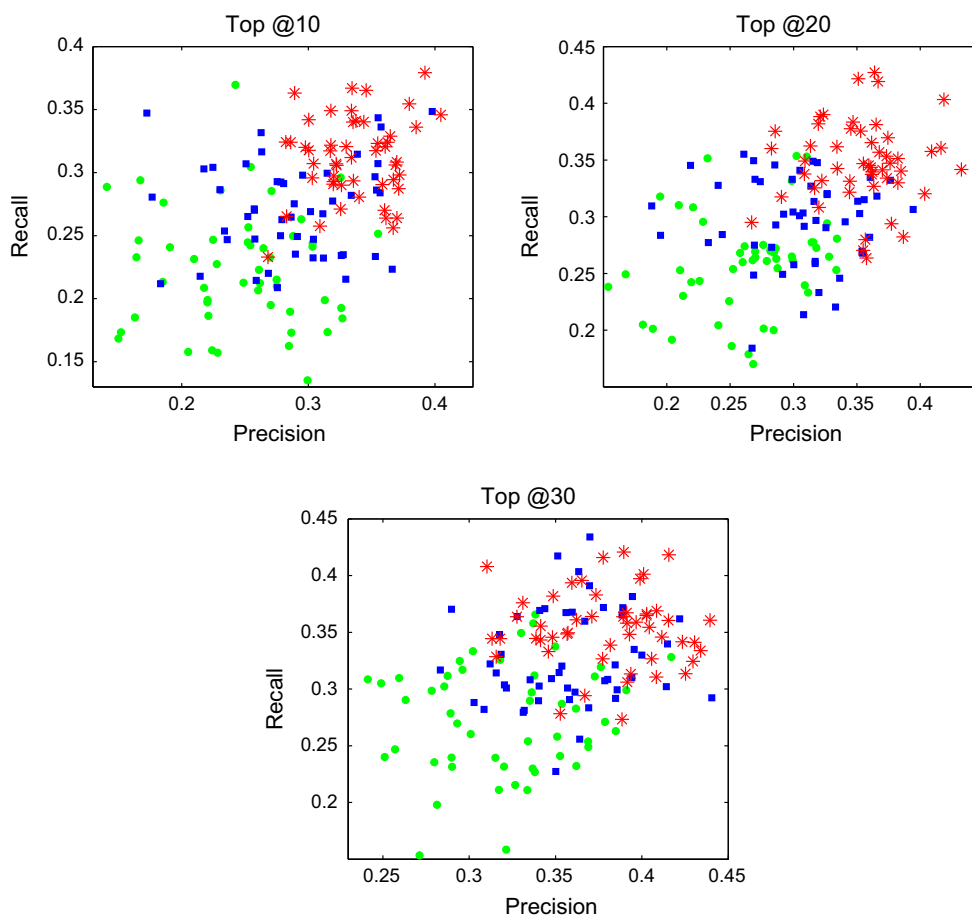
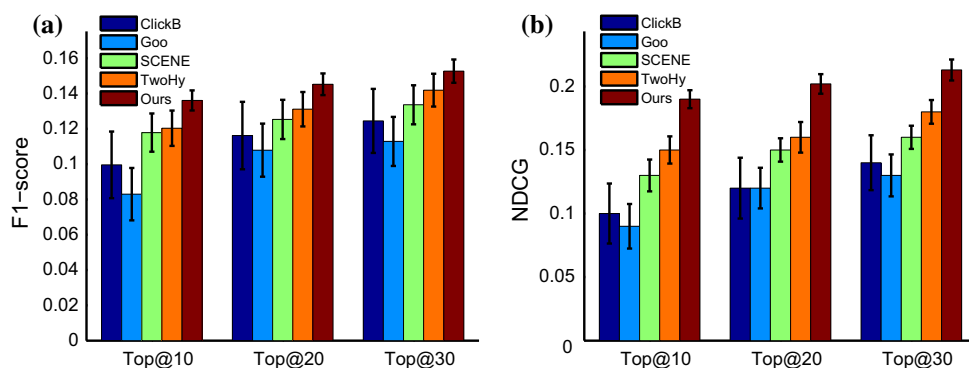
For the stability experiment, we randomly choose 50 users to provide news articles for them in three recommended approaches (ClickB, SCENE and ours). Figure 7 shows the comparison results as Top@10, Top@20 and Top@30 lists for each user.

In this experiment from sample users, we can observe that, besides the higher accuracy, our proposed method

Fig. 5 Comparison on F1-score of different approaches. a F1-score. b NDCG



**Fig. 6** Comparison on F1-score of different approaches on cold users. **a** F1-score. **b** NDCG



**Fig. 7** Precision–recall plot of different recommenders. Remark circle (green) denotes the ClickB; square (blue) denotes the SCENE; and asterisk (red) represents our method

outperforms the other baselines in stability due to the density of distribution. The straightforward reasons may be that: (1) the traditional approaches recommended news articles using the all user’s reading history. For example, the ClickB and SCENE methods constructed user profile based on the user’s reading history by VSM model. Therefore, the user profile will be equivalent to a multi-dimensional center vector in Euclidean space. In this way, it would lead to recommend the same result in spite of the user read lots of different news

topics. (2) Here may have the application reason. We observe that there are many cold users in our data set. Cold users may lead to more unstable factors in traditional approaches.

### 5.7 Diversity evaluation

A user would not prefer to read the recommended list which contains many similar news articles. Therefore, a recommendation will be more popular if it performs more



**Table 1** Diversity evaluation on different sizes of recommendation lists. Remarks: the significance of bold shows our approach outperforms the other baselines

Methods	Top@10	Top@20	Top@30
ClickB	0.4231	0.3128	0.0765
Goo	0.4598	0.3210	0.0841
SCENE	0.6014	0.4517	0.3104
TwoHy	0.5984	0.4124	0.3487
Ours	<b>0.6514</b>	<b>0.6012</b>	<b>0.5278</b>

diversity. For this metric evaluation, we choose ClickB [2] (a content-based method), Goo [3] (a collaborative filtering-based method), SCENE [1] (a hybrid method using LSH for clustering and greedy algorithm for news selection) and TwoHy [4] (a hypergraph ranking-based news recommendation) as the baselines. Table 1 shows the comparison result of the *Top@10*, *Top@20* and *Top@30*.

From the comparison result above, we can observe that our proposed method outperforms the other approaches significantly. Besides, our method slightly descends with the size of recommendation list increases. The straightforward reason may be that we diverse the news articles by two times, in which the first time is that we select the candidate set from each day and construct the different chains from different endpoints, and the second time is the greedy precise selection strategy we used. We can also observe that the diversity decreases as the recommended list enlarges. It is straightforward that when more news documents are selected, the topic distribution of the list becomes closer to the user's long-term preference. Therefore, the diversity of methods ClickB and Goo drops dramatically compared to the other approaches. The reason may be that they fail to consider the diversity factor into the recommendation model.

## 6 Conclusion

In this work, we propose the news personalized recommendation via a news story chain building and focus on the application problem by providing the news article for helping user gain a better understanding of the development of the whole news story chain. We model this recommendation problem as a news story chain building problem over time. To the best of our knowledge, the idea of the proposed method is first used in news recommendation. Some studies or applications constructed the news chain based on the word bag and vector space model. Different from the previous work, we proposed to use the random walk for linking relationship between key word and the news document. In this way, the user's reading preference and context can be naturally captured. Besides,

we use a greedy selection method for providing more stability, accuracy and diversity recommendation list in this framework. Based on this framework, we treat the news personalized recommendation as a *critical narrative missing link provider*. The experimental results based on a real-world data set have showed that our proposed approach significantly outperforms the traditional baseline approaches.

For future work, to process massive news articles and users, we plan to deploy our recommendation onto the Map-Reduce framework on our distributed system. Moreover, we also plan to integrate this news story chain building approach into our news search engine due to the effectiveness in this study. Another remarkable point is the interest effectiveness of user's reading chain on news recommender, which is able to insights on the exploration of personalized news reading behaviors. If the news story chain or news story network can be built perfectly, the news recommended research or other applications about network news could be studied in this background which connecting to the real-world news development.

## References

- Li L, Wang DD, Li T, Knox D, Padmanabhan B (2011) Scene: a scalable two-stage personalized news recommendation system. In: ACM conference on information retrieval (SIGIR)
- Liu J, Dolan P, Pedersen ER (2010) Personalized news recommendation based on click behavior. In: Proceedings of the 15th international conference on Intelligent user interfaces. ACM, pp 31–40
- Das AS, Datar M, Garg A, Rajaram S (2007) Google news personalization: scalable online collaborative filtering. In: Proceedings of the 16th international conference on World Wide Web. ACM, pp 271–280
- Li L, Li T (2013) News recommendation via hypergraph learning: encapsulation of user behavior and news content. In: Proceedings of the sixth ACM international conference on Web search and data mining. ACM, pp 305–314
- Schafer JB, Konstan J, Riedi J (1999) Recommender systems in e-commerce. In: Proceedings of the 1st ACM conference on electronic commerce. ACM, pp 158–166
- Billsus D, Pazzani MJ (1999) A personal news agent that talks, learns and explains. In: Proceedings of the third annual conference on Autonomous Agents. ACM, pp 268–275
- Pazzani M, Billsus D (1997) Learning and revising user profiles: the identification of interesting web sites. Mach Learn 27(3):313–331
- Cota RG, Ferreira AA, Nascimento C, Gonçalves MA, Laender AHF (2010) An unsupervised heuristic-based hierarchical method for name disambiguation in bibliographic citations. J Am Soc Inf Sci Technol 61(9):1853–1870
- Breese JS, Heckerman D, Kadie C (1998) Empirical analysis of predictive algorithms for collaborative filtering. In: Proceedings of the fourteenth conference on Uncertainty in artificial intelligence. Morgan Kaufmann Publishers Inc., pp 43–52
- Claypool M, Gokhale A, Miranda T, Murnikov P, Netes D, Sartin M (1999) Combining content-based and collaborative filters in an

- online newspaper. In: Proceedings of ACM SIGIR workshop on recommender systems, vol 60, Citeseer
11. Chu W, Park S-T (2009) Personalized recommendation on dynamic content using predictive bilinear models. In: Proceedings of the 18th international conference on World wide web. ACM, pp 691–700
  12. Li L, Chu W, Langford J, Schapire RE (2010) A contextual-bandit approach to personalized news article recommendation. In: Proceedings of the 19th international conference on World wide web. ACM, pp 661–670
  13. Best C, van der Goot E, de Paola M, Garcia T, Horby D (2002) Europe media monitor-emm. JRC Technical Note No I, 2
  14. Gabrilovich E, Dumais S, Horvitz E (2004) Newsjunkie: providing personalized newsfeeds via analysis of information novelty. In: Proceedings of the 13th international conference on World Wide Web. ACM, pp 482–490
  15. Turner SR (2014) The creative process: a computer model of storytelling and creativity. Psychology Press, London
  16. Niehaus J, Michael YR (2009) A computational model of inferring in narrative. In: AAAI Spring Symposium: Intelligent Narrative Technologies II, pp 83–90
  17. Masand B, Linoff G, Waltz D (1992) Classifying news stories using memory based reasoning. In: Proceedings of the 15th annual international ACM SIGIR conference on Research and development in information retrieval. ACM, pp 59–65
  18. Kleinberg J (2003) Bursty and hierarchical structure in streams. *Data Mining Knowl Discov* 7(4):373–397
  19. Yang Y, Carbonell JG, Brown RD, Pierce T, Archibald BT, Liu X (1999) Learning approaches for detecting and tracking news events. *IEEE Intell Syst* 14(4):32–43
  20. Rowe JP, McQuiggan SW, Robison JL, Marcey DR, Lester JC (2009) Storyeval: an empirical evaluation framework for narrative generation. In: AAAI Spring Symposium: Intelligent Narrative Technologies II, pp 103–110
  21. Shahaf D, Guestrin C (2010) Connecting the dots between news articles. In: Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, pp 623–632
  22. Kempe D, Kleinberg J, Tardos É (2003) Maximizing the spread of influence through a social network. In: Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, pp 137–146
  23. Kleinberg JM (1999) Authoritative sources in a hyperlinked environment. *J ACM* 46(5):604–632
  24. Brin S, Page L (1998) The anatomy of a large-scale hypertextual web search engine. *Comput Netw ISDN Syst* 30(1):107–117
  25. Cunningham H, Maynard D, Bontcheva K, Tablan V (2002) Gate: an architecture for development of robust HLT applications. In: Proceedings of the 40th annual meeting on association for computational linguistics. Association for Computational Linguistics, pp 168–175
  26. Gu W, Dong S, Zeng Z, He J (2014) An effective news recommendation method for microblog user. *Sci World J* 2014:907515