

# Mining User Preferences from Online Reviews with the Genre-aware Personalized neural Topic Model

Anonymous Author(s)  
Submission Id: 2800

## Abstract

Customer-generated reviews on e-commerce websites often contain valuable insights into users' interests in product genres and provide a rich source for mining user preferences. However, most existing neural topic models tend to generate meaningless topics that have low correlations with product genres. Furthermore, they often fail to mine user preferences and discover personalized topic profiles due to the absence of explicit user modeling. To address these limitations, we propose a novel Genre-aware Personalized neural Topic Model (GPTM), which incorporates product genre information into the topic modeling process to ensure the relevance between mined topics and product genres. Moreover, it could produce a personalized topic profile for each user by performing user preference modeling. Extensive experimental results on three publicly available Amazon review corpora validate the effectiveness of the proposed GPTM in genre-aware topic modeling. Furthermore, GPTM surpasses state-of-the-art baselines in user preference mining and generating high-quality personalized topic profiles.

## CCS Concepts

• Computing methodologies → Information extraction; • Information systems → Document topic models.

## Keywords

Neural Topic Modeling, User Preference Discovery, Text Mining

### ACM Reference Format:

Anonymous Author(s). 2024. Mining User Preferences from Online Reviews with the Genre-aware Personalized neural Topic Model. In . ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

With the development of e-commerce, customers are increasingly accustomed to shopping online and sharing their experiences and opinions on websites. For example, Figure 1 shows review posts for books (e.g., 'Everyone Communicates, Few Connect', 'Your Health Destiny', etc.) from different genres. Such review content often reflects users' genre interests and provides a valuable source for mining user preferences and personalized topics.

As the primary data mining tool, conventional topic models, such as the Latent Dirichlet Allocation (LDA) [2], and emerging Neural

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
Conference'17, July 2017, Washington, DC, USA

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM  
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

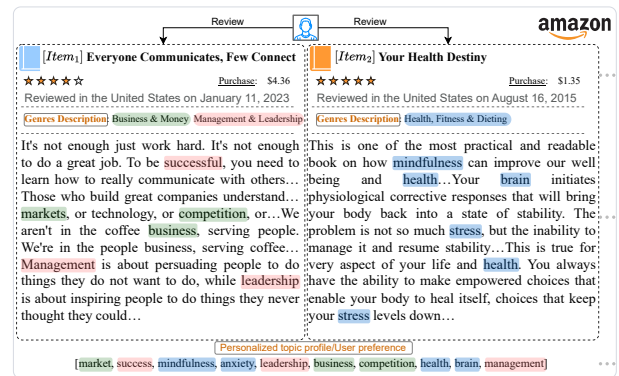


Figure 1: Review examples posted by a user to multiple books with various genres, words in green and red are related to 'Business & Money' and 'Management & Leadership' genres, and blue words denote the 'Health, Fitness & Dieting' genre.

Topic Models (NTM) like the Embedding Clustering Regularization Topic Model (ECRTM) [41] have been extensively explored. Notably, contextualized neural topic models, such as the Contextualized Topic Modeling with Negative sampling (CTMNeg) [1] and the Contextualized Word Topic Model (CWTM) [7], significantly boost the topic quality by incorporating pre-trained language models [6, 29]. Nevertheless, none of these approaches is capable of mining user preferences or discovering personalized topic profiles due to the absence of explicit user modeling.

To mine user preferences and personalized topic profiles, Liu et al. propose the Neural Personalized Topic Model (NPTM) [24], which models personalized topics with a mixture of topic word distributions weighted by user preference distribution. However, it falls short in the following aspects when dealing with Amazon<sup>1</sup> reviews: 1). It may produce genre-irrelevant topics, as it is less capable of incorporating genre description information shown in Figure 1. 2). It uses the Gaussian prior in the latent topic space, which is unsuitable for text modeling [37] and leads to incoherent topics. 3). It follows an autoencoding framework, which often faces mode collapse [33], thereby sacrificing topic diversity.

Thus, to address the above limitations, we propose the Genre-aware Personalized neural Topic Model (GPTM), which incorporates product genre descriptions into the modeling process to generate genre-aware topics and ensure accurate user preference mining. Specifically, GPTM utilizes a topic-inference network that creates a projection from genre-aware text representations to the document-topic distributions to capture genre-aware topics. To incorporate genre information, it first employs a pre-trained transformer [6] to form the genre-aware document representation. Also, topics

<sup>1</sup><https://amazon.com>

are modeled with Dirichlet distribution to ensure interpretability. Then, GPTM leverages a user-inference network to produce user preference distribution over topics, guided by document-topic distributions of user-generated reviews. After the asynchronous genre-aware and user-aware contrastive learning, optimized inference networks (topic and user) could produce genre-aware topics and user preference distributions over topics, which could be further employed to construct personalized topic profiles.

The main contributions of this paper could be summarized as:

- We propose the novel Genre-aware Personalized neural Topic Model (GPTM), which could mine genre-aware topics and produce personalized topic profiles to indicate user preference, based on contrastive learning.
- GPTM incorporates genre information into the topic modeling process and ensures relevance between topics and product genres. Moreover, it utilizes user-aware contrastive learning for user preference mining.
- Experimental results on three Amazon review datasets reveal that the proposed GPTM outperforms state-of-the-art baselines in terms of topic coherence and diversity, while maintaining stronger correlations between topics and genres. Moreover, GPTM surpasses the competitive NPTM on the proposed Personalized Hit Rate (*PHR*) and Personalized Genre Correlation (*PGC*) metrics, demonstrating its superiority in personalized topic modeling and user preference mining.

## 2 RELATED WORK

In this section, we briefly review two related lines of research which are neural topic modeling and contrastive representation learning.

### 2.1 Neural Topic Modeling

Neural topic modeling [13], a recently emerged research topic, has attracted a lot of interest in the Natural Language Processing (NLP) community and made some progress.

Early pioneers, such as the Neural Variational Document Model (NVDM) [27] and the Adversarial-neural Topic Model (ATM) [39], often follow the Bag-of-Words (BOW) assumption and generative models like VAE [16] and GAN [10]. Followed by NVDM, Srivastava et al. employed the Logistic-Normal distribution to model topics and proposed the Neural Variational Latent Dirichlet Allocation (NVLDA) [34]. On the other hand, Hu et al. extended the ATM and proposed the Topic Modeling with Cycle-consistent Adversarial Training (ToMCAT) [14].

Furthermore, scholars have also explored how to boost modeling performance by incorporating word embeddings and contextualized language models. Wu et al. proposed the Embedding Clustering Regularization Topic Model (ECRTM) [41] by forcing topic embeddings to be centers of the word embeddings clusters. To capture context information among texts, Adhya et al. utilized a pre-trained language model to conduct topic inference and proposed the Contextualized Topic Model with Negative sampling (CTMNeg) [1] based on contrastive learning. Fang et al. leveraged the contextualized word embedding from Bert and proposed the Contextualized Word Topic Model (CWTM) [7]. However, these approaches could

not model user preferences as they do not explicitly incorporate user information into the modeling process.

To mine personalized topics that reflect user preferences, Liu et al. incorporated contextualized document representations and user preferences into the modeling process and proposed the Neural Personalized Topic Model (NPTM) [24]. And our work differs from NPTM in the following aspects: 1). Unlike NPTM which is prone to extract genre-irrelevant topics, GPTM incorporates product genre descriptions into the modeling process to mine genre-aware topics and ensure accurate user preference mining. 2). Unlike NPTM which models topics with Gaussian, GPTM utilizes the Dirichlet prior in the topic space to enhance the interpretability. 3). GPTM follows the contrastive learning framework which could tackle the topic collapse problem.

### 2.2 Contrastive Representation Learning

Contrastive Representation Learning (CRL), a popular unsupervised learning paradigm, has recently achieved state-of-the-art performance for visual [5] and textual [42] representation learning.

The intuitive idea of CRL is to pull similar samples closer and push dissimilar samples apart by maximizing the similarities of similar samples and minimizing those of dissimilar pairs in a shared representation space [19, 21, 32]. Recently, contrastive representation learning has gained significant traction in the NLP community, and it has been proven to be effective in learning sentence representation [3], multi-modal sentiment analysis [26], neural machine translation [22] and fake news fact checking [44].

Scholars have also explored whether CRL could alleviate mode collapse [18] and generate diverse samples [20]. Su et al. proposed a contrastive-based framework [35] for diverse text generation. Zhong et al. proposed Graph Contrastive Clustering (GCC) [45] to address the dimension collapse in graph clustering. Thus, we follow their idea and formulate personalized topic modeling as a contrastive learning task to alleviate topic collapse.

## 3 PROBLEM FORMULATION

Given a document corpus  $D = \{x_{i^1, u^1}, x_{i^2, u^2}, \dots, x_{i^N, u^N}\}$ , collected from the reviews posted by a set of  $N_u$  users  $U = \{u_1, u_2, \dots, u_{N_u}\}$  to a set of  $N_i$  product items  $I = \{i_1, i_2, \dots, i_{N_i}\}$  (each item is associated with multiple genre categories in the genre set  $G = \{g_1, g_2, \dots, g_K\}$ , as shown in Figure 1). For the item  $i$  in  $I$ ,  $G_i = \{g_i^1, g_i^2, \dots, g_i^{N_i}\}$  represents its genre set which contains  $N^i \geq 1$  genres in  $G$ . For the  $n$ -th ( $n \in \{1, 2, \dots, N\}$ ) document  $x_{i^n, u^n}$  in  $D$ , it is the review content posted by  $u^n$  to the item  $i^n$ . Here,  $u^n \in U$  and  $i^n \in I$  mean user and item attached to  $n$ -th review. The aims of our work are: 1). Mining a set of  $K$  genre-aware topics that are semantically consistent with genres in  $G$ . 2). For each user  $u \in U$ , inferring the user preference distribution  $\vec{p}_u$  over topics and producing a personalized topic profile  $\vec{\phi}_u$  that reflects his/her interests.

## 4 METHODOLOGIES

As shown in Figure 2 (a), our proposed Genre-aware Personalized neural Topic Model (GPTM) contains four components which are: 1). Text Augmentation and Representation module (top-left): It first conducts text augmentation for each review  $x_{i, u} \in X$  to build semantically consistent pair  $x_{i, u}^a$  and  $x_{i, u}^b$ . Then, a transformer  $\mathcal{T}$

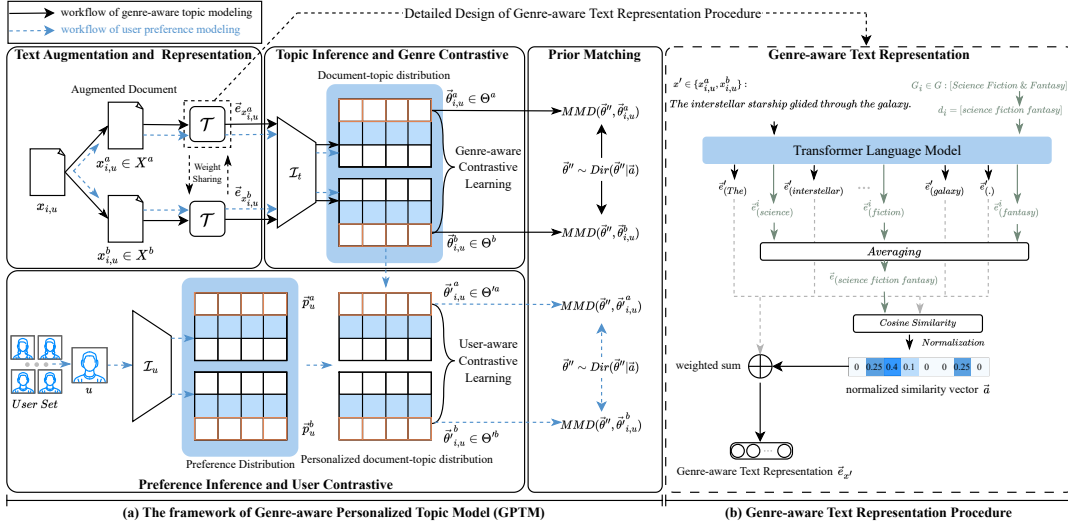


Figure 2: The framework of GPTM (a) and details of Genre-aware Text Representation (b). In sub-figure (a), black arrows denote the workflow of genre-aware topic modeling, and blue dashed arrows represent the workflow of user preference modeling.

Table 1: Key notations and illustrations.

Symbol	Description
Data Representation and Distribution	
$D$	a collection of $N$ reviews posted by users
$V$	vocabulary size of corpus $D$
$X$	a batch of reviews in the corpus $D$
$G_i$	a set of $N^i$ genres in $G$ associated with item $i$
$U, I, G$	a set of $N_u, N_i, K$ users, items, genre categories
$N_u, N_i$	the number of users, product items in the corpus
$X^* = X^a \cup X^b$	two batches of augmented reviews from $X$
$\Theta = \Theta^a \cup \Theta^b$	inferred document-topic distributions of $X^*$
$\Theta' = \Theta'^a \cup \Theta'^b$	inferred personalized document-topic distributions of $X^*$
Model Parameters	
$x_{i,u}$	review posted by user $u$ to the item $i$
$\mathcal{T}$	transformer language model for text representation
$H_V$	dimension of transformer embeddings
$H$	hidden units of inference networks $\mathcal{I}_i$ and $\mathcal{I}_u$
$\mathcal{I}_i, \mathcal{I}_u$	topic-inference network, user-inference network
$\tilde{\phi}_u$	personalized topic profile/distribution of user $u$
$\tilde{e}_{d_i}$	contextualized representation of genre descriptions of item $i$
$x_{i,u}^a, x_{i,u}^b$	augmented semantically consistent text pair of $x_{i,u}$
$\tilde{\theta}_{i,u}^a, \tilde{\theta}_{i,u}^b$	genre-aware text representations of $x_{i,u}^a$ and $x_{i,u}^b$
$\tilde{\theta}_{i,u}^a, \tilde{\theta}_{i,u}^b$	document-topic distributions of augmented $x_{i,u}^a$ and $x_{i,u}^b$
$\tilde{p}_u^a, \tilde{p}_u^b$	preference distribution of user $u$ (correspond to $\tilde{\theta}_{i,u}^a, \tilde{\theta}_{i,u}^b$ )
$\tilde{\theta}_{i,u}^a, \tilde{\theta}_{i,u}^b$	personalized document-topic distributions of $\tilde{\theta}_{i,u}^a$ and $\tilde{\theta}_{i,u}^b$
$W_s = \{w_s^1, w_s^2, w_s^3\}$	selected word set from review $x_{i,u}$ for augmentation
$d_i = [w_i^1, w_i^2, \dots, w_i^{N_{d_i}}]$	genres description of the item $i$ (contain $N_{d_i}$ words)
$I_g(\cdot, \cdot), J_u(\cdot, \cdot)$	genre, user indicator function
$\mathcal{L}_{G_c}, \mathcal{L}_{U_c}$	genre-aware, user-aware contrastive objectives of $X^*$
$\mathcal{L}_{M_g}, \mathcal{L}_{M_u}$	genre-aware, user-aware matching objectives of $X^*$
$\lambda_1, \lambda_2$	coefficient hyper-parameters in Eq. 17 and Eq. 18
$\mathcal{L}_{TM}, \mathcal{L}_{UM}$	objectives of genre-aware topic modeling and user modeling
$C \in \mathbb{R}^{K \times V}$	correlation matrix between topics and words
$\Phi \in \mathbb{R}^{K \times V}$	topic word distribution matrix

is utilized to build genre-aware text representations  $\tilde{e}_{i,u}^a$  and  $\tilde{e}_{i,u}^b$ , which convey genre information of item  $i$ , for augmented pair. 2). Topic Inference and Genre Contrastive module (top-middle): Feeding with text representations  $\tilde{e}_{i,u}^a$  and  $\tilde{e}_{i,u}^b$ , it infers document-topic distributions  $\tilde{\theta}_{i,u}^a$  and  $\tilde{\theta}_{i,u}^b$  with the topic-inference network  $\mathcal{I}_i$ . Besides, it conducts genre-aware contrastive learning to capture

genre-aware topics among texts. 3). Preference Inference and User Contrastive module (bottom-left): Firstly, it infers preference distributions  $\tilde{p}_u^a$  and  $\tilde{p}_u^b$  for user  $u$  with the user-inference network  $\mathcal{I}_u$ . Then, together with inferred document-topic distributions, it constructs personalized document-topic distributions  $\tilde{\theta}_{i,u}^a$  and  $\tilde{\theta}_{i,u}^b$  to conduct user-aware contrastive learning for user preference mining. 4). Dirichlet Prior Matching module (right): It matches the inferred document-topic distributions and personalized document-topic distributions to the Dirichlet prior  $\text{Dir}(\tilde{\theta}^i | \bar{\alpha})$  in the latent topic space. This will ensure the interpretability of mined topics during (genre and user) contrastive learning. Also, Figure 2 (b) depicts the design of the genre-aware text representation mechanism. The functionalities of each component will be discussed in more detail below. For the sake of presentation, Table 1 lists the key notations and illustrations, the left column lists appeared symbols, and the right column is the corresponding illustrations.

#### 4.1 Text Augmentation and Representation

Since maintaining semantic consistency is crucial for contrastive representation learning, we follow Feng et al. [8] and use a WordNet<sup>2</sup> based text augmentation procedure. In detail, for each document  $x_{i,u} = [w_1, w_2, \dots, w_{N_x}]$  in  $D$ , which contains  $N_x$  words, its augmentation process could be summarized as:

- (1) Randomly select three words  $W_s = \{w_s^1, w_s^2, w_s^3\}$  from  $x_{i,u}$ ;
- (2) For each word  $w_s$  in the selected  $W_s$ , obtain its synonym set  $\text{synset}(w_s)$  with WordNet;
- (3) From each  $\text{synset}(w_s)$ , randomly select a substitute word to replace the  $w_s$  in the document  $x_{i,u}$ .

Thus, GPTM could build a pair of semantically similar augmented documents  $x_{i,u}^a$  and  $x_{i,u}^b$  by conducting the above augmentation twice.

<sup>2</sup><https://wordnet.princeton.edu/>

To incorporate genre-categories information  $G_i = \{g_1^i, \dots, g_{N_i}^i\}$  of reviewed item  $i$ , we devise a genre-aware text representation mechanism as shown in Figure 2 (b). Specifically, for the document  $x_{i,u}$ , we first construct its genre description  $d_i = [w_1^i, w_2^i, \dots, w_{N_{d_i}}^i]$  by concatenating genre names in  $G_i$  and obtain its embedding  $\vec{e}_{d_i}$  with:

$$[\vec{e}_1^i, \vec{e}_2^i, \dots, \vec{e}_{N_{d_i}}^i] = \mathcal{T}([w_1^i, w_2^i, \dots, w_{N_{d_i}}^i]) \quad (1)$$

$$\vec{e}_{d_i} = \frac{1}{N_{d_i}} \sum_{l=1}^{N_{d_i}} \vec{e}_l^i \quad (2)$$

where  $N_{d_i}$  means the number of words in genre description of item  $i$ , and  $\vec{e}_l^i$  is the contextualized word representation of the  $l$ -th word in  $d_i$ .

Then, under the guidance of  $\vec{e}_{d_i}$ , we generate the genre-aware text representation  $\vec{e}_{x'}$  for the augmented document  $x' \in \{x_{i,u}^a, x_{i,u}^b\}$ , which contains word sequence  $[w'_1, w'_2, \dots, w'_{N_{x'}}]$ , by weighting their contextualized word representations with the semantic similarities between words and genre description. Concretely,  $\vec{e}_{x'}$  could be calculated with formulas:

$$[\vec{e}'_1, \vec{e}'_2, \dots, \vec{e}'_{N_{x'}}] = \mathcal{T}([w'_1, w'_2, \dots, w'_{N_{x'}}]) \quad (3)$$

$$[a_1, \dots, a_{N_{x'}}] = \text{softmax}([\cos(\vec{e}'_1, \vec{e}_{d_i}), \dots, \cos(\vec{e}'_{N_{x'}}, \vec{e}_{d_i})]) \quad (4)$$

$$\vec{e}_{x'} = \sum_{l=1}^{N_{x'}} a_l \cdot \vec{e}'_l \quad (5)$$

where  $N_{x'}$  denotes the number of words in  $x'$ ,  $\vec{e}'_l$  represents the contextualized word representation of the  $l$ -th word in  $x'$ ,  $a_l$  is the normalized similarity between word  $w'_l$  and genre description  $d_i$ .

## 4.2 Topic Inference and Genre Contrastive

To infer document-topic distributions for augmented pairs, GPTM projects their genre-aware text representations into the topic space with a topic-inference network  $\mathcal{I}_t$ . Concretely, for the augmented document  $x_{i,u}^a$ , its document-topic distribution  $\vec{\theta}_{i,u}^a$  could be inferred with formula:

$$\vec{\theta}_{i,u}^a = \mathcal{I}_t(\vec{e}_{x_{i,u}^a}) \quad (6)$$

where  $\vec{e}_{x_{i,u}^a}$  is the genre-aware text representation of  $x_{i,u}^a$ , and  $\vec{\theta}_{i,u}^b$  could be generated similarly.

Besides, to mine genre-aware topics and ensure that documents with different genres are inferred to distinctive document-topic distributions, the topic-inference network  $\mathcal{I}_t$  is trained with genre-aware contrastive learning.

Specifically, given a batch of documents  $X = \{x_{i_1, u_1}, x_{i_2, u_2}, \dots, x_{i_M, u_M}\}$ , we perform text augmentation twice for each document and generate a set of  $2M$  augmented documents  $X^* = X^a \cup X^b = \{x_{i_1, u_1}^a, x_{i_1, u_1}^b, x_{i_2, u_2}^a, x_{i_2, u_2}^b, \dots, x_{i_M, u_M}^a, x_{i_M, u_M}^b\}$ . For each document  $x_{i_l, u_l}^a$  ( $1 \leq l \leq M$ ), we choose the  $x_{i_l, u_l}^b$  to construct the positive pair  $(x_{i_l, u_l}^a, x_{i_l, u_l}^b)$ . On the other hand, we only select the document  $x_{i^*, u^*}$  with different genre in  $X^*$  to form the negative pairs  $(x_{i_l, u_l}^a, x_{i^*, u^*})$ , where  $G_{i^*} \cap G_{i_l} = \emptyset$ .

Thus, the genre-aware contrastive objective  $L_l^a(g_c)$  of document  $x_{i_l, u_l}^a$  could be computed with:

$$L_l^a(g_c) = -\log \frac{e^{s(\vec{\theta}_{i_l, u_l}^a, \vec{\theta}_{i_l, u_l}^b)/\tau_g}}{\sum_{j=1}^M [I_g(l, j) \cdot e^{s(\vec{\theta}_{i_l, u_l}^a, \vec{\theta}_j^a)/\tau_g} + I_g(l, j) \cdot e^{s(\vec{\theta}_{i_l, u_l}^a, \vec{\theta}_j^b)/\tau_g}]} \quad (7)$$

where  $s(\cdot, \cdot)$  means cosine similarity,  $\tau_g$  is the temperature parameter of genre-aware contrastive learning,  $\vec{\theta}_l^a$  is abbreviated to  $\vec{\theta}_{i_l, u_l}^a$  which means the inferred document-topic distribution of  $x_{i_l, u_l}^a$ , and  $I_g(l, j)$  is the genre indicator function defined as:

$$I_g(l, j) = I_g(i_l, i_j) = \begin{cases} 1, & G_{i_l} \cap G_{i_j} = \emptyset \\ 0, & \text{otherwise} \end{cases} \quad (8)$$

Generally, the genre-aware contrastive objective of the augmented document set  $X^*$  could be calculated with:

$$\mathcal{L}_{G_c} = \frac{1}{2M} \sum_{l=1}^M [L_l^a(g_c) + L_l^b(g_c)] \quad (9)$$

## 4.3 Preference Inference and User Contrastive

Aiming at inferring user preference distributions over topics, GPTM utilizes a user-inference network  $\mathcal{I}_u$  to capture their interests. For each document  $x_{i,u}$  posted by user  $u$ , the user preference distribution  $\vec{p}_u$  could be inferred with:

$$\vec{p}_u = \mathcal{I}_u(\text{one-hot}[u]) \quad (10)$$

where  $\text{one-hot}[u]$  represents the one-hot encoding of user  $u$ , preference distributions  $\vec{p}_u^a$  and  $\vec{p}_u^b$  of augmented documents ( $x_{i,u}^a$  and  $x_{i,u}^b$ ) are equal to  $\vec{p}_u$ . Furthermore, the personalized document-topic distribution  $\vec{\theta}_{i,u}^a$  of augmented document  $x_{i,u}^a$  could be computed with formulation:

$$\vec{\theta}_{i,u}^a = \text{Multinomial}(\vec{p}_u^a \odot \vec{\theta}_{i,u}^a) \quad (11)$$

where  $\odot$  represents the element-wise product, and  $\vec{\theta}_{i,u}^b$  could be generated via Eq. 11.

To enable  $\mathcal{I}_u$  to capture the user's interests, GPTM employs user-aware contrastive learning to help  $\mathcal{I}_u$  distinguish the preferences of different users. In detail, given a set of  $2M$  augmented documents  $X^*$ , for each document  $x_{i_l, u_l}^a$ , we select the matched  $x_{i_l, u_l}^b$  to build the positive pair  $(x_{i_l, u_l}^a, x_{i_l, u_l}^b)$ . Contrarily, we only select the document  $x_{i^+, u^+}$ , posted by other users to a different genre item, in  $X^*$  to construct negative pairs  $(x_{i_l, u_l}^a, x_{i^+, u^+})$ . Here,  $G_{i^+} \cap G_{i_l} = \emptyset$  and  $u^+ \neq u_l$ .

Thus, the user-aware contrastive objective  $L_l^a(u_c)$  of augmented document  $x_{i_l, u_l}^a$  could be calculated with:

$$L_l^a(u_c) = -\log \frac{e^{s(\vec{\theta}_{i_l, u_l}^a, \vec{\theta}_{i_l, u_l}^b)/\tau_u}}{\sum_{j=1}^M [I_u(l, j) \cdot e^{s(\vec{\theta}_{i_l, u_l}^a, \vec{\theta}_j^a)/\tau_u} + I_u(l, j) \cdot e^{s(\vec{\theta}_{i_l, u_l}^a, \vec{\theta}_j^b)/\tau_u}]} \quad (12)$$

where  $\tau_u$  denotes the temperature parameter of user-aware contrastive learning,  $\vec{\theta}_l^a$  is abbreviated to the personalized document-topic distribution  $\vec{\theta}_{i_l, u_l}^a$ , computed with Eq. 11, and  $I_u(l, j)$  is the user indicator function defined with:

$$I_u(l, j) = I_u(i_l, u_l, i_j, u_j) = \begin{cases} 1, & G_{i_l} \cap G_{i_j} = \emptyset \text{ and } u_l \neq u_j \\ 0, & \text{otherwise} \end{cases} \quad (13)$$

Likewise, the user-aware contrastive objective of augmented set  $X^*$  could be computed with:

$$\mathcal{L}_{U_c} = \frac{1}{2M} \sum_{l=1}^M [L_l^a(u_c) + L_l^b(u_c)] \quad (14)$$



#### 4.4 Dirichlet Prior Matching

As Wallach et al. [38] argue that modeling topics with Dirichlet distribution helps to capture multiple patterns in texts and ensure interpretability, GPTM employs the Maximum Mean Discrepancy (MMD) [11] to match inferred document-topic distributions to the Dirichlet prior during genre-aware contrastive learning procedure.

Concretely, given a set of  $2M$  inferred document-topic distributions  $\Theta = \Theta^a \cup \Theta^b = \{\vec{\theta}_1, \vec{\theta}_2, \dots, \vec{\theta}_{2M-1}, \vec{\theta}_{2M}\}$ <sup>3</sup> and a randomly sampled set  $\Theta'' = \{\vec{\theta}''_1, \vec{\theta}''_2, \dots, \vec{\theta}''_{2M}\}$  from the  $Dir(\vec{\theta}''|\vec{\alpha})$ , we follow [11] to estimate the genre-aware matching objective  $\mathcal{L}_{M_g}$  with:

$$\mathcal{L}_{M_g} = \frac{1}{2M(2M-1)} \sum_{p \neq q} [k(\vec{\theta}_p, \vec{\theta}_q) + k(\vec{\theta}''_p, \vec{\theta}''_q)] - \frac{1}{2M^2} \sum_{p, q} k(\vec{\theta}_p, \vec{\theta}''_q) \quad (15)$$

where  $M$  is batch size. Likewise, in user-aware contrastive learning, the inferred personalized document-topic distributions  $\Theta' = \Theta'^a \cup \Theta'^b = \{\vec{\theta}'_1, \vec{\theta}'_2, \dots, \vec{\theta}'_{2M-1}, \vec{\theta}'_{2M}\}$ <sup>4</sup> are matched to  $Dir(\vec{\theta}'|\vec{\alpha})$ . The user-aware matching objective  $\mathcal{L}_{M_u}$  is computed with:

$$\mathcal{L}_{M_u} = \frac{1}{2M(2M-1)} \sum_{p \neq q} [k(\vec{\theta}'_p, \vec{\theta}'_q) + k(\vec{\theta}''_p, \vec{\theta}''_q)] - \frac{1}{2M^2} \sum_{p, q} k(\vec{\theta}'_p, \vec{\theta}''_q) \quad (16)$$

where  $k(\cdot, \cdot)$  means kernel function, we follow Nan et al. [28] and utilize the diffusion kernel [17].

#### 4.5 Training Objective

As mentioned, the key aims of our proposed GPTM are:

- (1) Mining topics relevant to genres with Genre-aware Topic Modeling;
- (2) Mining user preference and producing personalized topic profiles for each user with User Preference Modeling.

In genre topic modeling, we employ genre-aware contrastive learning to help topic-inference network  $\mathcal{I}_t$  capture word patterns of genres. Meanwhile, we also match inferred document-topic distributions to the Dirichlet prior for improving topic interpretability. Thus, the training objective of genre-aware topic modeling could be formulated as:

$$\mathcal{L}_{TM} = \mathcal{L}_{G_c} + \lambda_1 \mathcal{L}_{M_g} \quad (17)$$

where  $\mathcal{L}_{G_c}$  is the genre-aware contrastive objective, computed with Eq. 9. And the genre-aware matching objective  $\mathcal{L}_{M_g}$  is calculated with Eq. 15,  $\lambda_1$  is the coefficient hyper-parameter.

On the other hand, to meet the requirement of user preference modeling, we conduct user-aware contrastive learning to guide the user-inference network  $\mathcal{I}_u$  to discover the user's interests. Likewise, we match personalized document-topic distributions to the Dirichlet distribution to ensure the interpretability of personalized topic profiles. And the training objective of user preference modeling  $\mathcal{L}_{UM}$  could be formulated as:

$$\mathcal{L}_{UM} = \mathcal{L}_{U_c} + \lambda_2 \mathcal{L}_{M_u} \quad (18)$$

Here,  $\mathcal{L}_{U_c}$  and  $\mathcal{L}_{M_u}$  are user-aware contrastive and matching objectives, which could be computed with Eq. 14 and Eq. 16. And  $\lambda_2$  is the coefficient hyper-parameter. Detailed training procedure of GPTM and hyper-parameter settings please refer to Appendix A.1.

<sup>3</sup> $\Theta$  equals to  $\{\vec{\theta}_{i_1, u_1}^a, \vec{\theta}_{i_1, u_1}^b, \dots, \vec{\theta}_{i_M, u_M}^a, \vec{\theta}_{i_M, u_M}^b\}$ , we modify notations for simplicity.

<sup>4</sup> $\Theta'$  equals to  $\{\vec{\theta}'_{i_1, u_1}, \vec{\theta}'_{i_1, u_1}, \dots, \vec{\theta}'_{i_M, u_M}, \vec{\theta}'_{i_M, u_M}\}$

#### 4.6 Topic Generation

As learned inference networks  $\mathcal{I}_t$  (topic) and  $\mathcal{I}_u$  (user) build projections from the shared word/text representations and the user spaces to the latent topic space, they could be utilized to mine genre-aware topics and personalized topic profiles for users.

**4.6.1 Genre-aware Topic Generation.** For the  $v$ -th ( $v \in \{1, 2, \dots, V\}$ ) word in the vocabulary, we first collect its contextualized word representations of  $N^v$  appearance in the corpus with an embedding matrix  $E^v \in \mathbb{R}^{H_v \times N^v}$ . Then, the semantic correlation between the  $v$ -th word and topics could be computed with:

$$\vec{c}^v = \text{avg}_c(\mathcal{I}_t(E^v)) \quad (19)$$

where  $\text{avg}_c(\cdot)$  denotes column-wise averaging. Similarly, the semantic correlation matrix  $C \in \mathbb{R}^{K \times V}$  between words and topics could be calculated. And the topic-word distribution matrix  $\Phi \in \mathbb{R}^{K \times V}$  could be obtained with:

$$\Phi = \text{norm}_c(C) \quad (20)$$

where  $\text{norm}_c(\cdot)$  means column-wise normalization, and the  $k$ -th row  $\vec{\phi}_k$  is the word distribution of the  $k$ -th topic.

**4.6.2 Personalized Topic Generation.** For each user  $u$ , the learned user-inference network  $\mathcal{I}_u$  could produce his/her preference distribution  $\vec{p}_u$  over topics with Eq.10. Together with mined topic-distributions in  $\Phi$ , the personalized topic profile  $\vec{\phi}_u$  of user  $u$  could be computed with:

$$\vec{\phi}_u = \sum_{k=1}^K p_u^k \cdot \vec{\phi}_k \quad (21)$$

where  $p_u^k$  is the  $k$ -th dimension of  $\vec{p}_u$ , and it indicates user  $u$ 's preference to the  $k$ -th topic.

### 5 EXPERIMENTS

In this section, we first describe the experimental setup, which contains descriptions of datasets, evaluation metrics and baselines. Then, we provide comparison results and discussions of genre-aware topic modeling and user preference modeling. Following this, the ablation study will be presented lastly.

#### 5.1 Experimental Setup

**5.1.1 Datasets.** To verify the effectiveness of GPTM on genre-aware topic modeling and user preference modeling, three Amazon review datasets ('Books'<sup>5</sup>, 'Sports'<sup>6</sup> and 'Movies'<sup>7</sup>) are utilized. For dataset construction, we discard reviews of product items with low frequency and similar genres, and we only retain users who posted more than 50/75/200 comments for Movies/Sports/Books datasets. Besides, we conduct pre-processing like lemmatization and spell-checking with spaCy<sup>8</sup>. Moreover, special characters, certain punctuations and reviews fewer than 15 words are omitted. The statistics of processed datasets are presented in Table. 2.

<sup>5</sup>[https://huggingface.co/datasets/McAuley-Lab/Amazon-Reviews-2023/blob/main/raw/review\\_categories/Books.jsonl](https://huggingface.co/datasets/McAuley-Lab/Amazon-Reviews-2023/blob/main/raw/review_categories/Books.jsonl)

<sup>6</sup>[https://huggingface.co/datasets/McAuley-Lab/Amazon-Reviews-2023/blob/main/raw/review\\_categories/Sports\\_and\\_Outdoors.jsonl](https://huggingface.co/datasets/McAuley-Lab/Amazon-Reviews-2023/blob/main/raw/review_categories/Sports_and_Outdoors.jsonl)

<sup>7</sup>[https://huggingface.co/datasets/McAuley-Lab/Amazon-Reviews-2023/blob/main/raw/review\\_categories/Movies\\_and\\_TV.jsonl](https://huggingface.co/datasets/McAuley-Lab/Amazon-Reviews-2023/blob/main/raw/review_categories/Movies_and_TV.jsonl)

<sup>8</sup><https://spacy.io/>

**Table 2: The statistics of processed datasets.**

Dataset	# Doc	# Genre	# User	# Vocab
Books	11,116	31	30	23,001
Sports	9,374	25	98	10,626
Movies	4,530	20	97	12,670

5.1.2 *Topic Evaluation Metrics.* We choose four widely utilized coherence metrics ( $C_P$ ,  $C_A$ ,  $NPMI$  and  $UCI$ ) [30], computed with the Palmetto<sup>9</sup> library, and Unique Term (UT) [40] to assess semantic quality and diversity of topics.

Besides, to evaluate the correlations between mined topics and genres, we design the Genre Discovered Rate ( $GDR$ ) and the Genre Topic Correlation ( $GTC$ ) metrics based on OpenAI embeddings<sup>10</sup>. Specifically, the  $GDR$  value could be computed with:

$$GDR = \frac{1}{K} \sum_{k=1}^K \max(0, \text{sgn}(\cos^*(g_k, \Phi) - \sigma)) \quad (22)$$

where  $\sigma$  is the threshold hyper-parameter,  $\text{sgn}(\cdot)$  is the sign function which outputs 1 for positive input,  $\cos^*(g_k, \Phi)$  is the maximum cosine similarity between OpenAI embeddings of the  $k$ -th genre  $g_k$  and topic words. And the  $GTC$  follows the computation:

$$GTC = \frac{1}{K} \sum_{k=1}^K \cos^*(\vec{\phi}_k, G) \quad (23)$$

where  $G$  denotes the genre set, and  $\cos^*(\vec{\phi}_k, G)$  is the maximum cosine similarity between OpenAI embeddings of the  $k$ -th topic and genres set  $G$ . Here, we use the top 10 words to represent the topic, and higher values indicate a higher quality of extracted topics.

5.1.3 *User Preference Evaluation Metrics.* To assess user preference modeling performance, we design the Personalized Hit Rate ( $PHR$ ) and Personalized Genre Correlation ( $PGC$ ) metrics. Concretely, the  $PHR$  value could be calculated with:

$$PHR = \frac{1}{N_u} \sum_{u=1}^{N_u} \frac{\text{count}(\vec{\phi}_u^{20} \cap W_u)}{20} \quad (24)$$

where  $N_u$  means the number of users,  $W_u$  denotes a set of words posted by user  $u$ ,  $\vec{\phi}_u^{20}$  means top 20 words obtained from personalized topic profile  $\vec{\phi}_u$ . And  $\text{count}(\cdot)$  is the counting function which returns the size of the input word set. The  $PGC$  value is defined as:

$$PGC = \frac{1}{N_u} \sum_{u=1}^{N_u} \sum_{k=1}^K \hat{p}_u^k \cdot \cos(\vec{\phi}_u^{20}, g_k) \quad (25)$$

where  $\hat{p}_u^k$  is the proportion of reviews about  $g_k$  posted by user  $u$ , which is computed with  $\#D_u^k / \#D_u$ . Here,  $\#D_k$  ( $\#D_k^u$ ) is the number of reviews posted by  $u$  (related to genre  $g_k$ ). And  $\cos(\vec{\phi}_u^{20}, g_k)$  is the cosine similarity between top 20 words of  $\vec{\phi}_u$  and genre  $g_k$ .

5.1.4 *Baselines.* We choose the following approaches as baselines:

- **LDA** [2], is the most widely used topic model, which views document is generated from a mixture of topics<sup>11</sup>.
- **ASTM** [23], is an autoencoding-based approach which utilizes sinkhorn distance to match topic distribution<sup>12</sup>.
- **CTMNeg** [1], is a neural topic model based on contrastive learning and negative sampling<sup>13</sup>.

<sup>9</sup><https://github.com/dice-group/Palmetto>

<sup>10</sup><https://platform.openai.com/docs/guides/embeddings>

<sup>11</sup><https://mimno.github.io/Mallet/>

<sup>12</sup><https://github.com/AEGISEdge/ASTM>

<sup>13</sup><https://github.com/adhyasuman/ctmneg>

- **BERTopic** [12], is a topic mining approach based on contextualized text representation and clustering<sup>14</sup>.
- **ECRTM** [41], is a neural topic model that regularizes topics with the optimal transport distance between topic embeddings and centers of word embedding clusters<sup>15</sup>.
- **vONT** [43], is a VAE-based approach which models topics with a mixture of von Mises-Fisher distributions<sup>16</sup>.
- **BertSenClu** [31], is a topic mining approach based on the Bag-of-Sentences (BoS) assumption<sup>17</sup>.
- **CWTM** [7], is a neural topic model which incorporates contextualized word representations for topic inference<sup>18</sup>.
- **NPTM** [24], is the first personalized neural topic model which uses a hybrid generative process to combine user preferences and contextualized document representations<sup>19</sup>.

## 5.2 Genre-aware Topic Evaluation

To validate topic modeling performance of GPTM, we conduct experiments with two different topic numbers for each dataset. In small topic number settings (15 topics for *Books*, 15 topics for *Sports* and 10 topics for *Movies*), we only utilize reviews associated with selected 15/15/10 genres for *Books/Sports/Movies* dataset.

The comparison results on four topic coherence metrics ( $C_P$ ,  $C_A$ ,  $NPMI$  and  $UCI$ ) and topic diversity metric  $UT$  are presented in Table 3. From statistics, we could observe that the proposed GPTM outperforms all the baseline approaches for all three datasets on all five metrics. This may be attributed to the following factors: 1). Injecting contextualized information among texts into the modeling process helps to improve topic quality; 2). The Dirichlet prior is suitable for modeling topics and enhancing interpretability; 3). Contrastive learning could alleviate mode collapse and result in a higher-level topic diversity. The detailed comparison of hyper-parameter analysis please refer to Appendix A.2.

Besides, to assess the correlations between mined topics and genres, we compute the Genre Discovered Rate ( $GDR$ ) and Genre Topic Correlation ( $GTC$ ) values, with the topic number set to 31 for *Books*, 25 for *Sports* and 20 for *Movies*. For the  $GDR$  metric, we set the threshold parameter  $\sigma$  to 0.4, 0.45, 0.5, and 0.55 respectively. And the detailed comparative results are shown in Figure 3. We could observe that our proposed GPTM outperforms almost all the baselines. Additionally, with the increase in  $\sigma$ , GPTM's performance drops more slowly than the compared approaches. For the  $GTC$  metric, we present the comparison results in Table 4 which also reveal the superior performance of GPTM. And these improvements on  $GDR$  and  $GTC$  may be attributed to factors: 1). Incorporating genre description information helps GPTM ensure relevance between topics and genres. 2). Genre-aware contrastive learning in GPTM helps to capture diverse genre-relevant topics, resulting in higher values on the  $GDR$  metric. We also provide several topic examples in Appendix A.3 for comparison.

<sup>14</sup><https://github.com/MaartenGr/BERTopic>

<sup>15</sup><https://github.com/BobXWu/ECRTM>

<sup>16</sup><https://github.com/xuweijieshuai/Neural-Topic-Modeling-vmf>

<sup>17</sup><https://github.com/JohnTailor/BertSenClu>

<sup>18</sup><https://github.com/Fitz-like-coding/CWTM>

<sup>19</sup><https://github.com/AEGISEdge/NPTM>

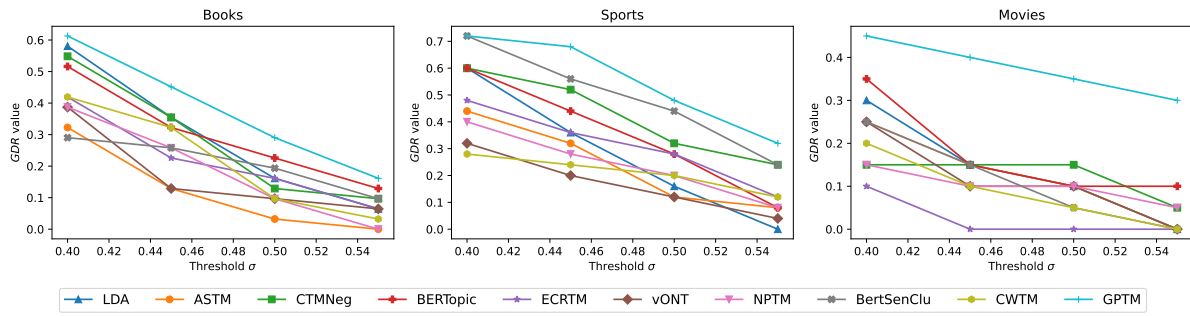


Figure 3: Comparison results on the Genre Discovered Rate (GDR) metric.

Table 3: The comparison on Coherence and Diversity metrics.

Dataset	Books									
	K = 15					K = 31				
# Topic	$C_P$	$C_A$	$NP MI$	$UCI$	$UT$	$C_P$	$C_A$	$NP MI$	$UCI$	$UT$
LDA	0.311	0.190	0.066	0.662	0.853	0.256	0.172	0.044	0.292	0.797
ASTM	0.127	0.183	0.024	-0.445	0.693	0.045	0.173	0.006	-0.755	0.677
CTMNeg	0.287	0.196	0.049	0.077	0.967	0.313	0.192	0.057	0.280	0.906
BERTopic	0.355	0.253	0.069	0.230	0.840	0.346	0.230	0.065	0.165	0.842
ECRTM	0.408	0.230	0.077	0.487	0.980	0.099	0.190	0.034	-0.427	0.977
vONT	0.230	0.150	0.040	0.250	0.700	0.168	0.132	0.025	0.150	0.713
NPTM	0.309	0.232	0.082	0.232	0.753	0.251	0.209	0.065	-0.092	0.794
BertSenClu	-0.100	0.163	-0.030	-2.017	0.987	-0.221	0.138	-0.044	-2.165	0.990
CWTM	0.446	0.221	0.073	0.258	0.913	0.360	0.196	0.062	0.286	0.835
GPTM	<b>0.517</b>	<b>0.268</b>	<b>0.103</b>	<b>0.803</b>	<b>1.000</b>	<b>0.492</b>	<b>0.269</b>	<b>0.086</b>	<b>0.293</b>	<b>1.000</b>

Dataset	Sports									
	K = 15					K = 25				
# Topic	$C_P$	$C_A$	$NP MI$	$UCI$	$UT$	$C_P$	$C_A$	$NP MI$	$UCI$	$UT$
LDA	0.328	0.200	0.039	-0.229	0.853	0.288	0.188	0.029	-0.307	0.760
ASTM	0.224	0.190	0.010	-0.880	0.853	0.214	0.196	0.023	-0.575	0.812
CTMNeg	0.131	0.167	0.000	-0.921	0.893	0.163	0.179	0.021	-0.492	0.864
BERTopic	0.279	0.189	0.010	-1.016	0.880	0.264	0.197	0.012	-0.890	0.860
ECRTM	0.228	0.180	0.007	-1.000	0.687	0.079	0.167	-0.015	-1.520	0.912
vONT	0.193	0.153	-0.013	-1.084	0.727	0.171	0.145	0.001	-0.565	0.800
NPTM	0.141	0.169	0.009	-1.040	0.960	0.169	0.197	0.027	-0.824	0.868
BertSenClu	0.023	0.160	-0.043	-2.292	0.993	-0.117	0.150	-0.056	-2.479	0.996
CWTM	0.323	0.199	0.028	-0.785	0.853	0.320	0.207	0.028	-0.558	0.860
GPTM	<b>0.476</b>	<b>0.241</b>	<b>0.062</b>	<b>-0.218</b>	<b>1.000</b>	<b>0.407</b>	<b>0.253</b>	<b>0.059</b>	<b>-0.285</b>	<b>1.000</b>

Dataset	Movies									
	K = 10					K = 20				
# Topic	$C_P$	$C_A$	$NP MI$	$UCI$	$UT$	$C_P$	$C_A$	$NP MI$	$UCI$	$UT$
LDA	0.253	0.149	0.036	0.088	0.870	0.180	0.140	0.020	-0.131	0.800
ASTM	0.272	0.195	0.045	0.150	0.810	0.207	0.194	0.042	-0.052	0.710
CTMNeg	-0.014	0.134	-0.020	-1.410	0.930	0.014	0.135	-0.012	-0.978	0.795
BERTopic	0.265	0.209	0.048	0.198	0.730	0.266	0.199	0.031	-0.344	0.770
ECRTM	0.037	0.134	-0.004	-0.930	0.780	-0.276	0.117	-0.063	-2.305	0.925
vONT	0.195	0.140	0.030	0.197	0.650	0.141	0.132	0.011	-0.151	0.675
NPTM	0.068	0.151	-0.013	-1.382	0.950	-0.232	0.113	-0.044	-2.080	0.935
BertSenClu	-0.412	0.100	-0.139	-4.157	0.990	-0.647	0.084	-0.145	-4.173	0.980
CWTM	0.181	0.146	0.015	-0.092	0.640	0.215	0.145	0.018	-0.303	0.690
GPTM	<b>0.459</b>	<b>0.241</b>	<b>0.076</b>	<b>0.224</b>	<b>1.000</b>	<b>0.358</b>	<b>0.216</b>	<b>0.060</b>	<b>0.186</b>	<b>1.000</b>

Table 4: Comparative results on the GTC metric.

Model	LDA	ASTM	CTMNeg	BERTopic	ECRTM	vONT	NPTM	BertSenClu	CWTM	GPTM
Books	0.44	0.37	0.43	0.42	0.39	0.40	0.35	0.38	0.40	<b>0.45</b>
Sports	0.43	0.38	0.42	0.41	0.41	0.39	0.36	0.39	0.43	<b>0.45</b>
Movies	0.40	0.31	0.37	0.38	0.33	0.36	0.36	0.33	0.37	<b>0.41</b>

### 5.3 User Preference Evaluation

To evaluate the performance of user preference modeling, we utilize the Personalized Hit Rate (PHR) and Personalized Genre Correlation (PGC) metrics, computed with Eq. 24 and Eq. 25, which reflect

Table 5: Comparative results on PHR and PGC metrics.

Dataset	Books				Sports				Movies			
	K = 15		K = 31		K = 15		K = 25		K = 10		K = 20	
# Topic	PHR	PGC	PHR	PGC	PHR	PGC	PHR	PGC	PHR	PGC	PHR	PGC
NPTM	0.60	0.59	0.27	0.58	0.25	0.51	0.24	0.47	0.20	0.53	0.19	0.58
GPTM	<b>0.72</b>	<b>0.65</b>	<b>0.88</b>	<b>0.60</b>	<b>0.62</b>	<b>0.60</b>	<b>0.60</b>	<b>0.60</b>	<b>0.38</b>	<b>0.57</b>	<b>0.43</b>	<b>0.59</b>

the agreement between mined personalized topic profiles and users' interests.

The corresponding results are listed in Table 5. We only present NPTM's results for comparison as other baselines could not provide personalized topic profiles. The PHR and PGC scores in Table 5 indicate that the personalized topic profiles generated by GPTM share a higher level of agreement with the user's interests. And such improvements may be caused by factors: 1). Genre-aware contrastive learning helps GPTM to mine genre-relevant topics; 2). The user-inference network  $I_u$ , trained with user-aware contrastive learning, could build preference distribution for users accurately. We also present several personalized topic profiles in Appendix A.3, generated by GPTM and NPTM, for comparison.

**5.3.1 User Preference Visualization.** To provide a more direct comparison of user preference modeling between GPTM and NPTM, we conduct a user-preference visualization experiment.

In Figure 4, subplots in each line represent the comparison on a dataset. For each dataset, we randomly choose six users and generate their preference distributions. Meanwhile, we also infer document-topic distributions for reviews posted by selected users. Also, t-SNE [36] is leveraged for dimension reduction, embeddings associated with different users are painted with different colors. Dots denote inferred document-topic distributions of reviews, and user preference distributions are represented with markers with distinctive shapes.

From Figure 4, we could observe the following findings: 1). GPTM could produce distinguishable document-topic distributions for different users, while NPTM entangles the document-topic distributions generated by different users together. 2). User-preference distributions inferred by GPTM are located close to the document-topic distributions which means GPTM could produce more accurate preference distributions than NPTM. These two observations indicate that the proposed GPTM exhibits competitive ability in user preference modeling. In more detail, for the Sport dataset in

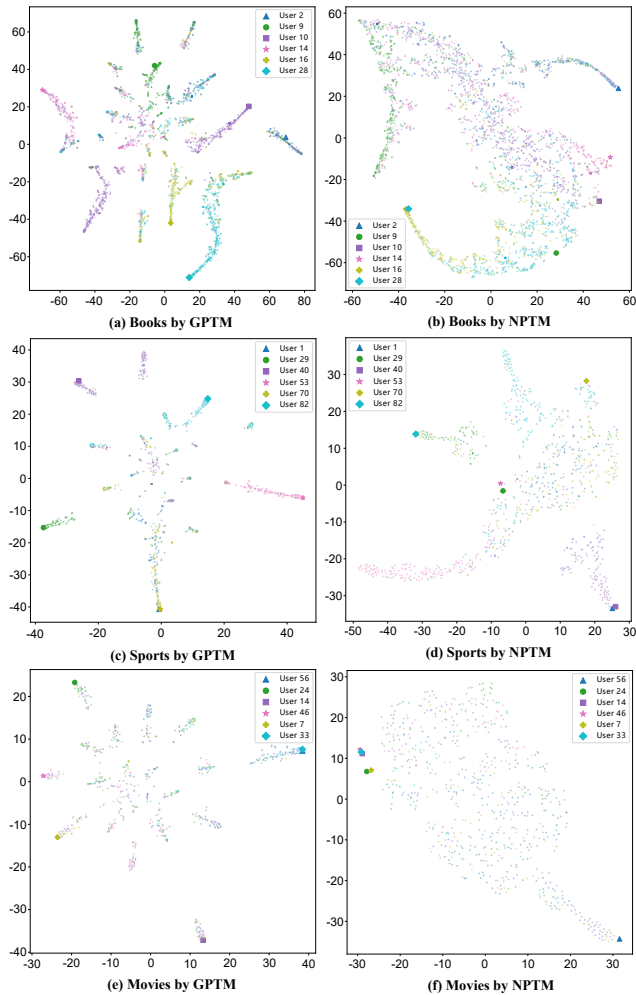


Figure 4: Visualization of user-preference distributions and document-topic distributions provided by GPTM and NPTM.

subplot (c), the blue triangle (user 1) and the yellow plus (user 70) are located together, which is because they have similar preferences. A similar case also can be found (user 56 and user 33) in subplot (e).

**5.3.2 User Preference Modeling Dynamics.** To exhibit the dynamic process of user preference modeling, we visualize the preference distributions of users associated with ‘*Hunting & Fishing*’ and ‘*Cycling*’ genres at different stages of user-aware contrastive learning in Figure 5.

For each subplot, the horizontal axis represents the user’s preference for ‘*Hunting & Fishing*’, while the vertical axis represents ‘*Cycling*’, and each point denotes a user associated with these two genres<sup>20</sup>. At the early stage of modeling, we observe that users are entangled together. However, as the training process to 20 epochs, users could be separated by GPTM according to their preferences. Lastly, GPTM predicts user preferences with higher confidence, with more colored points located around (0,1) and (1,0).

<sup>20</sup>We label a user with ‘*Cycling*’ if most of her/his reviews are tagged with ‘*Cycling*’.

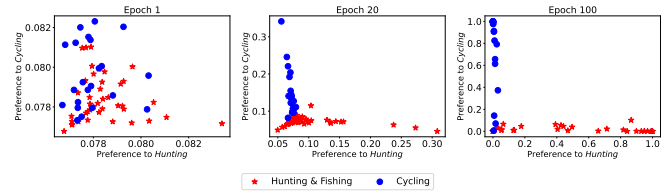


Figure 5: Dynamic process of user preference modeling on ‘*Hunting & Fishing*’ and ‘*Cycling*’ genres in *Sports* dataset.

Table 6: Comparison of ablation on transformer variants.

Dataset	Model	$C_P$	$C_A$	$NPMI$	$UCI$	$UT$
Books	NPTM	0.309	0.232	0.082	0.232	0.753
	CWTM	0.446	0.221	0.073	0.258	0.913
	GPTM-Bert	0.478	0.275	0.090	0.560	<b>1.000</b>
	GPTM-RoBERTa	0.515	0.278	<b>0.109</b>	<b>0.922</b>	<b>1.000</b>
	GPTM-SimCSE	0.512	<b>0.280</b>	0.090	0.388	<b>1.000</b>
	GPTM-FLAN-T5	0.455	0.278	0.085	0.599	<b>1.000</b>
GPTM-Sentence-Bert	<b>0.517</b>	0.268	0.103	0.803	<b>1.000</b>	

## 5.4 Ablation Study with Transformer variants

To explore the impact of the transformer language model on GPTM, we conduct an ablation study in the experiment on the *Books* dataset with 15 topic settings.

Specifically, we examine five transformer variants which are Bert [6]<sup>21</sup>, RoBERTa [25]<sup>22</sup>, SimCSE [9]<sup>23</sup>, FLAN-T5 [4]<sup>24</sup> and Sentence-Bert [29]<sup>25</sup>. And the detailed results are presented in Table 6. Optimal results are marked in bold, we only list the values generated by CWTM and NPTM due to their superior performance among baselines. Statistics show that GPTM could consistently produce competitive results with different transformers considering coherence and diversity metrics.

## 6 CONCLUSION AND FUTURE WORK

In this paper, we have proposed a novel Genre-aware Personalized neural Topic Model (GPTM) to mine genre-aware topics from Amazon reviews and conduct user preference modeling. GPTM incorporates genre information into the topic modeling process to ensure correspondence between topics and product genres. Besides, it employs user-aware contrastive learning to help the user-inference network to capture users’ preferences. The experimental comparisons on three review datasets with state-of-the-art baselines show that GPTM achieves improved coherence and diversity while maintaining stronger correlations between topics and genres. Moreover, it also surpasses baselines in user-preference modeling and could produce high-quality personalized topic profiles. In the future, we will explore leveraging large language models for user-preference modeling and generating personalized topic profiles.

<sup>21</sup><https://huggingface.co/bert-base-uncased>

<sup>22</sup><https://huggingface.co/sentence-transformers/all-roberta-large-v1>

<sup>23</sup><https://huggingface.co/princeton-nlp/sup-simcse-roberta-large>

<sup>24</sup><https://huggingface.co/google/flan-t5-base>

<sup>25</sup><https://huggingface.co/sentence-transformers/all-mpnet-base-v2>



## References

- [1] Suman Adhya, Avishek Lahiri, Debarshi Kumar Sanyal, and Partha Pratim Das. 2022. Improving Contextualized Topic Models with Negative Sampling. In *Proceedings of the 19th International Conference on Natural Language Processing, ICON 2022, New Delhi, India, December 15-18, 2022*. Association for Computational Linguistics, 128–138.
- [2] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent Dirichlet Allocation. *J. Mach. Learn. Res.* 3 (2003), 993–1022.
- [3] Sihao Chen, Hongming Zhang, Tong Chen, Ben Zhou, Wenhao Yu, Dian Yu, Baolin Peng, Hongwei Wang, Dan Roth, and Dong Yu. 2024. Sub-Sentence Encoder: Contrastive Learning of Propositional Semantic Representations. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, NAACL 2024, Mexico City, Mexico, June 16-21, 2024. Association for Computational Linguistics, 1596–1609.
- [4] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling Instruction-Finetuned Language Models. *CoRR abs/2210.11416* (2022). arXiv:2210.11416
- [5] Elijah Cole, Xuan Yang, Kimberly Wilber, Oisín Mac Aodha, and Serge J. Beaulieu. 2022. When Does Contrastive Visual Representation Learning Work?. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 1–10.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, 4171–4186.
- [7] Zheng Fang, Yulan He, and Rob Procter. 2024. CWTM: Leveraging Contextualized Word Embeddings from BERT for Neural Topic Modeling. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy*. ELRA and ICCL, 4273–4286.
- [8] Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. 2021. A Survey of Data Augmentation Approaches for NLP. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 968–988.
- [9] Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple Contrastive Learning of Sentence Embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*. Association for Computational Linguistics, 6894–6910.
- [10] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. 2014. Generative Adversarial Nets. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*. 2672–2680.
- [11] Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander J. Smola. 2012. A Kernel Two-Sample Test. *J. Mach. Learn. Res.* 13 (2012), 723–773.
- [12] Maarten Grootendorst. 2022. BERTopic: Neural topic modeling with a class-based TF-IDF procedure. *CoRR abs/2203.05794* (2022). arXiv:2203.05794
- [13] Pankaj Gupta, Yatin Chaudhary, Thomas A. Runkler, and Hinrich Schütze. 2020. Neural Topic Modeling with Continual Lifelong Learning. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event (Proceedings of Machine Learning Research, Vol. 119)*. PMLR, 3907–3917.
- [14] Xuemeng Hu, Rui Wang, Deyu Zhou, and Yuxuan Xiong. 2020. Neural Topic Modeling with Cycle-Consistent Adversarial Training. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*. Association for Computational Linguistics, 9018–9030.
- [15] Diederik P Kingma. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [16] Diederik P. Kingma and Max Welling. 2014. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*.
- [17] John Lafferty, Guy Lebanon, and Tommi Jaakkola. 2005. Diffusion kernels on statistical manifolds. *Journal of Machine Learning Research* 6, 1 (2005).
- [18] Kwot Sin Lee, Ngoc-Trung Tran, and Ngai-Man Cheung. 2021. InfoMax-GAN: Improved Adversarial Image Generation via Information Maximization and Contrastive Learning. In *IEEE Winter Conference on Applications of Computer Vision, WACV 2021, Waikoloa, HI, USA, January 3-8, 2021*. IEEE, 3941–3951.
- [19] Yunfan Li, Peng Hu, Jerry Zitao Liu, Dezhong Peng, Joey Tianyi Zhou, and Xi Peng. 2021. Contrastive Clustering. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*. AAAI Press, 8547–8555.
- [20] Yuheng Li, Yijun Li, Jingwan Lu, Eli Shechtman, Yong Jae Lee, and Krishna Kumar Singh. 2022. Contrastive Learning for Diverse Disentangled Foreground Generation. In *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XVI (Lecture Notes in Computer Science, Vol. 13676)*. Springer, 334–351.
- [21] Yunfan Li, Mouxiang Yang, Dezhong Peng, Taihao Li, Jiantao Huang, and Xi Peng. 2022. Twin Contrastive Learning for Online Clustering. *Int. J. Comput. Vis.* 130, 9 (2022), 2205–2221.
- [22] Yunlong Liang, Fandong Meng, Jiaan Wang, Jinan Xu, Yufeng Chen, and Jie Zhou. 2024. Continual Learning with Semi-supervised Contrastive Distillation for Incremental Neural Machine Translation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*. Association for Computational Linguistics, 10914–10928.
- [23] Luyang Liu, Heyan Huang, Yang Gao, and Yongfeng Zhang. 2022. Improving neural topic modeling via Sinkhorn divergence. *Inf. Process. Manag.* 59, 3 (2022), 102864.
- [24] Luyang Liu, Qunyang Lin, Haonan Tong, Hongyin Zhu, Ke Liu, Min Wang, and Chuang Zhang. 2023. Neural Personalized Topic Modeling for Mining User Preferences on Social Media. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, CIKM 2023, Birmingham, United Kingdom, October 21-25, 2023*. ACM, 1545–1555.
- [25] Yihan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *CoRR abs/1907.11692* (2019). arXiv:1907.11692
- [26] Sijie Mai, Ying Zeng, Shuangjia Zheng, and Haifeng Hu. 2023. Hybrid Contrastive Learning of Tri-Modal Representation for Multimodal Sentiment Analysis. *IEEE Trans. Affect. Comput.* 14, 3 (2023), 2276–2289.
- [27] Yishu Miao, Lei Yu, and Phil Blunsom. 2016. Neural Variational Inference for Text Processing. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016 (JMLR Workshop and Conference Proceedings, Vol. 48)*. JMLR.org, 1727–1736.
- [28] Feng Nan, Ran Ding, Ramesh Nallapati, and Bing Xiang. 2019. Topic Modeling with Wasserstein Autoencoders. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics.
- [29] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*. Association for Computational Linguistics, 3980–3990.
- [30] Michael Röder, Andreas Both, and Alexander Hinneburg. 2015. Exploring the Space of Topic Coherence Measures. In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining, WSDM 2015, Shanghai, China, February 2-6, 2015*. ACM, 399–408.
- [31] Johannes Schneider. 2024. Efficient and Flexible Topic Modeling Using Pretrained Embeddings and Bag of Sentences. In *Proceedings of the 16th International Conference on Agents and Artificial Intelligence, ICAART 2024, Volume 2, Rome, Italy, February 24-26, 2024*. SCITEPRESS, 407–418.
- [32] Vivek Sharma, Makarand Tapaswi, M. Saquib Sarfraz, and Rainer Stiefelhagen. 2020. Clustering based Contrastive Learning for Improving Face Representations. In *15th IEEE International Conference on Automatic Face and Gesture Recognition, FG 2020, Buenos Aires, Argentina, November 16-20, 2020*. IEEE, 109–116.
- [33] Tianbao Song, Jingbo Sun, Xin Liu, and Weiming Peng. 2024. Scale-VAE: Preventing Posterior Collapse in Variational Autoencoder. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy*. ELRA and ICCL, 14347–14357.
- [34] Akash Srivastava and Charles Sutton. 2017. Autoencoding Variational Inference For Topic Models. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. Open-Review.net.
- [35] Yixuan Su, Tian Lan, Yan Wang, Dani Yogatama, Lingpeng Kong, and Nigel Collier. 2022. A Contrastive Framework for Neural Text Generation. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- [36] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, 11 (2008).

- [37] Hanna M. Wallach, David M. Mimno, and Andrew McCallum. 2009. Rethinking LDA: Why Priors Matter. In *Advances in Neural Information Processing Systems 22: 23rd Annual Conference on Neural Information Processing Systems 2009. Proceedings of a meeting held 7-10 December 2009, Vancouver, British Columbia, Canada*. 1973–1981.
- [38] Hanna M. Wallach, David M. Mimno, and Andrew McCallum. 2009. Rethinking LDA: Why Priors Matter. In *Advances in Neural Information Processing Systems 22: 23rd Annual Conference on Neural Information Processing Systems 2009. Proceedings of a meeting held 7-10 December 2009, Vancouver, British Columbia, Canada*. Curran Associates, Inc., 1973–1981.
- [39] Rui Wang, Deyu Zhou, and Yulan He. 2019. ATM: Adversarial-neural Topic Model. *Inf. Process. Manag.* 56, 6 (2019).
- [40] Rui Wang, Deyu Zhou, Haiping Huang, and Yongquan Zhou. 2024. MIT: Mutual Information Topic Model for Diverse Topic Extraction. *IEEE Transactions on Neural Networks and Learning Systems* (2024).
- [41] Xiaobao Wu, Xinchuai Dong, Thong Thanh Nguyen, and Anh Tuan Luu. 2023. Effective Neural Topic Modeling with Embedding Clustering Regularization. In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA (Proceedings of Machine Learning Research, Vol. 202)*. PMLR, 37335–37357.
- [42] Sishi Xiong, Yu Zhao, Jie Zhang, Mengxiang Li, Zhongjiang He, Xuelong Li, and Shuangyong Song. 2024. Dual Prompt Tuning based Contrastive Learning for Hierarchical Text Classification. In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*. Association for Computational Linguistics, 12146–12158.
- [43] Weijie Xu, Xiaoyu Jiang, Srinivasan Sengamedu Hanumantha Rao, Francis Iannacci, and Jinjin Zhao. 2023. vONTSS: vMF based semi-supervised neural topic modeling with optimal transport. In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*. Association for Computational Linguistics, 4433–4457.
- [44] Yongcheng Zhang, Lingou Kong, Sheng Tian, Hao Fei, Changpeng Xiang, Huan Wang, and Xiaomei Wei. 2024. Multi-view Counterfactual Contrastive Learning for Fact-checking Fake News Detection. In *Proceedings of the 2024 International Conference on Multimedia Retrieval, ICMR 2024, Phuket, Thailand, June 10-14, 2024*. ACM, 385–393.
- [45] Huasong Zhong, Jianlong Wu, Chong Chen, Jianqiang Huang, Minghua Deng, Liqiang Nie, Zhouchen Lin, and Xian-Sheng Hua. 2021. Graph Contrastive Clustering. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, 9204–9213.

## A APPENDIX

### A.1 Training Procedure and Parameter Settings

The detailed training procedure of GPTM is shown in Algorithm 1. In the experiment, we set the learning rate  $\eta$  to 3e-5. The temperature parameters  $\tau_g$  and  $\tau_u$  are set to 0.5, hidden units  $H$  of inference networks  $\mathcal{I}_g$  and  $\mathcal{I}_u$  is set to 200, batch size  $M$  is set to 32, coefficient parameters  $\lambda_1$  and  $\lambda_2$  are set to 1.0, the hyper-parameter of Dirichlet prior  $\vec{\alpha}$  is set to 0.1. And our GPTM is optimized by Adam [15] optimizer.

### A.2 Hyper-Parameter Analysis

To investigate the impact of learning rate  $\eta$ , the number of hidden units  $H$ , and the Dirichlet prior  $\vec{\alpha}$  on the performance of GPTM, we conduct experiments on the *Books* dataset, with the topic number set to 15. Specifically,  $\eta$ ,  $H$ ,  $\vec{\alpha}$  are set to various values listed below:

- $\eta \in \{1e-5, 2e-5, 3e-5, 4e-5, 5e-5\}$ ;
- $H \in \{100, 150, 200, 250, 300\}$ ;
- $\vec{\alpha} \in \{0.08, 0.09, 0.10, 0.11, 0.12\}$

We present comparisons of  $\eta$ ,  $H$ ,  $\vec{\alpha}$  on coherence and diversity metrics in Table 7. Here, we only list results obtained by CWTM and NPTM, which outperform other baselines according to Table 3, for simplicity. The results reveal that GPTM could surpass the competitive CWTM with different combinations of hyper-parameters.

Table 7: Hyper-parameter analysis results on  $\eta$ ,  $H$  and  $\vec{\alpha}$ .

Parameter	Setting	$C_P$	$C_A$	$NPMI$	$UCI$	$UT$
$\eta$	NPTM	0.309	0.232	0.082	0.232	0.753
	CWTM	0.446	0.221	0.073	0.258	0.913
	1e-5	0.529	<b>0.282</b>	<b>0.105</b>	0.797	<b>1.000</b>
	2e-5	0.519	0.259	0.097	0.685	<b>1.000</b>
	3e-5	0.517	0.268	0.103	<b>0.803</b>	<b>1.000</b>
	4e-5	0.552	0.254	0.100	0.746	<b>1.000</b>
5e-5	<b>0.558</b>	0.265	0.102	0.789	<b>1.000</b>	
$H$	NPTM	0.309	0.232	0.082	0.232	0.753
	CWTM	0.446	0.221	0.073	0.258	0.913
	100	<b>0.562</b>	0.264	0.103	0.791	<b>1.000</b>
	150	0.512	0.268	<b>0.104</b>	0.836	<b>1.000</b>
	200	0.517	0.268	0.103	0.803	<b>1.000</b>
	250	0.517	<b>0.273</b>	0.101	0.765	<b>1.000</b>
300	0.540	0.257	<b>0.104</b>	<b>0.925</b>	<b>1.000</b>	
$\vec{\alpha}$	NPTM	0.309	0.232	0.082	0.232	0.753
	CWTM	0.446	0.221	0.073	0.258	0.913
	0.08	0.502	0.259	0.102	0.854	<b>1.000</b>
	0.09	0.521	0.260	0.100	0.759	<b>1.000</b>
	0.10	0.517	<b>0.268</b>	0.103	0.803	<b>1.000</b>
	0.11	0.515	0.262	0.101	0.828	<b>0.993</b>
0.12	<b>0.539</b>	0.262	<b>0.108</b>	<b>0.955</b>	<b>1.000</b>	

Algorithm 1: The training procedure of GPTM.

**Input:** Corpus  $D$ ; genre-inference network  $\mathcal{I}_g$ ; user-inference network  $\mathcal{I}_u$ ; batch size  $M$ ; genre contrastive epoch  $E_g$ ; user contrastive epoch  $E_u$ .  
**Output:** The trained inference networks  $\mathcal{I}_g$  and  $\mathcal{I}_u$ .

- 1: /\*\*\*\*\* Genre-aware Contrastive Phase \*\*\*\*\*/
- 2: **for** each genre contrastive epoch  $e_g \in \{1, 2, \dots, E_g\}$  **do**
- 3:   **for** each batch of reviews  $X$  in  $D$  **do**
- 4:     **for** each review  $x_{i,u}$  in  $X$  **do**
- 5:       Conduct augmentation and obtain  $\{x_{i,u}^a, x_{i,u}^b\}$ .
- 6:       Construct text representations for  $\{x_{i,u}^a, x_{i,u}^b\}$  with Eq. 5.
- 7:       Infer topic distributions with Eq. 6 and obtain  $\{\hat{\theta}_{i,u}^a, \hat{\theta}_{i,u}^b\}$ .
- 8:     **end for**
- 9:     Draw random samples  $\Theta'' = \{\hat{\theta}''_n\}_{n=1}^{2M}$  from the  $Dir(\hat{\theta}''|\vec{\alpha})$ .
- 10:     Estimate genre-aware matching objective  $\mathcal{L}_{Mg}$  via Eq. 15.
- 11:     Compute genre-aware contrastive objective  $\mathcal{L}_{Gc}$  via Eq. 9.
- 12:     Compute objective of genre-aware topic modeling  $\mathcal{L}_{TM}$  via Eq. 17.
- 13:     update  $\mathcal{I}_g$  by minimizing  $\mathcal{L}_{TM}$  with gradient descent.
- 14:   **end for**
- 15: **end for**
- 16: Freezing the parameters of topic-inference network  $\mathcal{I}_g$ .
- 17: /\*\*\*\*\* User-aware Contrastive Phase \*\*\*\*\*/
- 18: **for** each user contrastive epoch  $e_u \in \{1, 2, \dots, E_u\}$  **do**
- 19:   **for** each batch of reviews  $X$  in  $D$  **do**
- 20:     **for** each review  $x_{i,u}$  in  $X$  **do**
- 21:       Infer preference distribution  $\hat{p}_u$  for attached user  $u$  via Eq. 10.
- 22:       Infer personalized document-topic distributions  $\hat{\theta}_{i,u}^a$  and  $\hat{\theta}_{i,u}^b$  for augmented documents pair via Eq. 11.
- 23:     **end for**
- 24:     Draw random samples  $\Theta'' = \{\hat{\theta}''_n\}_{n=1}^{2M}$  from the  $Dir(\hat{\theta}''|\vec{\alpha})$ .
- 25:     Estimate the user-aware matching objective  $\mathcal{L}_{Mu}$  via Eq. 16.
- 26:     Compute user-aware contrastive objective  $\mathcal{L}_{Uc}$  via Eq. 14.
- 27:     Obtain objective of user preference modeling  $\mathcal{L}_{UM}$  via Eq. 18.
- 28:     update  $\mathcal{I}_u$  by minimizing  $\mathcal{L}_{UM}$  with gradient descent.
- 29:   **end for**
- 30: **end for**

### A.3 Topics and Personalized Topic Profiles

To intuitively compare the quality of topics and personalized topic profiles, we also present several examples extracted by NPTM and our proposed GPTM.

In Table 8, we present the mined topics and matched genres from the *Books* dataset with the topic number set to 15. It could be observed that topics extracted by GPTM often have a higher level of interpretability than those of NPTM. Additionally, GPTM tends to produce genre-aware topics, whereas NPTM generates nine meaningless topics (labeled with ‘-’) that are not semantically

**Table 8: Mined topics from the *Books* dataset by GPTM and NPTM on 15 topic settings.**

Model	Matched Genre	Topic Words
GPTM	Cookbooks, Food & Wine	Topic 1: recipe cookbook dessert dish chef meal delicious cook salad ingredient
	Business & Money	Topic 2: economy economic finance investment financial investor trader debt loan income
	-	Topic 3: sentence paragraph dialogue prose translation passage text narrator translate convey
	Mystery, Thriller & Suspense	Topic 4: mystery corpse clue puzzle whodunit cozy skeleton grave bury dead
	Science Fiction & Fantasy	Topic 5: interstellar solar alien mars planet galaxy ship fleet imperial sky
	Arts & Photography	Topic 6: artwork photographer pencil photographic museum illustration photograph photo photography artist
	Military	Topic 7: confederate war battlefield civil military soldier campaign army troop battle
	Mystery, Thriller & Suspense	Topic 8: cop police policeman detective criminal crime homicide attorney enforcement deputy
	Christian Books & Bibles	Topic 9: christian biblical bible scripture pastor ministry god gospel devotional testament
	Computers & Technology	Topic 10: software computer scientist google web manual internet digital user technology
	Romance	Topic 11: marry marriage married bride romantic girlfriend romance mate boyfriend husband
	Health, Fitness & Dieting	Topic 12: doctor medicine physician health yoga meditation heal anxiety healthy medical
	Science Fiction & Fantasy	Topic 13: magic mage fairy wizard fantasy werewolf magical supernatural witch paranormal
	Music	Topic 14: pupil musical composer musician piano song holocaust music teacher educational
	Business & Money	Topic 15: organizational workplace startup employee organization leadership business enterprise innovation company
NPTM	-	Topic 1: daughter father child mother family wife husband life married widow
	Mystery, Thriller & Suspense	Topic 2: suspect evidence bomb motive gunman defendant plot whereabouts unidentified attacker
	-	Topic 3: book write review insight experience learn focus publish business read
	-	Topic 4: story episode mystery character detail stories book evidence novel description
	Mystery, Thriller & Suspense	Topic 5: murder death sentence rape robbery killing trial guilty killer murderer
	Mystery, Thriller & Suspense	Topic 6: suspect character plot police dead bomb series story murder family
	-	Topic 7: father night mother brother kill vampire friend wife daughter king
	-	Topic 8: understaffed celibate invincible reintegrate droid inglorious cloistered eyeball vizer marriageable
	-	Topic 9: mother daughter child husband family father wife girl love woman
	-	Topic 10: author book read reader writer character story love novel novelist
	Computers & Technology	Topic 11: system computer technology systems software data program internet customer user
	-	Topic 12: book reader read write publish books fiction essay novel author
	-	Topic 13: miles nautical mile yards highway southwest coastline kilometer epicenter coast
	Mystery, Thriller & Suspense	Topic 14: sentence conviction sentencing trial prison jail murder manslaughter guilty criminal
	Cookbooks, Food & Wine	Topic 15: recipe recipes finder republish algorithm sauce cake menu cookbook dish

**Table 9: Personalized topic profiles from the *Books* dataset extracted by GPTM and NPTM, words related to the 1-st and 2-nd genres are marked with underline and wave, words without notations are irrelevant words.**

User	Preferred Genres	Model	Personalized topic profile
User 1	1-st: <u>History</u>	GPTM	<u>invention</u> <u>civilization</u> <u>evolution</u> <u>ancestor</u> <u>artifact</u> <u>origin</u> <u>dinosaur</u> <u>history</u> <u>century</u> <u>genesis</u> <u>gene</u> <u>mankind</u> <u>biology</u> <u>discovery</u> <u>technology</u> <u>anatomy</u> <u>chronology</u> <u>renaissance</u> <u>plant</u> <u>timeline</u>
	2-nd: <u>Science &amp; Math</u>	NPTM	murder death wife daughter brother husband mother kill killing married girlfriend dead father friend boyfriend woman sister widow cousin murderer
User 2	1-st: <u>Children's Books</u>	GPTM	<u>kindergarten</u> <u>kids</u> <u>grader</u> <u>learner</u> <u>grade</u> <u>youngster</u> <u>granddaughter</u> <u>grandson</u> <u>activity</u> <u>age</u> <u>alphabet</u> <u>entice</u> <u>caregiver</u> <u>motor</u> <u>classroom</u> <u>range</u> <u>daughter</u> <u>seven</u> <u>listener</u> <u>adore</u>
	2-nd: <u>Growing Up &amp; Facts of Life</u>	NPTM	book illustration story read text reader books <u>stories</u> novel reading copy manuscript poem author <u>color</u> photo paragraph graphic essay bible
User 3	1-st: <u>Business &amp; Money</u>	GPTM	<u>leadership</u> <u>enterprise</u> <u>employee</u> <u>organization</u> <u>management</u> <u>manager</u> <u>company</u> <u>initiative</u> <u>boss</u> <u>employer</u> <u>leader</u> <u>corporation</u> <u>strategy</u> <u>innovation</u> <u>entrepreneur</u> <u>productivity</u> <u>jobs</u> <u>business</u> <u>collaboration</u> <u>transformation</u>
	2-nd: <u>Management &amp; Leadership</u>	NPTM	<u>economic</u> <u>industry</u> <u>growth</u> <u>business</u> <u>billion</u> <u>economy</u> <u>market</u> <u>china</u> <u>investment</u> <u>global</u> <u>financial</u> <u>sector</u> <u>technology</u> <u>trade</u> <u>manufacturing</u> <u>consumer</u> <u>retail</u> <u>revenue</u> <u>sales</u> <u>company</u>
User 4	1-st: <u>Comics &amp; Graphic Novels</u>	GPTM	<u>comics</u> <u>superman</u> <u>superhero</u> <u>issues</u> <u>batman</u> <u>humor</u> <u>laughter</u> <u>cartoon</u> <u>comedy</u> <u>spider</u> <u>magazine</u> <u>flash</u> <u>costume</u> <u>wit</u> <u>hero</u> <u>adventures</u> <u>panel</u> <u>arc</u> <u>surfer</u> <u>skull</u>
	2-nd: <u>Science Fiction &amp; Fantasy</u>	NPTM	detective murder police crime mystery murders homicide suspect killer robbery plot murderer unsolved episode serial prosecutor assassination killing terrorist kidnapping

related to product genres. Such improvement is attributed to the incorporation of genre supervision.

To directly compare the performance of user preference mining, we present four personalized topic profiles in Table 9. The 'Preferred Genres' column lists the top two genres that the user is interested in.

These examples indicate that personalized topic profiles generated by GPTM could accurately match users' preferences, while NPTM may generate background topic profiles (User 2) and irrelevant topic profiles (User 1 and User 4). This is attributed to the usage of user-aware contrastive learning in GPTM.