

Portland State University

PDXScholar

Electrical and Computer Engineering Faculty
Publications and Presentations

Electrical and Computer Engineering

5-30-2020

DIM: Adaptively Combining User Interests Mined at Different Stages Based on Deformable Interest Model

Xiaoru Wang

Beijing University of Posts and Telecommunications

Yueli Li

Beijing University of Posts and Telecommunications

Zhihong Yu

Intel China Research

Fu Li

Portland State University, lif@pdx.edu

Heng Zhang

Beijing University of Posts and Telecommunications

Follow this and additional works at: https://pdxscholar.library.pdx.edu/ece_fac



See next page for additional authors
Part of the [Electrical and Computer Engineering Commons](#)

Let us know how access to this document benefits you.

Citation Details

Wang, X., Li, Y., Yu, Z., Li, F., Zhang, H., Cai, Y., & Li, L. (2020). DIM: Adaptively Combining User Interests Mined at Different Stages Based on Deformable Interest Model. *Mathematical Problems in Engineering*, 2020.

This Article is brought to you for free and open access. It has been accepted for inclusion in Electrical and Computer Engineering Faculty Publications and Presentations by an authorized administrator of PDXScholar. Please contact us if we can make this document more accessible: pdxscholar@pdx.edu.

Authors

Xiaoru Wang, Yueli Li, Zhihong Yu, Fu Li, Heng Zhang, Yali Cai, and Lixian Li

Research Article

DIM: Adaptively Combining User Interests Mined at Different Stages Based on Deformable Interest Model

Xiaoru Wang ¹, Yueli Li ¹, Zhihong Yu ², Fu Li ³, Heng Zhang ¹, Yali Cai ¹,
and Lixian Li ¹

¹Beijing University of Posts and Telecommunications, Beijing, China

²Intel China Research Center, Beijing, China

³Department of Electrical and Computer Engineering, Portland State University, Portland, OR 97207-0751, USA

Correspondence should be addressed to Xiaoru Wang; wxr@bupt.edu.cn

Received 13 January 2020; Revised 12 April 2020; Accepted 11 May 2020; Published 30 May 2020

Academic Editor: Ioannis Kostavelis

Copyright © 2020 Xiaoru Wang et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

User interest mining is widely used in the fields of personalized search and personalized recommendation. Traditional methods ignore the formation of user interest which is a process that evolves over time. This leads to the inability to accurately describe the distribution of user interest. In this paper, we propose the interest tracking model (ITM). To add the timing, ITM uses Dirichlet distribution and multinomial distribution to describe the evolutionary process of interest topics and frequent patterns, which well adapts to the evolution of user interest hidden in short texts between different time slices. In addition, it is well known that user interest is composed of long-term interest and situational interest including short-term interest and social hot topics. State-of-the-art methods simply regard the users' long-term interest as the users' final interest, which makes those unable to completely describe the user interest distribution. To solve this problem, we propose the deformable interest model (DIM) which designs an objective function to combine users' long-term interest and situational interest and more comprehensively and accurately mine user interest. Furthermore, we present the degree of deformation which measures the subinterest's degree of influence on final interest and propose in DIM the influence real-time update mechanism. The mechanism adaptively updates the degree of deformation through the linear iteration and reduces the degree of dependence of the interest model on training sets. We present results via a dataset consisting of Flickr users and their uploaded information in three months, a dataset consisting of Twitter users and their tweets in three months, and a dataset consisting of Instagram users and their uploaded information in three months, showing that the perplexity is reduced to 0.378, the average accuracy is increased to 94%, and the average NMI is increased to 0.20, which prove better interest prediction.

1. Introduction

User interest mining refers to establishing a user interest model by analyzing a large amount of user behavior data. Through user models of high quality, it is able to describe the real interest of users, making it possible to implement personalized services for users. In recent years, user interest mining has been widely used in the fields of personalized search and personalized recommendation.

Describing the distribution of user interest is the core of user interest mining. Psychologists believe that the formation of user interest is a process that evolves over time [1]. Therefore, tracking and describing the evolution of user

interest is the biggest challenge in describing the distribution of user interest. In previous research works, static topic models were often used to describe the distribution of user interest, such as the latent Dirichlet allocation (LDA) model proposed by Blei et al. in 2003 [2]. However, in real life, people's subjective intentions often change with time, and user interest will continue to evolve with time. Static topic models are difficult to meet this demand. Therefore, some literature studies [3–10] attempt to introduce time dimension to track the dynamic changes of user interest. The distribution of user interest described by the current dynamic topic models is a Gaussian distribution centered on the superparameter α of the interest distribution of the

aforementioned time slice and cannot adapt to user interest that suddenly changes between different time slices [4].

In addition, psychologists divide user interest into long-term interest and situational interest [11–13]. Long-term interest refers to a relatively stable and persistent individual tendency that develops over time. Situational interest is considered as a relatively passive and transient emotional state triggered by certain conditions in the environment, including short-term interest of users and social hot topics. The traditional methods of user interest mining simply regarded the user's long-term interest as the user's final interest, which makes the traditional methods unable to completely describe the distribution of user interest.

On this basis, it should also be considered that the influence of long-term interest, short-term interest, and hot topics on user interest is updated in real time. This variability requires that the new algorithm must have an adaptive mechanism. The stable influence of three interests not only ignores the possibility that the user generates a new interest but also the effect of interest mining which is strongly dependent on training sets. Once the training set is changed, it is necessary to retrain the model in order to update the influence of three interests, which is undoubtedly very time consuming.

In order to solve the aforementioned problems, this paper proposes a user mining method based on the deformable interest model to adaptively integrate users' long-term interest and situational interest. The contributions of this paper mainly include the following aspects:

- (i) For tracking and describing the evolution of user interest, this paper introduces the time dimension and proposes the interest tracking model (ITM), which maps annotated words to the frequent pattern space and uses the Dirichlet distribution and the multinomial distribution to describe the evolutionary process of user interest and frequent patterns between different time slices, respectively.
- (ii) For solving the problem of the integrity of user interest, this paper proposes the deformable interest model (DIM), which integrates users' long-term interest and situational interest and more comprehensively and accurately mines user interests.
- (iii) For solving the problem that the influence of long-term interest and situational interest on user interest needs to be updated in real time, this paper proposes the deformable interest model (DIM), which uses the real-time update mechanism to adaptively update the influence of long-term interest, short-term interest, and hot topics on user interest. The real-time update mechanism not only considers the possibility of interest change but also reduces the dependence of the interest model on training sets.

2. Related Work

Building a topic model is the primary means of mining user interest. The topic model is a language model that uses the Bayesian statistics and machine learning methods to discover the underlying semantic content of nonlabeled

documents and uses these latent semantics to predict the future characteristics of the document set.

The topic models used for interest mining in the early days were static topic models of which the establishment was independent of time. David et al. proposed a probabilistic topic model called LDA (latent Dirichlet allocation) in 2003 [2]. Because of its good mathematical foundation and flexible scalability, LDA has been widely paid attention and used in various research fields since its introduction. However, due to the semantic gap, the application of LDA on sparse short texts makes it difficult to confirm the semantic consistency of words. As a result, some methods aggregated short texts into long texts to reduce inaccuracy. Other methods enriched original data through external knowledge basis. These methods are not always effective because there may be semantic inconsistencies between the pseudo-long-text after aggregation and the original short text. Another recent work finds the embedding topic model [14] that combines traditional topic models and word embedding which can well analyze the semantic connotation of large text sets with many long-tail words and low-frequency words. However, using word embedding to represent documents makes the features too low-level, and it is often difficult to obtain satisfactory results. In addition, LDA treats a word as a unit, which undoubtedly reduces semantic accuracy. Wallach [15] proposed the bigram LDA model, which treats bigram as a unit. On this basis, Wang et al. [16] proposed topical N-grams (TNG). Jhnichen et al. [17] proposed scalable generalized dynamic topic models which used stochastic processes to introduce stronger correlations. These methods can break the limitations of the models with bag-of-words to a certain extent and find common phrases and potential topics in the text, but the model is complex. In addition, association rule mining technology is also an effective data mining technology, such as association rule mining applied to streaming data [18, 19] and association rule mining applied to dynamic databases [20, 21]. Compared with other association rule mining techniques, the frequent pattern mining model has the simplest structure. This inspired us to introduce frequent pattern mining into static topic modeling so that topic modeling based on bag-of-words can be transformed into topic modeling based on pattern sets. Nevertheless, user interest is dynamically formed over time, and static topic models cannot satisfy this need.

In order to solve the problem that the static topic model cannot change the subject content over time, topic models with dynamics have also been widely studied. These include DTM [3], cDTM [6], TTM [4], and D-ETM [22]. The methods solved the problem that the static topic model cannot respond to the change of the user's related information in time, but they all work in the context of long texts. Moreover, in DTM [3] and cDTM [6], the distribution of user interest comes from the Gaussian distribution centered on the superparameter α of the interest distribution in the last time slice and cannot adapt to the sudden change of user interest. At the same time, the Gaussian distribution and the multinomial distribution are not conjugated, and the model is not interpretative and practical. Some recent works

[23, 24] found embedding representations that vary over time. Therefore, D-ETM came into being on the basis of ETM. However, how to process higher-level features on the basis of the dense vector space is the key to improving the effect of the algorithm. In addition, Liang et al. proposed a dynamic topic model called UCIT-L for short texts [25]. UCIT-L infers the user's interest based on the user's information in multiple time periods and the information of its followers. The disadvantage is that the amount of calculation is too large, and the degree of coincidence between the user's interest and the followers' interest is uncertain.

In order to alleviate the key problem of nonconjugation caused by potential Gaussian variables and their subsequent nonlinear transformation of count values of the model, Linderman et al. proposed the Polya-gamma augmentation in 2015 [5]. This approach helped to alleviate the problems with DTM [3] and cDTM [6], but it did not necessarily improve the performance. The interest tracking model (ITM) proposed in this paper can essentially solve the nonconjugation problem caused by Gaussian distribution and multidistribution and accurately track and describe the user interest that evolves over time.

The above interest mining methods were devoted to describing the dynamic process of user interest and simply regarded the user's long-term interest as the user's final interest, ignoring the situational interests consisting of short-term interest of the user and social hot topics, which cannot completely describe the distribution of user interest [1, 11–13]. The variable interest model proposed in this paper is a model that fully considers the user's real interest, adaptively integrates long-term interest, short-term interest, and current hot topics, and depicts the real interest formation process.

3. Mining User Interest Based on the Deformable Interest Model

In this paper, we consider that user interest is composed of long-term interest and situational interest including short-term interest and the current social hot topics. Therefore, we propose a method of mining user interest based on the deformable interest model (DIM) which is aimed at fusing the above three interests. Long-term interest is mined by the interest tracking model (ITM), short-term interest is mined by LDA-FP [2, 26], and the current social hot topics are obtained by the knowledge base [27, 28]. Using DIM, three interests are adaptively combined to obtain user interest.

3.1. Problem Definition. A social network has a set of users $U = \{u_1, u_2, \dots, u_a, \dots, u_n\}, n \in Z$; the user u_a uploads a set of pictures $I = \{i_{a1}, i_{a2}, \dots, i_{am}\}, m \in Z$, where u_a means the a -th user. At the same time, when uploading a picture, the user adds a set of annotated words $W = \{w_{i1}, w_{i2}, \dots, w_{ia}, \dots, w_{ip}\}, p \in Z$, to the i -th picture according to their own interests, where w_{ia} means the a -th annotated word of the i -th picture. All annotated words for different pictures marked by user u_a in time slice T are denoted as $d_a^T = \{w_1, w_2, \dots, w_{(m \times p)}\}$. It is known that user interest

changes as time elapsed, and user interest is easily changed by the situational interest including the user's short-term interest and the current social hot topics.

The tasks of DIM are as follows: (1) give the distribution of user's short-term interest in the last time slice, the distribution of user's long-term interest that changes over time, and the distribution of the current hot topics; (2) combine the user's long-term interest and situational interest including short-term interest and current social hot topics; and (3) adaptively update the degree of deformation which measures the subinterest's degree of influence on final interest.

First, we process different categories of annotation words into a corpus; then, we preprocess the corpus; and finally, we mine frequent patterns from the corpus by using FP-growth algorithm [26] and establish a frequent pattern library. This pattern library is defined as $C = \{c_1, c_2, \dots, c_n\}, n \in Z, c_i = \{w_1, w_2, \dots, w_{np_i}\}$, where $w_t (t = 1, 2, \dots, np_i)$ is the t -th word and np_i is the number of words in the frequent pattern i . LDA [2] represents the text by mapping the text to bag-of-words, but this method does not apply to short texts with sparsity problems, so we represent the user's annotated words by mapping words to the frequent pattern library. We denote the user's annotated words as $d_a^T = \{c_i | \forall w_t \in c_i, \forall w_t \in d_a^T\}$. We consider $d_a = (d_a^1, d_a^2, \dots, d_a^T)$ as the input of ITM to obtain long-term interest of users $L = (a_{L1}w_{L1}, a_{L2}w_{L2}, \dots, a_{Ln}w_{Ln}), n \in Z$. At the same time, we consider d_a^T as the input of LDA-FP to obtain short-term interest of users $S = (a_{S1}w_{S1}, a_{S2}w_{S2}, \dots, a_{Sn}w_{Sn}), n \in Z$, and obtain the current hot topics based on the knowledge base $H = (a_{H1}w_{H1}, a_{H2}w_{H2}, \dots, a_{Hn}w_{Hn}), n \in Z$. Each element in L, S , and H is represented by a weight coefficient $(a_{i1}, i \in S, L, H)$ and the corresponding labeled-word $(w_{i1}, i \in S, L, H)$. Eventually, we consider (L, S, H) as the input of DIM to adaptively combine long-term interest, short-term interest, and current hot topics.

3.2. Mining User Long-Term Interest Based on the Interest Tracking Model. User interest is known to change over time, but this change does not happen suddenly; it has some kinds of continuity between time periods. ITM defines the "user-interest" vector with a first-order Markov property and considers that user's interest distribution in the current time slice Θ_t is basically around that in the last time slice Θ_{t-1} . Therefore, we define the "user-interest" distribution Θ_t of the current time slice as

$$P(\Theta_t | \hat{\Theta}_{t-1}, \alpha_t) \propto \prod_z \theta_{t,z}^{\alpha_t \hat{\Theta}_{t-1,z}^{-1}}, \quad (1)$$

where Z is the number of interests, α_t is the hyperparameter of the "user-interest" distribution in the current time slice, and $\theta_{t,z} = P(z | t)$ represents interests of a user, which is the probability that the user is interested in interest z at time t , where $\theta_{t,z} \geq 0, \sum_Z \theta_{t,z} = 1$.

Correspondingly, the latent semantics of each interest in the current time slice will also change. We confirm the latent semantics of interest by finding the "interest-frequent

pattern” distribution. Similarly, we define the “interest-frequent pattern” distribution $\Phi_{t,z} = \{\phi_{t,z,i}\}_i^C$ in the current time slice as

$$P(\Phi_{t,z} | \hat{\Phi}_{t-1,z}, \beta_{t,z}) \propto \prod_i \phi_{t,z,i}^{\beta_{t,z} \hat{\phi}_{t-1,z,i} - 1}, \quad (2)$$

where $\beta_{t,z}$ is the hyperparameter of the “interest-frequent pattern” distribution in the current time slice and $\phi_{t,z,i} = P(i|z, t)$ represents trends in an interest, which is the probability that the frequent pattern p_i is selected from interest z at time t , where $\phi_{(t,z,i)} \geq 0$, $\sum_i \phi_{t,z,i} = 1$.

Based on (1) and (2), the process of generating each frequent pattern in the model is described as follows: for each frequent pattern, the “user-interest” distribution Θ_t in the current time slice is determined by the prior knowledge α_t in the current time slice and the user-interest distribution Θ_{t-1} in the previous time slice together. Next, extract an interest from Θ_t , and then deduce the “interest-frequent pattern” distribution $\Phi_{t,z}$ of the current time slice according to the prior knowledge $\beta_{t,z}$ in the current time slice and the “interest-frequent pattern” distribution $\Phi_{t-1,z}$ in the previous time slice. Finally, extract a frequent pattern from $\Phi_{t,z}$ corresponding to the interest Θ_t .

It can be seen from the above that the ITM maps annotated words to a frequent pattern set to represent a collection of annotated words for each user that is considered short text. Each frequent pattern contains a set of annotated

words that occur frequently at the same time. In addition, the ITM uses Θ_{t-1} and $\Phi_{t-1,z}$ of the previous time slice to correct the Dirichlet parameter in the current time slice to achieve the purpose of tracking the evolution of user interest. At the same time, the ITM maintains the conjugate distribution of Dirichlet-multinomial. This design reflects the mathematical nature of user interest’s evolution and makes the model interpretable. The probability model diagram for this model is shown in Figure 1.

How to reverse user interest based on the known frequent pattern of each user’s labeled words is the purpose of constructing the interest tracking model. In this work, we estimate parameters in the ITM based on a random EM algorithm [29], where the Gibbs sampling of latent topics and the maximum joint likelihood estimation of parameters are alternately iterated.

The ultimate goal of building ITM is to get the posterior probability:

$$P(Z_t | W_t, \alpha_t, \beta_t) = \frac{P(Z_t, W_t | \alpha_t, \beta_t)}{P(W_t | \alpha_t, \beta_t)}. \quad (3)$$

Thus, the problem of solving the posterior probability $P(Z_t | W_t, \alpha_t, \beta_t)$ is transformed into the problem of solving the joint distribution of user interest and patterns. From the definition of Dirichlet distribution and multinomial distribution [2], we make the following inference:

$$\begin{aligned} P(Z_t, W_t | \hat{\Theta}_{t-1}, \hat{\Phi}_{t-1}, \alpha_t, \beta_t) &= P(Z_t | \hat{\Theta}_{t-1}, \alpha_t) P(W_t | \hat{\Phi}_{t-1}, Z_t, \beta_t) \\ &= \int P(Z_t, \Theta_t | \hat{\Theta}_{t-1}, \alpha_t) d\Theta_t \times \int P(W_t, \Phi_t | \hat{\Phi}_{t-1}, Z_t, \beta_t) d\Phi_t \\ &= \frac{\Gamma(\alpha_t)}{\prod_z \Gamma(\alpha_t \hat{\theta}_{t-1,z})} \frac{\prod_z \Gamma(n_{t,z} + \alpha_t \hat{\theta}_{t-1,z})}{\Gamma(n_t + \alpha_t)} \times \prod_z \frac{\Gamma(\beta_{t,z})}{\prod_i \Gamma(\beta_{t,z} \hat{\phi}_{t-1,z,i})} \frac{\prod_i \Gamma(n_{t,z,i} + \beta_{t,z} \hat{\phi}_{t-1,z,i})}{\Gamma(n_{t,z} + \beta_{t,z})}, \end{aligned} \quad (4)$$

where $n_{t,z}$ is the number of patterns that have been assigned to interest z at time t , $n_{t,z,i}$ is the number of times for which pattern p_i has been assigned to interest z at time t , and $\Gamma(x)$ is the gamma function.

Solving the parameters in this joint distribution by Gibbs sampling and maximum likelihood estimation [29] yields the following results:

$$\hat{\theta}_{t,z} = \frac{n_{t,z} + \alpha_t \hat{\theta}_{t-1,z}}{n_t + \alpha_t}, \quad (5)$$

$$\hat{\phi}_{t,z,i} = \frac{n_{t,z,i} + \beta_{t,z} \hat{\phi}_{t-1,z,i}}{n_{t,z} + \beta_{t,z}}. \quad (6)$$

When inferring the current interests Θ_t and trends Φ_t , ITM only uses the current data. Therefore, compared with

the traditional model, it not only describes the evolution process of the user’s long-term interest without increasing the potential variables but also reduces the calculation amount and improves the calculation speed.

3.3. Adaptively Fusing User Subinterests Based on the Deformable Interest Model. User interest is not instantaneously formed by long-term interest, short-term interest, and current hot topics but gradually changes to the final state. To describe this process of change, we consider user interest as a deformable interest model consisting of three variables of sudden interest (hot topics), short-term interest, and long-term interest. In the deformable interest model (DIM), the degree of deformation of each interest controls the deformation of the entire system. The degree of deformation of each interest is determined by the interaction of the

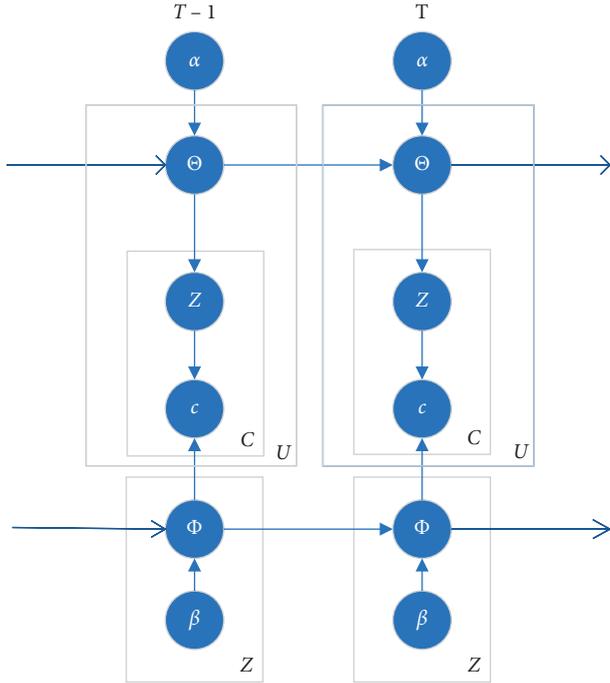


FIGURE 1: The structure of the interest tracking model.

respective interest and its similarity to user interest, thus ensuring that important interest has a significant impact. The essence of the deformable interest model is the spring model. The structure diagram of the spring model is shown in Figure 2 [30].

The deformable interest model can be defined as a quaternion (L, S, H, b) , where L represents the interest tracking model, S represents LDA-FP [2, 26], H represents the knowledge base of the current hot topics [27, 28], and b represents the bias value. Each submodule is represented by a multigroup $(a_{i1}c_{i1}, a_{i2}c_{i2}, a_{in}c_{in})$, $i = \{L, S, H\}$, where the number of elements is determined by the number of interest set by the submodel, c_{it} ($t = 1, 2, \dots, n$) is the t -th pattern of subinterest i , and a_{it} is the corresponding weight coefficient.

The score of the target hypothesis is equal to the sum of the similarity between long-term interest and real interest, short-term interest and real interest, and hot topics and real interest minus the difference between short-term interest and long-term interest and hot topics and long-term interest:

$$\begin{aligned} \text{score}(L, S, H, b) &= \text{similar}(R, L) \\ &+ \sum_{i=S}^H [\text{similar}(R, i) - \text{similar}(L, i)] + b, \end{aligned} \quad (7)$$

where i belongs to $\{S, H\}$, and R is the user interest. R is also represented by a multigroup $(a_{R1}w_{R1}, a_{R2}w_{R2}, \dots, a_{Rn}w_{Rn})$. The number of elements is determined by the total number of patterns. The weight coefficient is determined by the frequency of occurrence of the pattern, which means $a_{Rj} = tf_{c_j} \times idf_{c_j}$, $j = 1, 2, \dots, n$, $tf_{c_j} = (n_j/n_{\text{tag}})$,

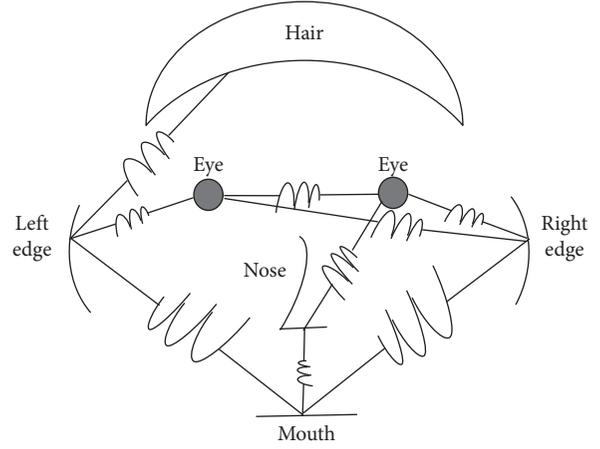


FIGURE 2: The structure of the spring model.

$idf_{c_j} = \lg(n_{\text{user}}/n_{\text{user}}^j + 1)$, where n_j is the number of occurrences of the word with the fewest occurrences in the pattern c_j , n_{tag} is the total number of words for the user, n_{user} is the number of users, and n_{user}^j is the number of users including the pattern c_j . Inspired by cosine similarity [31], the similarity calculation is shown in Algorithm 1.

The degree of deformation of the deformable interest model is

$$\begin{aligned} K &= (k_1, k_2, k_3) \\ &= (\text{similar}(L, R), \text{similar}(R, S) \\ &\quad - \text{similar}(L, S), \text{similar}(R, H) - \text{similar}(L, H)). \end{aligned} \quad (8)$$

Use the linear iterative method to maximize score (L, S, H, b) , and the distribution of user interest is obtained:

$$\begin{aligned} \text{result}(L, S, H) &= K \cdot I^T \\ &= (\text{similar}(L, R), \text{similar}(R, S) \\ &\quad - \text{similar}(L, S), \text{similar}(R, H) \\ &\quad - \text{similar}(L, H)) \cdot (L, S, H)^T. \end{aligned} \quad (9)$$

User interest gained by using DIM is not only more comprehensive but also more interpretable in describing the interactions of deformable interests.

4. Experiments

To verify the validity of the user interest mining algorithm based on the deformable interest model, this section will show the following three experiments:

- (1) The experiment of evaluating performance for the algorithm of users' long-term interest mining based on the interest tracking model
- (2) The experiment of evaluating effectiveness for the algorithm of combining user long-term interest and situational interest based on the deformable interest model
- (3) The experiment of evaluating effectiveness for the algorithm of adaptively adjusting subinterests' influence based on the deformable interest model

Data:
 $A = (\alpha_1 c_{A1}, \alpha_2 c_{A2}, \dots, \alpha_n c_{An}) = \alpha C_A$
 $B = (\beta_1 c_{B1}, \beta_2 c_{B2}, \dots, \beta_m c_{Bm}) = \beta C_B$
Result:
 Similarity between A and B $\text{sim}(A, B)$
Procedure:
 $C'_A =$
 for each item in C_B :
 if item in C_A :
 $C'_A = C'_A \cdot \text{append}(\alpha_{\text{index}(\text{item})})$
 else:
 $C'_A = C'_A \cdot \text{append}(0)$
 $\text{sim}(A, B) = (\sum_{i=1}^m (\alpha_i \times \beta_i)) / \sqrt{\sum_{i=1}^m (\alpha_i)^2} \times \sqrt{\sum_{i=1}^m (\beta_i)^2}$, ($\alpha_i \in C'_A, \beta_i \in \beta$)

ALGORITHM 1: The algorithm of calculating similarity between interests.

4.1. Datasets and Metrics

4.1.1. Datasets. Experiments in this section use three real datasets to test user interest mining performance. We processed the original data by the following steps: (1) converting letters to lowercase; (2) deleting the stop words; and (3) deleting words whose document frequency is less than 5.

(1) *Flickr* [32]. The dataset is based on the Flickr website (<http://www.Flickr.com>) and contains 354,531 pieces of personal information and 2,222,379 image annotations which were uploaded from October to December in 2012. These users are from 20 interest categories.

(2) *TWEETS*. Based on Twitter (<https://twitter.com>), we selected 30 hashtags as interest categories and sampled 253,159 tweets for three consecutive months from July to September in 2009 under these hashtags. After pre-processing, we obtained 16,753 words.

(3) *Instagram*. The dataset is based on the Instagram website (<http://www.instagram.com>) and contains 163,479 pieces of personal information and 1,048,575 image annotations which were uploaded from October to December in 2016. These users are from 20 interest categories.

Three datasets are divided into two parts for training and testing. Among them, 10% of the annotated words for each user in the third month are considered as test data, and the rest are considered as training data. Table 1 gives a summary of the dataset.

4.1.2. Metrics. We use the following indicators to evaluate the model:

(1) *Perplexity* [2]. For quantitatively comparing multiple models with different hypotheses and inference mechanisms, perplexity of the pattern to which per-word on the test dataset belongs was computed, which is defined as

$$\text{Perplexity}(U_{\text{test}}) = \exp - \frac{\sum_{i=1}^N \ln p(c_{u_i})}{\sum_{i=1}^N N_{u_i}}, \quad (10)$$

where U_{test} represents the users of the test dataset, N represents the number of users on the test dataset, c_{u_i} represents the set of patterns of user u_i , $p(c_{u_i})$ represents the generation probability of the pattern of user u_i based on the proposed model, and N_{u_i} represents the total number of the patterns of user u_i . The smaller the perplexity is, the higher the likelihood estimation and the better the performance of the model is.

(2) *Classification Accuracy* [29]. One of the purposes of user interest modeling is to get the proportion of topics for each document, which provides a potential semantic representation of the user's interests. This indicator is intended to determine the accuracy and discriminability of the latent semantic representation of user interest. The classification accuracy is defined as follows:

$$\text{Accuracy}(U_{\text{test}}) = \frac{1}{N} \times \sum_{d \in U_{\text{test}}} I(C_d = P_d), \quad (11)$$

where U_{test} is the users of the test dataset, N represents the number of users on the test dataset, d is the document composed of some user's annotated words, I is an indicator function, C_d is the user's actual interest category, and P_d is the predicted user category. The larger the value of accuracy, the more accurate the potential semantic representation of user interest generated by the model and the better the performance of the model.

(3) *Normalized Mutual Information* [33]. Perplexity is a commonly used metric for user interest models, but it does not directly measure the semantic consistency of learning user interest. Therefore, in order to further evaluate the quality of user interest generated by models, we used another evaluation metric, normalized mutual information (NMI), which is used to assess how well the predicted interest matches the actual interest. The definition of NMI is as follows:

$$\text{NMI}(C, P) = \frac{I(C, P)}{[H(C) + H(P)]/2}, \quad (12)$$

where C is the actual user interest set, P is the predicted user interest set, $H(C)$ and $H(P)$ are the entropy of the random variable, and $I(C, P)$ is the mutual information between C

TABLE 1: The information of the dataset.

Dataset	Number of users	Number of labels
Flickr	354,531	2,222,379
TWEETS	253,159	16,753
Instagram	163,479	1,048,575

and P . The value of NMI is between 0 and 1. The closer to 1, the more consistent the predicted interest and actual interest, on the contrary, the more independent the predicted interest and the actual interest.

4.2. Results and Analysis Based on the Interest Tracking Model.

The first experiment explores the user's long-term interest in three datasets. Comparison models are as follows:

- (1) LDA [2]: a standard static topic model based on bag-of-words representation.
- (2) LDA-U [15]: a standard static topic model based on unigram representation.
- (3) DTM [3]: a probabilistic time series model used to analyze the temporal evolution of topics in large document sets. DTM is a dynamic topic model that is commonly used nowadays. It looks for the connection of superparameters for each time slice topic based on LDA.
- (4) TTM [4]: based on the content of the document and the previous estimated distribution, a dynamic topic distribution of long texts at time slice t is captured.
- (5) PGMult [5]: an improved dynamic topic model that focuses on discrete data with certain dependency. PGMult reconstructs the multiple distribution using latent variables with joint Gaussian likelihood, which takes advantage of the logistic stick-breaking representation and the Polya-gamma augmentation.
- (6) UCIT-L [25]: an improved dynamic topic model that tracks users' dynamic interests based on users' topic distributions at not only the last time period but also other multiple time periods in the past. Parameter settings are as follows:
 - (1) Parameter setting for frequent pattern mining: \min_{up} , which represents the minimum support. We set the minimum support $\min_{\text{up}} \in \{0.01\%, 0.05\%, 0.1\%, 1\%, 3\%\}$ [26] and then select the best performance by cross-validation.
 - (2) Parameter setting of interest models: based on different datasets, for LDA and LDA-U, we consider annotated words of each user in three months as a document in the DTM, TTM, PGMult, UCIT-L, and ITM, we set time slice as $[U_{\text{train}}, U_{\text{train}}, U_{\text{train}}]$ which represents that the same users added annotations to the uploaded image based on their own interests in three months. The above five models all set the number of interests from 5 to 185. At the same time, the Gibbs sample of 1,000 iterations runs 10 times, and the average is calculated. The Dirichlet prior parameter of the "user-interest"

distribution is set as $\alpha = K/50$, where K is the number of interests, and the Dirichlet prior parameter of the "interest-frequent pattern" distribution is set as $\beta = 0.01$. Figure 3 shows the perplexity of Flickr, TWEETS, and Instagram for different topic numbers. It can be seen that the perplexity of ITM is always lower than LDA, LDA-U, DTM, TTM, PGMult, and UCIT-L models. In addition, since the distribution of topics in short texts is very sparse, perplexity does not increase as the number of interests increases for each model. This not only proves that mapping short texts to the frequent pattern space is effective but also proves that time series models are superior to static topic models in describing user interests over a longer period of time.

Then, we use the classification accuracy to evaluate the effectiveness of different models to describe the user's long-term interest. It can be observed from Table 2 that the classification accuracy of ITM is always higher than that of LDA, LDA-U, DTM, TTM, PGMult, and UCIT-L on different datasets. Especially for Flickr and Instagram, since each document consists of annotation words that includes multiple topics over a long period of time, ITM which represents each document with frequent patterns can fit the classifier well. LDA cannot solve the problem of sparsity and interest evolution. DTM and TTM cannot solve the problem of sparsity. Compared with representing each document with frequent patterns to alleviate the sparsity problem in ITM, clustering users based on short-text streams of multiple time periods of users and their followers to solve the sparsity problem in UCIT-L is not effective. PGMult cannot solve the problem of interest evolution. Therefore, they are not as good as ITM in describing users' long-term interests.

We also use NMI to evaluate the semantic consistency between user interest generated by different models and the actual interests of users. From Table 2, we can see that the NMI is generally lower. However, ITM is again significantly better than other models. At the same time, the NMI of LDA, DTM, TTM, and LDA-U on TWEETS is very low, indicating that these three models cannot model good topic representations for short-text tweets. From the NMI of each model on Flickr and Instagram, PGMult does not perform as well as ITM in finding the semantic connection of discrete data. The NMI of UCIT-L on three datasets is lower than ours, demonstrating that if the user's interest is not closely related to the followers' interest, UCIT-L's performance in mining long-term interest based on short texts is limited.

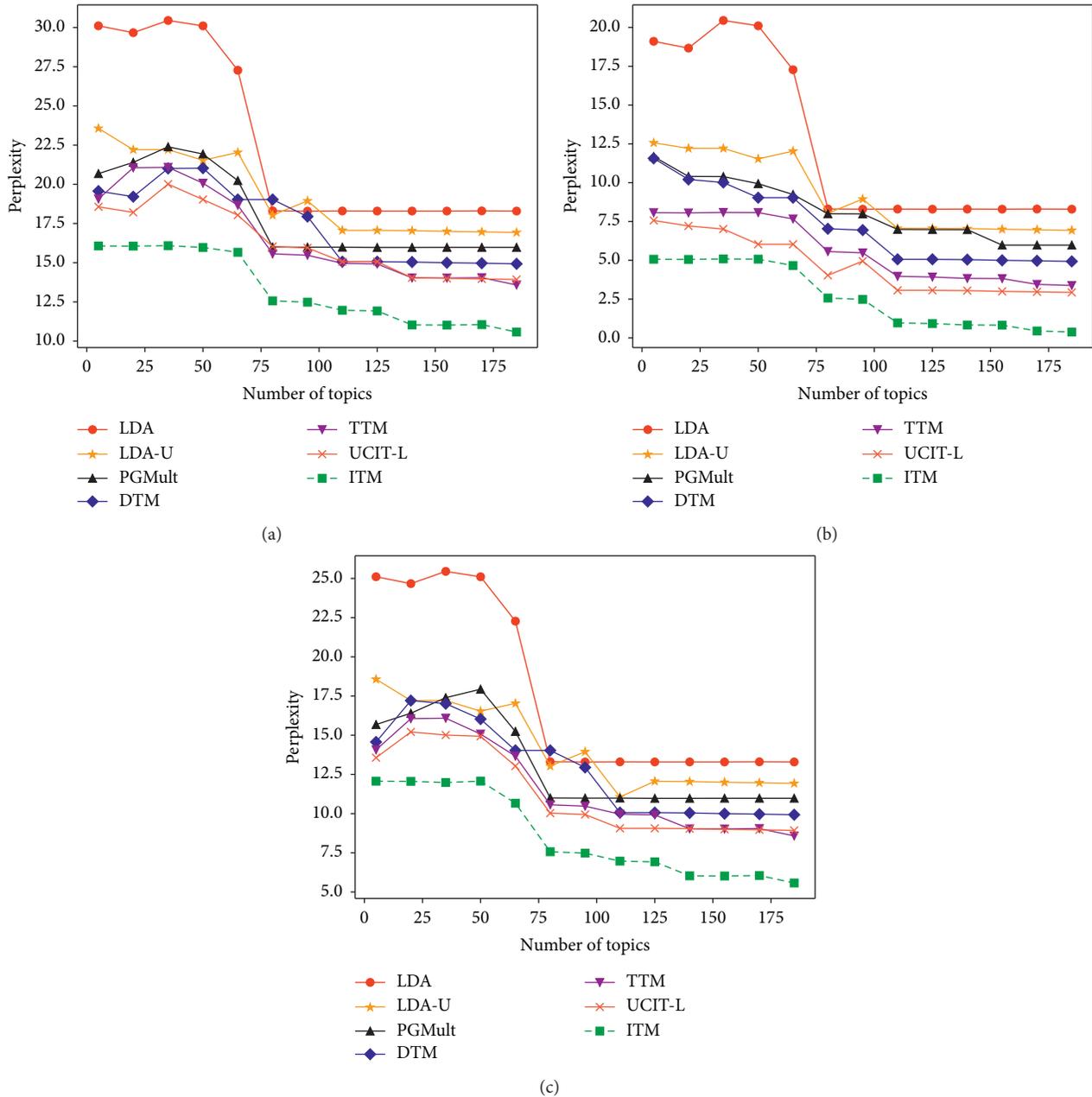


FIGURE 3: Perplexity of mining user long-term interest on different datasets. (a) Flickr. (b) TWEETS. (c) Instagram.

TABLE 2: The average accuracy and NMI of mining user long-term interest.

	Accuracy			NMI		
	Flickr	TWEETS	Instagram	Flickr	TWEETS	Instagram
LDA	0.58	0.71	0.58	0.03	0.05	0.03
LDA-U	0.65	0.79	0.69	0.07	0.08	0.08
PGMult	0.70	0.76	0.75	0.10	0.06	0.10
DTM	0.68	0.81	0.70	0.08	0.09	0.09
TTM	0.70	0.85	0.73	0.10	0.11	0.10
UCIT-L	0.72	0.86	0.75	0.10	0.13	0.10
ITM	0.78	0.90	0.81	0.13	0.15	0.13

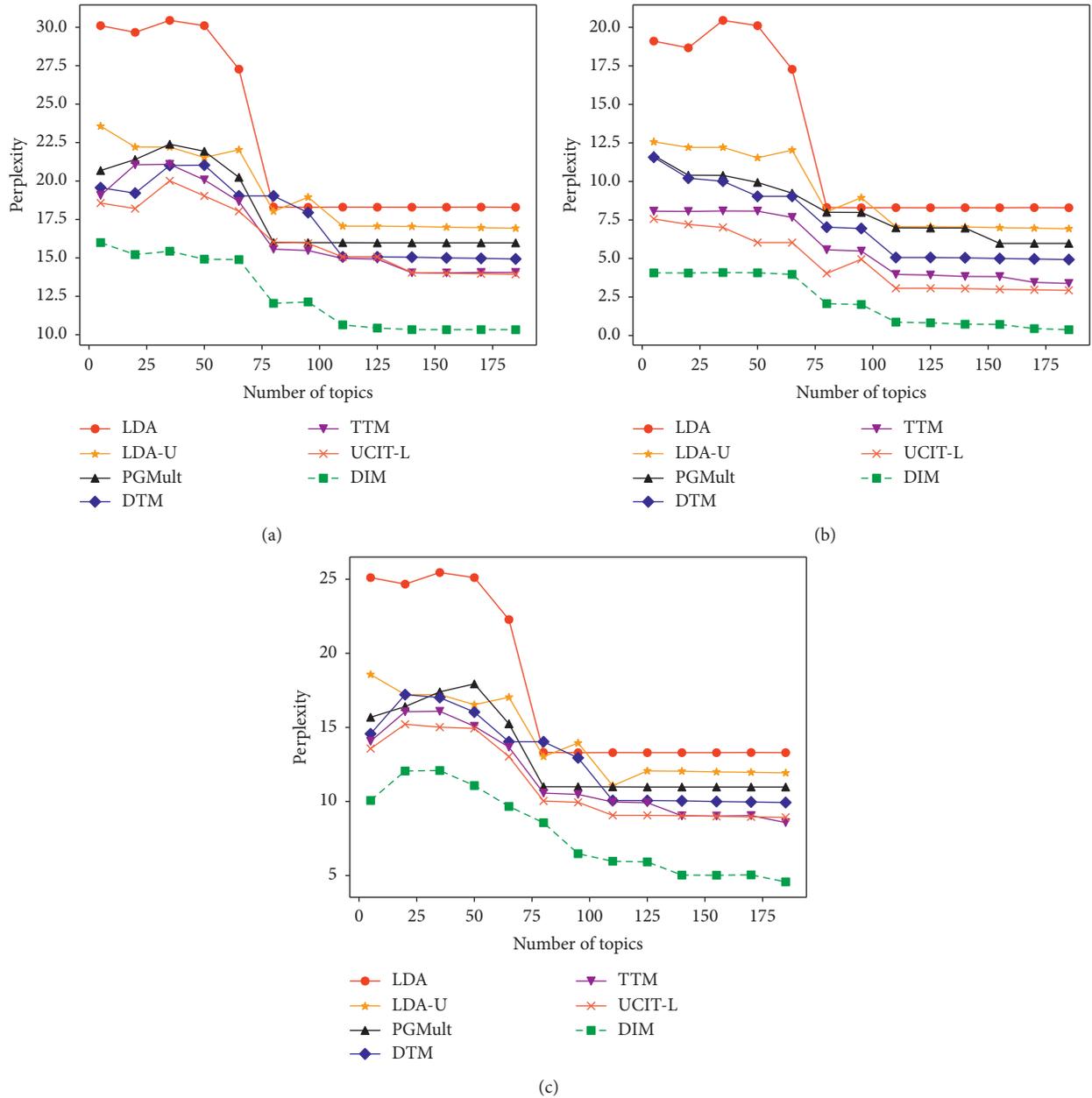


FIGURE 4: Perplexity of mining user interest on different datasets. (a) Flickr. (b) TWEETS. (c) Instagram.

4.3. Results and Analysis of Fusion Performance Based on the Deformable Interest Model. The second experiment combines user long-term interest and situational interest based on the deformable interest model in three datasets and compares the result with that of the LDA, LDA-U, DTM, TTM, PGMult, and UCIT-L. The experimental parameter setting of frequent pattern mining, LDA, LDA-U, DTM, PGMult, and UCIT-L is the same as that of the first experiment. The experimental parameters of DIM are set as follows:

- (1) DIM: for DIM, we set the initial deformation coefficient $k_1 = 1, k_2 = 0, k_3 = 0$ and the initial bias $b = 0$. The number of interests is set to be the same as the other four models. Figure 4 shows the perplexity of Flickr, TWEETS, and Instagram for different topic

numbers. It can be seen that the perplexity of DIM is always lower than LDA, LDA-U, DTM, TTM, PGMult, and UCIT-L. It can be observed from Table 3 that the classification accuracy of DIM is always higher than that of LDA, LDA-U, DTM, TTM, PGMult, and UCIT-L on different datasets. It can be observed from Table 3 that, from the perspective of NMI, the performance of DIM is also significantly better than other models. User interest consists of user long-term interest and situational interest. Situational interest is considered to be triggered by certain conditions or stimuli in the environment. It is a relatively passive and short-lived emotional state, so situational interest is not only the short-term interest of users but also includes the current hot

TABLE 3: The average accuracy and NMI of mining user interest.

	Accuracy			NMI		
	Flickr	TWEETS	Instagram	Flickr	TWEETS	Instagram
LDA	0.58	0.71	0.58	0.03	0.05	0.03
LDA-U	0.65	0.79	0.69	0.07	0.08	0.08
PGMult	0.70	0.76	0.75	0.10	0.06	0.10
DTM	0.68	0.81	0.70	0.08	0.09	0.09
TTM	0.70	0.85	0.73	0.10	0.11	0.10
UCIT-L	0.72	0.86	0.75	0.10	0.13	0.10
DIM	0.79	0.94	0.81	0.15	0.20	0.13

TABLE 4: The average accuracy of mining different interests.

	Flickr	TWEETS	Instagram
LDA-FP	0.70	0.81	0.73
KB	0.21	0.95	0.30
ITM	0.78	0.90	0.81
DIM	0.79	0.94	0.81

topics. Models such as LDA, LDA-U, DTM, TTM, PGMult, and UCIT-L default user interest are not affected by the environment, and users' long-term interest is regarded as user interests. DIM combines users' long-term interest, users' short-term interest, and current hot topics. One of the qualities that social networks attract users is the sharing and exchange of information. This trait determines the environment greatly affects the evolution of user interest. From Tables 2 and 3, on the TWEETS, the classification accuracy of DIM is greatly improved compared with ITM, while on the other two datasets, the classification accuracy of DIM is almost unchanged compared with ITM. This is because the content of TWEETS is more susceptible to situational interests [34], especially social hot topics, and Flickr and Instagram have gathered more photographers, and their interests are relatively stable. Table 4 shows the accuracy of ITM which describes long-term interest, LDA-FP which describes short-term interest, KB which describes the social hot topic, and DIM on different datasets. It can be seen that, on TWEETS, the tweets sent by users follow social hot topics. Therefore, DIM is a general-purpose model that describes user interests, and user interest using DIM mining is more in line with real user interest.

From Figures 4 and 5, it can be seen that although DIM brings a slightly higher time cost, perplexity is significantly reduced. When the model performed best ($T = 185$), the time spent on DIM is also no more than three minutes. With the development of industrial technology, the computing power of computers has increased rapidly, and this time is acceptable.

4.4. Results and Analysis of Adaptive Update Performance Based on the Deformable Interest Model. In the third experiment, adaptively adjusting the degree of deformation of

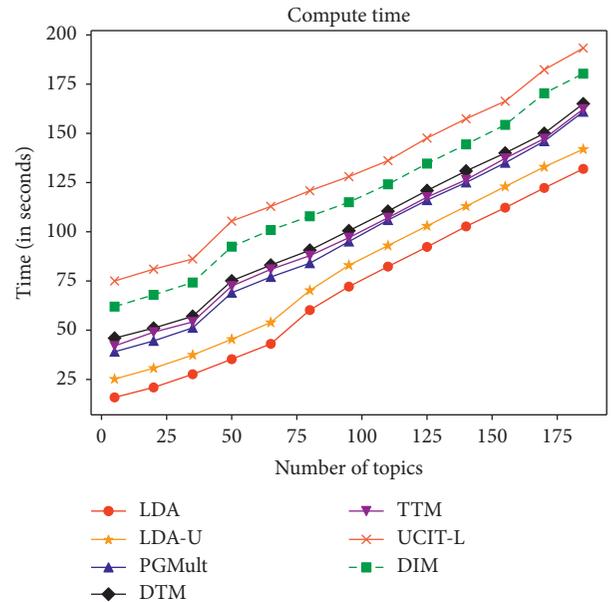


FIGURE 5: Compute times for different models.

subinterests based on the deformable interest model is carried out on three datasets, and the experimental results are compared with those of static-DIM. Subinterests of two models are produced in the same way. The DIM experimental parameter setting is the same as that in the second experiment. The remaining experimental parameters are set as follows:

- (1) Static-DIM: the model fixed interest weights and sets long-term interest weight as $k_1 = 0.7$, short-term interest weight as $k_2 = 0.2$, and hot topic weight as $k_3 = 0.1$ based on experience

Figure 6 shows the perplexity of Flickr, TWEETS, and Instagram for different topic numbers. It can be seen that the perplexity of DIM is always lower than static-DIM. It can be

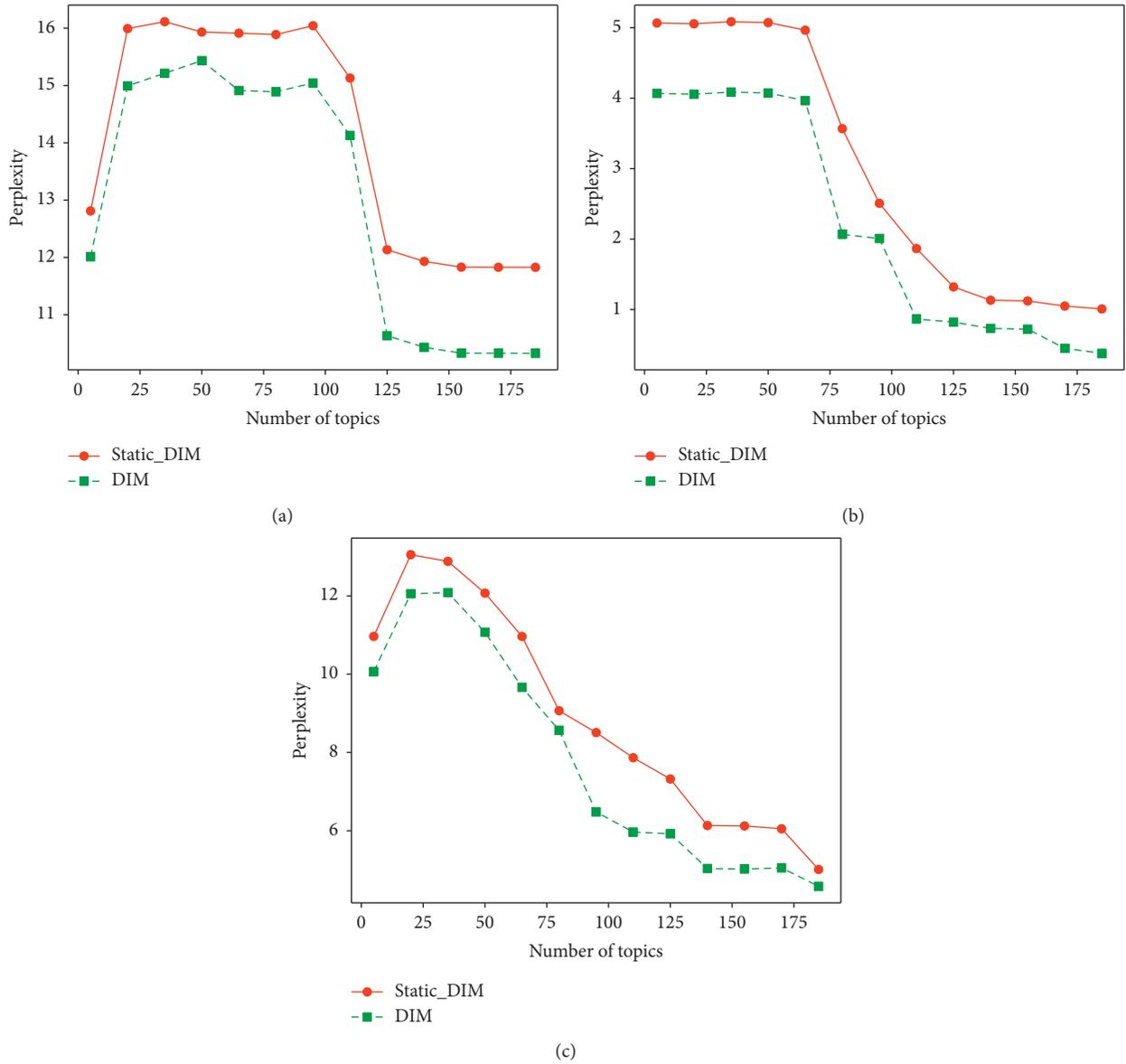


FIGURE 6: Perplexity of combining user interest on different datasets. (a) Flickr. (b) TWEETS. (c) Instagram.

TABLE 5: The average accuracy and NMI of combining user interests.

	Accuracy			NMI		
	Flickr	TWEETS	Instagram	Flickr	TWEETS	Instagram
Static-DIM	0.73	0.71	0.75	0.10	0.06	0.09
DIM	0.79	0.94	0.81	0.15	0.20	0.13

observed from Table 5 that the classification accuracy of DIM is always higher than that of static-DIM on different datasets. It can be observed from Table 5 that, from the perspective of NMI, the performance of DIM is also significantly better than static-DIM. We can conclude that the status of the three subinterests is equal, which means that their influence is not static. It is not difficult to find out that the user will take photos for social hot topics and even

generate new interests. The fixed interests can not only correctly describe the relationship between subinterests and real interest and subinterests and subinterests but also indicate user interest can only be changed slightly and impossible to update, which is obviously not realistic. In general, our experimental results demonstrated that the perplexity is reduced to 0.378, the average accuracy is increased to 94%, and the average NMI is increased to 0.20,

which proved the model can effectively mine user interest, especially for short-text streams with discrete semantics.

5. Conclusion

In this paper, we propose a novel method based on the deformable interest model (DIM) for modeling the evolution of user interest in dynamic social networks. We introduce the time factor and leverage sophisticated interest tracking model (ITM) which is based on two-layer Bayesian model to describe dynamic user long-term preferences. Compared with traditional models, it not only describes the evolution process of the user's long-term interest without increasing the number of the latent variables but also maps annotated words to the frequent pattern space to solve the sparsity problem of short texts. We then obtain user interest which combines long-term interest and situational interest by DIM. Compared with traditional models, DIM proposes an objective function which not only fully considers the composition of the user's real interests but also adaptively updates the influence of long-term interest, short-term interest, and hot topics on user interest. We evaluated the performance of the proposed model in terms of perplexity, accuracy, and NMI and made comparisons with state-of-the-art models. The experimental results demonstrate the effectiveness of the introduced model. This enlightens us that the model can be applied in the field of image retrieval or e-commerce so that users can quickly find pictures or commodities that match their interests. It can also be used in social networks to present users with information streams that match their interests.

In future work, we intend to use the deformable interest model (DIM) to annotate areas in the image which the user is interested in. Like most previous works, how to calculate the cross-modal similarity between images and text is also a challenge. Therefore, our future work is to study this problem by extending the model proposed in this paper.

Data Availability

The data used to support the findings of this study are included within the article.

Conflicts of Interest

The authors declare that they have no conflicts of interest.

Acknowledgments

This work was supported by the National Natural Science Foundation of China (no. 61672108).

References

- [1] A. Krapp, "Structural and dynamic aspects of interest development: theoretical considerations from an Ontogenetic Perspective," *Learning and Instruction*, vol. 12, no. 4, pp. 383–409, 2002.
- [2] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," *Journal of Machine Learning Research*, vol. 3, pp. 993–1022, 2003.
- [3] D. M. Blei and J. D. Lafferty, "Dynamic topic models," in *Proceedings of the 23rd International Conference on Machine Learning*, pp. 113–120, Pittsburg, PA, USA, 2006.
- [4] T. Iwata, S. Wanatabe, T. Yamada et al., "Topic tracking model for analyzing consumer purchase behavior," in *Proceedings of the 18th International Joint Conference on Artificial Intelligence*, pp. 1427–1432, New York, NY, USA, 2009.
- [5] S. W. Linderman, M. J. Johnson, and R. P. Adams, C. Cortes, "Dependent multinomial models made easy: stick-breaking with the polya-gamma augmentation," Edited by C. Cortes, Ed., pp. 3456–3464, NIPS, Vancouver, Canada, 2015.
- [6] C. Wang, D. Blei, and D. Heckerman, "Continuous time dynamic topic models," 2012.
- [7] S. Y. Song, Q. D. Li, and H. Y. Bao, "Detecting dynamic association among Twitter topics," in *Proceedings of the 21st International Conference on World Wide Web*, pp. 605–606, New York, NY, USA, 2012.
- [8] G. P. Nicholas, G. S. James, and W. Jesse, "Bayesian inference for logistic models using Polya-gamma latent variables," *Journal of the American Statistical Association*, vol. 108, no. 504, pp. 1339–1349, 2013.
- [9] A. Acharya, J. Ghosh, and M. Y. Zhou, "A dual Markov chain topic model for dynamic environments," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 1099–1108, New York, NY, USA, 2018.
- [10] X. Q. Gui, J. Zhang, X. M. Zhang et al., "Survey on temporal topic model methods and application," *Computer Science*, vol. 44, no. 2, pp. 46–55, 2017.
- [11] S. Tobias, "Interest, prior knowledge, and learning," *Review of Educational Research*, vol. 64, no. 1, pp. 37–54, 1994.
- [12] A. Chen, P. W. Darst, and R. P. Pangrazi, "An examination of situational interest and its sources," *British Journal of Educational Psychology*, vol. 71, no. 3, pp. 383–400, 2001.
- [13] S. Hidi, "Interest and its contribution as a mental resource for learning," *Review of Educational Research*, vol. 60, no. 4, pp. 549–571, 1990.
- [14] A. B. Dieng, F. J. Ruiz, and D. M. Blei, "Topic modeling in embedding spaces," 2019, <https://arxiv.org/abs/1907.04907>.
- [15] H. M. Wallach, "Topic modeling: beyond bag-of-words," in *Proceedings of the 23rd International Conference on Machine Learning*, pp. 977–984, Beijing, China, 2006.
- [16] X. R. Wang, A. McCallum, and X. Wei, "Topical n-grams: phrase and topic discovery, with an application to information retrieval," in *Proceedings of the 7th IEEE International Conference on Data Mining*, pp. 697–702, Beijing, China, 2007.
- [17] P. Jhnicen, F. Wenzel, M. Kloft, and S. Mandt, "Scalable generalized dynamic topic models," in *Proceedings of the International Conference on Artificial Intelligence and Statistics*, pp. 1427–1435, Canary Islands, Spain, 2018.
- [18] U. Yun, D. Kim, E. Yoon, and H. Fujita, "Damped window based high average utility pattern mining over data streams," *Knowledge-Based Systems*, vol. 144, pp. 188–205, 2018.
- [19] D. Kim and U. Yun, "Mining high utility itemsets based on the time decaying model," *Intelligent Data Analysis*, vol. 20, no. 5, pp. 1157–1180, 2016.
- [20] G. Lee and U. Yun, "Single-pass based efficient erasable pattern mining using list data structure on dynamic incremental databases," *Future Generation Computer Systems*, vol. 80, pp. 12–28, 2018.

- [21] U. Yun and G. Lee, "Incremental mining of weighted maximal frequent itemsets from dynamic databases," *Expert Systems With Applications*, vol. 54, pp. 304–327, 2016.
- [22] A. B. Dieng, F. J. R. Ruiz, and D. M. Blei, "The dynamic embedded topic model," 2019.
- [23] R. Bamler and S. Mandt, "Dynamic word embeddings," in *Proceedings of the 34th International Conference on Machine Learning*, pp. 380–389, Sydney, Australia, 2017.
- [24] M. Rudolph and D. M. Blei, "Dynamic embeddings for language evolution," in *Proceedings of the 2018 World Wide Web Conference*, pp. 1003–1011, Geneva, Switzerland, 2018.
- [25] S. Liang, E. Yilmaz, and E. Kanoulas, "Collaboratively tracking interests for user clustering in streams of short texts," *IEEE Transactions on Knowledge and Data Engineering*, vol. 31, no. 2, pp. 257–272, 2019.
- [26] J. Han, J. Pei, Y. Yin, and R. Mao, "Mining frequent patterns without candidate generation: a frequent-pattern tree approach," *Data Mining and Knowledge Discovery*, vol. 8, no. 1, pp. 53–87, 2004.
- [27] Q. Zhang, Q. Wang, J. X. Hao et al., "Mapping smart tourism research in China: a semantic and social network analysis using CiteSpace," in *Proceedings of the 13th International Conference on Service Systems and Service Management (ICSSSM)*, Kunming, China, 2016.
- [28] J. Yang, C. Cheng, S. Shen et al., "Comparison of complex network analysis software: citespace SCI2 and Gephi," in *Proceedings of the IEEE International Conference on Big Data Analysis (ICBDA)*, pp. 169–172, Beijing, China, 2017.
- [29] H. Li, *Statistical Learning methods*, Tsinghua University Press, Beijing, China, 2012.
- [30] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 9, pp. 1627–1645, 2010.
- [31] Q. Liu and S. Li, "Word similarity computing based on how-net," *IJCLCLP*, vol. 7, no. 2, 2002.
- [32] B. Thomee, D. A. Shamma, G. Friedland et al., "Yfcc100M," *Communications of the ACM*, vol. 59, no. 2, pp. 64–73, 2016.
- [33] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*, Cambridge University Press, Cambridge, UK, 2008.
- [34] L. Hong and B. D. Davison, "Empirical study of topic modeling in Twitter," in *Proceedings of the First Workshop on Social Media Analytics*, pp. 80–88, Washington, DC, USA, 2010.