

Toward Comprehensive User and Item Representations via Three-tier Attention Network

HONGTAO LIU, WENJUN WANG, QIYAO PENG, and NANNAN WU, Tianjin University
FANGZHAO WU, Microsoft Research Asia
PENGFEI JIAO, Tianjin University

Product reviews can provide rich information about the opinions users have of products. However, it is non-trivial to effectively infer user preference and item characteristics from reviews due to the complicated semantic understanding. Existing methods usually learn features for users and items from reviews in single static fashions and cannot fully capture user preference and item features. In this article, we propose a neural review-based recommendation approach that aims to learn comprehensive representations of users/items under a three-tier attention framework. We design a review encoder to learn review features from words via a word-level attention, an aspect encoder to learn aspect features via a review-level attention, and a user/item encoder to learn the final representations of users/items via an aspect-level attention. In word- and review-level attentions, we adopt the context-aware mechanism to indicate importance of words and reviews dynamically instead of static attention weights. In addition, the attentions in the word and review levels are of multiple paradigms to learn multiple features effectively, which could indicate the diversity of user/item features. Furthermore, we propose a personalized aspect-level attention module in user/item encoder to learn the final comprehensive features. Extensive experiments are conducted and the results in rating prediction validate the effectiveness of our method.

CCS Concepts: • **Information systems** → **Recommender systems**;

Additional Key Words and Phrases: Recommender system, context-aware, personalized, multi-aspect, attention

ACM Reference format:

Hongtao Liu, Wenjun Wang, Qiyao Peng, Nannan Wu, Fangzhao Wu, and Pengfei Jiao. 2021. Toward Comprehensive User and Item Representations via Three-tier Attention Network. *ACM Trans. Inf. Syst.* 39, 3, Article 25 (February 2021), 22 pages.
<https://doi.org/10.1145/3446341>

This work was supported by the National Key R&D Program of China (2018YFC0832100) and the National Natural Science Foundation of China (61902278).

Authors' addresses: H. Liu, W. Wang, Q. Peng, and N. Wu (corresponding author), College of Intelligence and Computing, Tianjin University, Tianjin China; emails: {htliu, wjwang, qypeng, nannan.wu}@tju.edu.cn; F. Wu, Microsoft Research Asia, Beijing, China; email: wufangzhao@tju.edu.cn; P. Jiao, Law School, Tianjin University, Tianjin, China; email: pjiao@tju.edu.cn.

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 ACM 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.

© 2021 Association for Computing Machinery.

1046-8188/2021/02-ART25 \$15.00

<https://doi.org/10.1145/3446341>

1 INTRODUCTION

Recommender Systems (RS) are information filtering systems that can help users to find products that they may be interested in [37]. In the age of information explosion, RS play an important role in many E-commerce platforms such as Amazon and Yelp. The key of most recommender systems is to learn accurate representations of users and items according to their historical records (e.g., ratings, reviews) so as to make a better recommendation product list for users.

There are many works proposed for recommender systems recently. Collaborative filtering (CF) technology [36] is one of the most widely used methods. Many of them are based on matrix factorization utilizing the ratings that users have scored toward items [18, 22, 30]. For example, Probabilistic Matrix Factorization (PMF) [30] models user and item based on the rating matrix and learns their latent factors via probabilistic matrix factorization. However, these rating-based methods suffer from the sparsity problem of the ratings, which will limit the ability of these methods to model users and items accurately. In many e-commerce platforms, users can post reviews to express their opinions of the products they have purchased. The reviews contain rich information that could fully reveal the user preference and item features [28]. As shown in Figure 1, from the text review we can observe the user attitude of the shoes and the ankle-support feature of the shoes. Therefore, many works have been proposed to exploit the review information to enhance the recommender systems and achieved better performance in rating prediction tasks [1, 4, 10, 21, 25–28, 44]. Previous methods usually adopt topic modelling to extract the semantic information of reviews and then incorporate it into user or item representation [1, 28]. HFT [28] utilizes matrix factorization method and regards review features as an additional regularizer to learn user and item latent features. These methods based on bag-of-word mechanism would be incapable of fully capturing the semantic information of reviews effectively [4]. Recently many neural network-based works are proposed and outperform those earlier methods significantly. These methods usually utilize the neural network (e.g., convolutional neural network (CNN)) to encode reviews and then learn representations of users and items from reviews, such as DeepCoNN [47], NARRE [4], TARMF [27], and CARP [21]. For example, DeepCoNN [47] first combines all the reviews of a user and an item into a long document, adopts word embedding to transform the document into word matrix, and utilizes convolutional neural network to learn semantic features as the representations of users and items.

Though these review-based methods have gained superior performance in rating prediction, it is still difficult to understand the complex rating behaviors between users and items. Due to the inherently multifaceted user preference and item features, it is nontrivial to learn precise representations of users and items. Specifically, in this article we suggest that there are three key characteristics to be considered in learning user and item features from reviews: (1) *Diverse*. Intuitively users always have multiple preference and the characteristics of items are always various (e.g., price, quality). Thus, to learn this diversity of users and items, it is imperative to extract the multi-aspect semantic information from their reviews instead of learning one single representation of reviews. (2) *Context-aware*. Most existing recommendation methods usually model user preference without considering the target item to be recommended, and learn a static user representation from reviews. Specifically, the same review from a user may have distinctive importance in learning user representation for disparate target items. For example, if a user has written reviews for a computer and a pair of shoes, we can conclude that the review about the computer is more useful than the shoes when the user would like to buy other digital products. Likewise, the same word in reviews of a user is different informative for different target items. In other words, this context information (user and the target item) can help describe users and items dynamically [24, 42]. (3) *Personalized*. It is intuitive that different users have different preference and



Fig. 1. A review example of a pair of shoes in Amazon.

different items have different characteristics. For example, suppose that User A cares more about price of item than the quality and User B cares more quality than price, the review “this phone is of high price but with good quality” differ in usefulness for User A and User B when learning their representations. Likewise, the same word in reviews may have different usefulness for different users and items. However, most existing recommendation methods usually use the same model to model all users/items from reviews and ignore the essential personalized features. Therefore, we suggest that all of the three characteristics above are essential to learn precise and comprehensive representations of users and items to improve the performance in recommendation.

To this end, in this article we propose a Multi-aspect neural Recommendation model with Context-aware Personalized attention (MRCP) by exploiting the textual reviews fully and deeply. In our approach, we devise a three-tier attention framework to learn more comprehensive representations of users and items. Specifically, our method consists of a review encoder with word-level attention to learn the multiple representations of reviews from words, an aspect encoder with review-level attention to extract semantic features of reviews, and a user/item encoder with an aspect-level attention to learn representations of users or items from their reviews. As denoted above, the same word or similar review of a user may have different informativeness with different target items. Hence in both review encoder and aspect encoder, we propose to utilize the context-aware attention whose query vectors are derived from the embedding of target item ID to select more informative words and reviews. In this way, we can learn dynamic and context-aware (i.e., candidate item specific) representations of users (vice versa for item modelling). Besides, the word-level attention of the review encoder and review-level attention of the aspect encoder are both multiple, where we utilize multiple query vectors to learn diverse features of words or reviews inspired by the multi-head mechanism in self attention [43]. Afterwards, considering the personalized characteristics of different users and items, we design a personalized aspect-level attention in the user/item encoder for each user and item to aggregate all the multiple representations of users and items. The personalized attention vector of each user/item is derived from the ID embeddings of the user or the item. Finally, we can obtain the comprehensive representations of users and items derived from reviews, and then in rating prediction layer, we adopt Latent Factor Model (LFM) [19] to predict the ratings that users would score toward items for recommendation. Note that the “aspect” in our model is implicit without any pre-processing tools to extract aspects in advance, which ensures that our method is end-to-end.

Our model is devised in the above attention fashions to indicate the three paradigms in recommendation. we adopt multiple paradigm in word- and review-level attention to demonstrate *Diverse*; the ID embeddings of the target items are used in the user attention networks (vice versa for item attention networks), which can be to depict *Context-aware*; the *Personalized* mechanism is

used in the aspect-level attention to estimate the different preference of different users (vice versa for items).

We conclude the contributions in this article as follows:

- We propose a neural recommendation model to learn more comprehensive and precise representations of users and items by exploiting their reviews to extract features for personalized recommendation. We design three encoders (i.e., review encoder, aspect encoder and user/item encoder) to extract hierarchical features of reviews, aspects and users/items.
- We propose to utilize three-tier attention modules (i.e., word-level attention, review-level attention, and aspect-level attention) to focus on those important and informative words, reviews and aspects in our model. The word- and review-level attentions are context-aware and multi-aspect, and the aspect-level attention is personalized. Through the three-tier attentions, the diverse, context-aware and personalized representations for users and items can be learned effectively.
- Extensive experiments are conducted on six real-world recommendation datasets to evaluate the performance of our method, the results demonstrate that our model MRCP can outperform many competitive state-of-the-art baselines.

This article is organized as follows: Section 2 presents the recent related works mainly in the review-based recommender systems. In Section 3, we describe our overall model in detail. The experimental settings and result analysis are reported in Section 4. We conclude our article in Section 5.

2 RELATED WORKS

Recommender systems play a crucial role in e-commerce platforms and a number of recommendation methods have been proposed. Collaborative filtering is one popular technology [11, 31, 35, 36, 40]. And many CF works are based on Matrix Factorization (MF) to learn latent features of users and items from the rating matrix data [13, 18, 20, 22, 30, 32]. For example, Mnih et al. [30] extended MF to PMF, and Koren et al. [19] introduced bias information of users and items into MF. Although these methods have shown good results in recommendation, the performance of these rating-based methods will degrade significantly when the rating matrix is very sparse. Considering that many e-commerce platforms allow users to post textual reviews toward the items, many works exploiting reviews have been proposed, since the reviews contain rich information about user preference and item features. Next we will present the related works in two paradigms, i.e., the traditional topic models and neural networks-based models for review-based recommender systems. Note that in our article we focus on rating prediction task in recommendation and leave other paradigms as our futures as sequential recommendation, top-N recommendation, and so on [9, 14, 15].

2.1 Review-based Recommendation with Topic Modelling

Topic modelling techniques such as Latent Dirichlet Allocation (LDA) [2] are widely used in learning semantic features of texts in natural language processing. Hence, some works exploit topic modelling to extract user and item latent factor features from review texts [1, 10, 23, 28]. For example, Mcauley et al. [28] used LDA model to extract the semantic topics of reviews, which were regarded as user and item representations for further recommendation. Yang et al. [1] proposed TopicMF and jointly learned latent features of users and items from ratings and reviews with MF and topic model technologies respectively. Diao et al. [10] designed a probabilistic model based on collaborative filtering and topic modeling to capture the interest distribution of users and the content distribution in movie recommendation. These methods utilizing reviews outperform models

that rely on rating matrix solely in recommendation. However, the topic modelling methods are based on bag of word mechanism, which would ignore the word order and hence could not learn accurate semantic features of reviews [45].

2.2 Deep Learning for Review-based Recommendation

2.2.1 Vanilla Neural Network Recommendation Methods. With the development of deep learning in various research fields, many works start to employ neural network such as CNN to learn deep representations of users and items from reviews [3, 16, 46, 47]. ConvMF [16] utilized a CNN network to extract semantic features from textual reviews for items, and integrated the features into matrix factorization framework for recommendation. DeepCoNN [47] combined all the reviews of users or items into long documents and then designed a parallel CNN model to learn the representations of users and items. These neural network-based methods have achieved better performance than the traditional topic modelling-based methods due to the superior capability of representation learning for reviews. However, these vanilla neural network methods always ignore that the different usefulness of different words or reviews of users and items, which would be insufficient to model users and items from their reviews.

2.2.2 Attention-based Recommendation Methods. Attention mechanism in deep learning has been popular in many areas e.g., natural language processing, computer vision. Attention is based on an intuition that the model should attend to a certain part when processing a large amount of information. Many attention-based works have been proposed in review-based recommendation [4, 5, 7, 21, 24–27, 38, 42]. For example, Seo et al. [38] first proposed a local and global attention network over word level to focus on the important words in review documents. CARL [45] utilized a context-aware attention to learn user-item pair-specific representation (i.e., context-aware) for rating prediction in recommendation. ANR [8] proposed an aspect-based co-attention network to extract different aspect information from the review documents. These methods usually focus on the word importance and ignore the different informativeness of different reviews. NARRE [4] considered the usefulness of reviews and introduced review-level attention to estimate the importance of different reviews for a user or an item. MPCN [42] utilized the review-by-review pointer-based co-attention network that could extract important reviews for users and items in a context-aware manner. However, these methods ignored the natural diverse characteristics of user preference or item features and only learned one aspect representations, which might not be sound. A very recent method, CARP [21], exploited a capsule network-based [34] recommendation system and regarded a user and an aspect of an item as a viewpoint to extract the diverse information of users and items from reviews, which achieved state-of-the-art performance in review-based recommendations. However, most existing methods learn the user/item representation with one same model and ignore the personalized features of users and items.

2.3 Comparison with Existing Works

In this article, our proposed method MRCP aims to learn comprehensive representations of users and items from reviews in recommendation. We summarize the recent deep learning-based recommendation methods under following fashions: word usefulness, review usefulness, diverse, context-aware, and personalized characteristics. In Table 1, we report a general comparison between the related works and our model MRCP in terms of the characteristics. We can find that our method MRCP equipped with all the five characteristics has the superiority to learn more accurate representations of users and items, which has the potential to improve the recommendation performance.

Table 1. Characteristic of Different Methods

	DeepCoNN	D-Attn	TARMF	NARRE	MPCN	ANR	DAML	CARP	MRCP
Word Usefulness	×	✓	✓	×	✓	×	✓	✓	✓
Review Usefulness	×	×	✓	✓	✓	×	×	×	✓
Diverse	×	×	✓	×	×	✓	×	✓	✓
Context-aware	×	×	×	×	✓	✓	✓	✓	✓
Personalized	×	×	×	×	×	×	×	×	✓

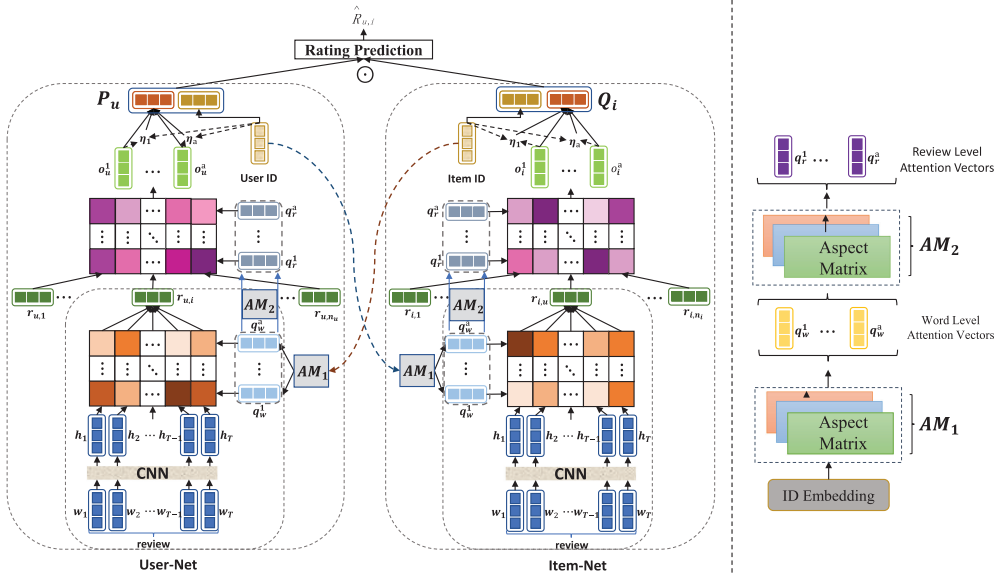


Fig. 2. Overview of our method MRCP, which learns the representation of user u or item i from reviews. Note that the AM_1 and AM_2 are the word-level aspect matrices containing a aspect matrix and the review-level aspect matrices including a aspect matrix, respectively.

3 PROPOSED METHOD

In this section, we will present our proposed method (MRCP) in detail. The overview of our approach is shown in Figure 2. There are three main modules in User Net (UN) and Item Net (IN), i.e., a word-level multi-aspect review encoder, an aspect encoder with multiple review-level attention, and the final user/item encoder with personalized attention. Considering that reviews can be regarded as a high abstract of words, the paradigm of our word-level and review-level attention is designed to be hierarchical shown in the right part in Figure 2. Then a prediction layer is devised to calculate the final rating score (i.e., $\hat{R}_{u,i}$) of the model. In the following sections, we will introduce the problem definition in this article first and then describe our method in detail. Since the User Net and the Item Net are similar in structure and only differ in their inputs, the details of the User Net will be illustrated in the following.

3.1 Problem Definition

Rating prediction is a fundamental problem in recommendation. Suppose that there are User set $U = \{u_1, \dots, u_N\}$ and Item set $I = \{i_1, \dots, i_M\}$, respectively, where N is the number of users and

Table 2. Notations and Their Definitions

Notations	Definitions
U, I	User Set and Item Set
N, M	Numbers of users and items
\mathbf{R}	Rating matrix
r_u	Reviews of user u
r_i	Reviews toward item i
n_u	Number of user u 's reviews
n_i	Number of item i 's reviews
$R_{u,i}$	The ground truth rating
$\hat{R}_{u,i}$	The prediction rating
ID^U	User ID embedding matrix
ID^I	Item ID embedding matrix
\mathbf{D}	Word embedding matrix
$r_{u,i}$	The review is posted by user u toward item i
\mathbf{H}	Semantic features of words in review
K	Number of convolution filters in CNN
d	Number of latent factors
a	Number of aspects
\mathbf{AM}_1	Aspect-specific word-level projection matrices
\mathbf{AM}_2	Aspect-specific review-level projection matrices
\mathbf{q}_w	Word-level attention vectors
\mathbf{q}_r	Review-level attention vectors
$\mathbf{r}_{u,i}$	Multi-aspect representations of the review
\mathbf{O}_u	Multi-aspect representations of user u
\mathbf{O}_i	Multi-aspect representations of item i
\mathbf{p}_u	The review-based representation of user u
ID^U_u	The ID embedding of user u
\mathbf{p}_u	The final representation of user u
\mathbf{p}_i	The review-based representation of item i
ID^I_i	The ID embedding of item i
\mathbf{Q}_i	The final representation of item i

M is the number of items. And we assume that the rating matrix is denoted as $\mathbf{R} \in \mathcal{R}^{N \times M}$, where user-item interaction (e.g., \mathbf{R}_{ui}) indicates the rating score of user $u \in U$ toward item $i \in I$ if the user u gives a rating to the item i . For a user u , all reviews written by u can be denoted by $r_u = \{r_{u,1}, \dots, r_{u,n_u}\}$ where n_u is the number of reviews. Similarly, for an item i , all reviews toward the item i can be represented as $r_i = \{r_{i,1}, \dots, r_{i,n_i}\}$, where n_i is the number of reviews. In this article, the primary objective is to predict the rating $\hat{R}_{u,i}$ for any unseen user-item pair, i.e., the rating of a given user u toward an item i that he/she has not interacted with before.

Table 2 summarizes the key notations used throughout the rest of this article.

3.2 User/Item ID Embedding

The ID embeddings of users and items are widely used in recommender systems [12, 33] and can be regarded as the identity information of corresponding users and items. Therefore, we project the

IDs of each user and item to low-dimensional vector representations as the ID embeddings in our method. Specifically, for the user u_1 , we associate the ID of user u_1 with a vector $\text{ID}_1^U \in \mathcal{R}^d$, where d is the ID embedding size. All users ID embeddings form an embedding matrix ID^U denoted as

$$\text{ID}^U = [\text{ID}_1^U, \text{ID}_2^U, \dots, \text{ID}_N^U], \quad (1)$$

where $\text{ID}^U \in \mathcal{R}^{N \times d}$ is the user ID embedding matrix and N is the number of users. The construction of the item ID embedding matrix is equivalent to the way the user ID embedding matrix is constructed, which is denoted as

$$\text{ID}^I = [\text{ID}_1^I, \text{ID}_2^I, \dots, \text{ID}_M^I], \quad (2)$$

where $\text{ID}^I \in \mathcal{R}^{M \times d}$ is the item ID embedding matrix, d is the dimension of item ID embedding, which is equal to that of user ID embedding, and M is the number of items.

The ID embeddings are randomly initialized for each user/item, and are updated during the training procedure, which can characterize their personalized intrinsic properties. We will incorporate the ID embeddings into our following three-tier attention networks. Specifically, the attention query vectors of word- and review-level in *User-Net* are learned from the candidate item ID embedding. Likewise, the attention query vectors of word- and review-level in *Item-Net* are learned from the corresponding user ID embedding.

3.3 Review Encoder with Word-level Attention

In this section, we design a review encoder for obtaining review representation from words. First, we utilize word embedding technology to project each word into a low-dimensional vector and then use CNN to extract the semantic features of words from user/item reviews. Besides, we propose to introduce the multi-aspect context-aware word-level attention mechanism into our method to learn the multifaceted review representations via the context-awareness paradigm.

3.3.1 Word Embedding and Convolution. As shown in User Net of Figure 2, the review posted by user u toward item i can be denoted as a sequence with T words, represented as $r_{u,i} = \{w_1, w_2, \dots, w_T\}$, $r_{u,i} \in \mathcal{R}^T$. Then we utilize word embedding to map each word in $r_{u,i}$ into a low-dimensional real-valued vector $\mathbf{w} \in \mathcal{R}^{d_w}$ with d_w as vector dimension. Unlike the topic-modelling methods, which relied on *bag of word* assumption and ignored the order and context of words, the embedding matrix \mathbf{D} is formed by stacking these words to preserve the order of words and computed as follows:

$$\mathbf{D} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_T], \quad (3)$$

where T is the length of the review $r_{u,i}$ and $\mathbf{D} \in \mathcal{R}^{T \times d_w}$ is the word embedding matrix of the review $r_{u,i}$. Recently some superior technologies have been proposed in natural language processing such as pre-trained language models (e.g., Bert). However these methods suffer from computation efficiency although they gain a rather slight improvement. Hence we still adopt static word embeddings in our method.

Afterwards, we utilize CNNs to extract the semantic feature of each word from the review embedding matrix. Compared with other popular neural models such as Recurrent Neural Networks and Transformers, CNNs are quite more efficient in computation complexity. Hence, in our model, we adopt CNNs as our basic encoder for reviews. Suppose that there are K different convolution filters, denoted as $\mathbf{F} = \{\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_K\}$, each filter is a parameter matrix $\mathbf{f}_j \in \mathcal{R}^{l \times d_w}$, where l is the filter window size. The convolutional result of the i th filter over the j th window on the whole word embedding matrix \mathbf{D} is computed as

$$z_{ij} = \text{ReLU}(\mathbf{D}_{j:j+1-1} * \mathbf{f}_i + b_i), \quad (4)$$

where $*$ is the convolution operation b_i is the bias and ReLU is a nonlinear activation function. Note that we adopt zero-paddings before the CNN layer to keep the dimension of the output consistent as the length of reviews. Then we can obtain the features $\mathbf{H} = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K\}$ produced by the K filters. Thus the m th row of $\mathbf{H} \in \mathcal{R}^{T \times K}$ is the feature of the m th word in the review $r_{u,i}$, denoted as $\mathbf{h}_m \in \mathcal{R}^K$.

3.3.2 Multi-aspect Context-aware Word Attention. Since we have obtained the semantic features of all words (i.e., the matrix \mathbf{H}) in the review, in this part, we will introduce how to learn review features from words. As mentioned above, when a user faces with different target items, the role of each word in the review may be different (i.e., context-aware information). Meanwhile considering the target item also has diverse aspects, such as quality, price, and so on, the importance of each word in the review is different based on different aspects, so we design a word-level multi-aspect review encoder based on context-aware mechanism to learn multi-aspect representations of review with context-awareness.

As shown in the right part of Figure 2, we first define the aspect-specific projection matrices over word level denoted as

$$\mathbf{AM}_1 = [\mathbf{W}_w^1, \dots, \mathbf{W}_w^a], \quad (5)$$

where \mathbf{AM}_1 represents a different potential aspects and a is the number of aspects. Each aspect is represented as a parameter matrix $\mathbf{W}_w^i \in \mathbf{AM}_1$ for word-level attention. To ensure our model is end to end, we define the aspect in our model as implicit features, and there is no need to pre-process with other tools to extract explicit aspects. And the number of aspects a is a hyperparameter in our model.

Since we would like to select important words under multiple aspects and different contexts, we generate the word-level attention query vectors from the target item ID embedding (to indicate the context-aware characteristic), and the aspects projection matrices \mathbf{W}_w , which is computed as

$$\mathbf{q}_w^i = \sigma(\mathbf{W}_w^i \mathbf{x}), i \in (1, a), \quad (6)$$

where the input \mathbf{x} is the item ID embedding (i.e., $\mathbf{x} \in \mathcal{R}^{d \times 1}$ contained in \mathbf{ID}^I), σ is an activation function, and \mathbf{q}_w^i can be regarded as the attention query vector under the i th aspect.

Then we use $\mathbf{q}_w = [\mathbf{q}_w^1, \dots, \mathbf{q}_w^a]^T$ as multiple attention vectors to focus on different aspects of the words in the review, the attention weight $\delta_{i,j}$ of the i th word in the review under the j th attention vector is denoted as follows:

$$\delta_{i,j} = \frac{\exp(Z_{i,j})}{\sum_{c=1}^T \exp(Z_{c,j})}, \delta_{i,j} \in (0, 1), j \in [1, 2, \dots, a], \quad (7)$$

$$Z_{i,j} = \mathbf{h}_i \odot \mathbf{q}_w^j, j \in [1, 2, \dots, T], j \in [1, 2, \dots, a], \quad (8)$$

where \mathbf{h}_i is the feature of the i th word in the review obtained in Section 3.3.1 (i.e., the i th row in \mathbf{H}), \odot is the inner product operator, and \mathbf{q}_w^j is the j th word-level attention vector derived by the target item and the aspect projection matrix.

For the j th aspect, we utilize weighted summation on all word features to obtain the representation of the review:

$$\mathbf{o}_r^j = \sum_{v=1}^T \delta_{v,j} \mathbf{h}_v, j \in [1, a], \quad (9)$$

\mathbf{h}_v is the feature of the v th word in the review. Note that the attentions in our model are designed with soft-attention paradigm. Compared with the hard/binary paradigm, soft attention could have the ability to select more important parts, and without losing any information. Besides, the hard attention is non-differentiable and requires more complicated techniques such as reinforcement

learning, while soft attention models are smooth and easy to train end to end. Thus, we adopt soft attention in our three-tier attentions.

Then we combine all representations of the review under a aspect attention vectors as follows:

$$\mathbf{r}_{u,i} = \mathbf{o}_r^1 \oplus \mathbf{o}_r^2 \oplus \cdots \oplus \mathbf{o}_r^a, \quad (10)$$

where $\mathbf{r}_{u,i} \in \mathcal{R}^{K*a}$ is the final feature vector of the review $r_{u,i}$ based on multi-aspect with context-awareness and \oplus is the concatenation operator. Afterwards, to match the dimension of ID embedding (i.e., d), we adopt a fully connected layer (FC) to transform the feature as

$$\mathbf{r}_{u,i} = \mathbf{W}_{fc}\mathbf{r}_{u,i} + \mathbf{b}, \quad \mathbf{r}_{u,i} \in \mathcal{R}^d, \quad (11)$$

where $\mathbf{W}_{fc} \in \mathcal{R}^{(K*a) \times d}$ is the matrix weight and \mathbf{b} is the bias item of the FC layer. $\mathbf{r}_{u,i}$ is the final feature vector of the review $r_{u,i}$.

3.4 Aspect Encoder with Review-level Attention

Since we have obtained the review features of the user u (i.e., $\mathbf{r}_u = [\mathbf{r}_{u,1}, \mathbf{r}_{u,2}, \dots, \mathbf{r}_{u,n_u}]$), in this section we will learn the aspect-specific representation of users via the aspect encoder.

As stated above, different reviews contribute differently to represent users, meanwhile considering a target item has nature diversity, each review is represented differently under disparate aspect while facing the specific item. Hence, similarly to the word-level attention, we apply another multi-aspect context-aware attention over review level to learn multiple context-aware representations of the users/items from their reviews.

Considering that review can be viewed as a high level abstract of words, the review-level attention vectors should be derived from word-level attention vectors. As shown in the Figure 2, review-level attention query vectors $\mathbf{q}_r = [\mathbf{q}_r^1, \mathbf{q}_r^2, \dots, \mathbf{q}_r^a]^T$ are derived from word-level attention vectors \mathbf{q}_w and another review-level aspect matrix \mathbf{W}_r (i.e., $\mathbf{AM}_2 = [\mathbf{W}_r^1, \dots, \mathbf{W}_r^a]$), denoted as

$$\begin{aligned} \mathbf{q}_r^i &= \sigma(\mathbf{W}_r^i \mathbf{q}_w^i), \quad i \in [1, a], \\ \mathbf{q}_r &= \text{stack}[\mathbf{q}_r^1, \dots, \mathbf{q}_r^a]^T, \end{aligned} \quad (12)$$

where $\mathbf{W}_r^i \in \mathcal{R}^{d \times d}$ is i th review-level aspect projection matrix, a is the number of review-level aspects, $\mathbf{q}_r \in \mathcal{R}^{d \times a}$ is the review-level attention matrix based on a aspects. The attention weight $\zeta_{i,j}$ of the review $\mathbf{r}_{u,i}$ via j th attention vector is denoted as

$$\zeta_{i,j} = \frac{\exp(X_{i,j})}{\sum_{b=1}^{n_u} \exp(X_{b,j})}, \quad \zeta_{i,j} \in (0, 1), \quad j \in [1, 2, \dots, a], \quad (13)$$

$$X_{i,j} = \mathbf{r}_{u,i} \odot \mathbf{q}_r^j, \quad i \in [1, 2, \dots, n_u], \quad j \in [1, 2, \dots, a], \quad (14)$$

where \odot is the inner product operator and \mathbf{q}_r^j is the j th review-level attention vector.

Afterwards, we utilize weighted summation to aggregate all review representations under j th aspect attention vector:

$$\mathbf{o}_u^j = \sum_{i=1}^{n_u} \zeta_{i,j} \mathbf{r}_{u,i}, \quad j \in [1, a]. \quad (15)$$

In this way, we can obtain all the a feature vectors of the user u , denoted as

$$\mathbf{O}_u = \text{stack}[\mathbf{o}_u^1, \mathbf{o}_u^2, \dots, \mathbf{o}_u^a], \quad (16)$$

where $\mathbf{O}_u \in \mathcal{R}^{d \times a}$ is the representation matrix of user u under a aspects. Likewise, from the aspect encoder in this section, we can obtain the representation of item i denoted as $\mathbf{O}_i \in \mathcal{R}^{d \times a}$. Next we will explore to derive the user/item final representation from \mathbf{O}_u and \mathbf{O}_i , respectively.

3.5 User/Item Encoder with Aspect-level Attention

As demonstrated in Section 1, users always have personalized preference and items have individual features; hence different aspects would have different importance in learning representations of users and items. Since we have obtained the multi-aspect representations (i.e., \mathbf{O}_u) of the user u , a simple way is to concatenate or sum the representations of all the aspects. However, different users are in favour of different aspects (*Personalized*), and, hence, we propose an aspect-level personalized attention to focus on the user-specific aspect and utilize user ID embedding as the query vector to guide the personalized attention learning.

To be specific, given the multi-aspect-based user representations (i.e., $\mathbf{O}_u = [\mathbf{o}_u^1, \mathbf{o}_u^2, \dots, \mathbf{o}_u^a]$), we compute the attention weight of j th representation under the guidance of user u ID embedding as follows:

$$\eta_j = \frac{\exp(l_j)}{\sum_{k=1}^a \exp(l_k)}, \eta_j \in (0, 1), \quad (17)$$

$$l_j = \mathbf{ID}_u^U \odot \mathbf{o}_u^j, j \in [1, a], \quad (18)$$

where \mathbf{ID}_u^U is the attention vector (i.e., the ID embedding of user u). Likewise for items, we could obtain the weight of j th representation under the guidance of item i ID embedding as follows:

$$\lambda_j = \frac{\exp(l_j)}{\sum_{k=1}^a \exp(l_k)}, \eta_j \in (0, 1), \quad (19)$$

$$l_j = \mathbf{ID}_i^I \odot \mathbf{o}_i^j, j \in [1, a], \quad (20)$$

Finally, we can obtain the user and item representation \mathbf{p}_u and \mathbf{p}_i via aggregating multi-aspect representations according to their weights:

$$\mathbf{p}_u = \sum_{j=1}^a \eta_j \mathbf{o}_u^j, \quad (21)$$

$$\mathbf{p}_i = \sum_{j=1}^a \lambda_j \mathbf{o}_i^j. \quad (22)$$

Considering that there are some users or items with very few or no reviews, and user and item ID embedding can describe their intrinsic attributes, we combine the above representation \mathbf{p}_u and \mathbf{p}_i learned from reviews and their ID embeddings,

$$\mathbf{P}_u = \mathbf{p}_u \oplus \mathbf{ID}_u^U, \quad (23)$$

$$\mathbf{Q}_i = \mathbf{p}_i \oplus \mathbf{ID}_i^I. \quad (24)$$

Here \mathbf{P}_u and \mathbf{Q}_i are the final representations of the user u and the item i .

3.6 Rating Prediction

In this section, we will report the final rating prediction process. We first combine the user feature \mathbf{P}_u and the target item feature \mathbf{Q}_i as follows:

$$\mathbf{b}_{u,i} = \mathbf{P}_u \cdot \mathbf{Q}_i, \quad (25)$$

where \cdot is the element-wise product.

Following previous work [4], we utilize LFM [19] to compute the rating that the user would score the item. The rating $\hat{R}_{u,i}$ is computed by

$$\hat{R}_{u,i} = \mathbf{W}_p^T(\mathbf{b}_{u,i}) + b_u + b_i + \mu, \quad (26)$$

where \mathbf{W}_p^T denotes the parameter matrix of the LFM model, \mathbf{b}_u denotes the user bias, \mathbf{b}_i denotes the item bias, and μ denotes the global bias.

3.7 Optimization

3.7.1 Complexity Analysis. The trained parameters of our model include: (1) three embeddings: user ID embeddings $\mathbf{ID}^U \in \mathcal{R}^{N \times d}$, the item ID embeddings $\mathbf{ID}^I \in \mathcal{R}^{M \times d}$, and the word embedding $\mathbf{D} \in \mathcal{R}^{T \times d_w}$; (2) two aspect projection matrices, $\mathbf{AM}_1 \in \mathcal{R}^{a \times d \times K}$ and $\mathbf{AM}_2 \in \mathcal{R}^{a \times d \times d}$; and (3) three parameter matrices: the convolution filter $\mathbf{F} \in \mathcal{R}^{K \times l \times d_w}$, the matrix $\mathbf{W}_{fc} \in \mathcal{R}^{(a \times K) \times d}$, and the $\mathbf{W}_p^T \in \mathcal{R}^{1 \times d}$ in rating prediction, and the corresponding bias items. Hence, the number of the whole parameters is $[M + N + a * (2K + d) + 1] * d + (T + K * l) * d_w + M + N + K + 1$. Generally, a, K is usually quite smaller than the numbers of users and items M and N , so the parameter complexity of our method is linear with the input size (i.e., M and N), which is similar with that in baseline methods. Since our method is an end-to-end neural network model, we can utilize mini-batch strategy to train the model effectively.

Considering that our method is context-aware, the final user representation would be related to the candidate item, which would indeed increase the computational complexity along with performance improvement. However, there are always two phases in real-world recommendation scenarios, including retrieval and ranking phases. The retrieval layer aims to retrieve a small candidate set of relevant items in low latency and computational cost from a large number of items (maybe all items); and the ranking layer targets to rank the most desired documents on the top with more complex models based on the retrieval results (maybe hundreds or thousands of items). In fact the complicated recommendation models are almost deployed in the ranking phase, such as DAML [24] and DIN [48], which utilize context-aware characteristic as well. Similarly, our proposed method aims to improve the recommendation performance in ranking phase instead of the retrieval phase. Hence, the impact of complexity would be small within tolerance.

3.7.2 Training. The task in this work is rating prediction, which is a regression problem, hence we utilize the square loss function to train our model:

$$L_{sq} = \sum_{u, i \in \Omega} (\hat{R}_{u,i} - R_{u,i})^2, \quad (27)$$

where Ω denotes the set of instances for training, and $R_{u,i}$ is the ground truth rating assigned by the user u to the item i .

4 EXPERIMENTS

In this section, we have conducted extensive experiments over several recommendation datasets to evaluate the performance of our proposed model. We will introduce the datasets and experimental settings first, then present the baseline algorithms selected for comparisons. Furthermore, we discuss performance evaluation and analysis of model effectiveness, respectively. Finally, hyperparameters analyses and case studies are presented.

4.1 Datasets and Experimental Settings

4.1.1 Datasets. We use the real-world datasets from *Yelp* and *Amazon* in our experiments. The first dataset is selected from Yelp Dataset Challenge,¹ and we extract the records in Yelp spanning from 2016 to 2017 denoted as *Yelp16-17*, since the whole Yelp dataset is too large and sparse. The other five datasets are selected from the Amazon dataset² and they come from five different domains, i.e., *Musical Instruments*, *Digital Music*, *Video Games*, *Office Products*, and *Tools Improvement*. Besides, we create another large Amazon dataset by combining all the five datasets from different

¹<https://www.yelp.com/dataset/challenge>.

²<http://jmcauley.ucsd.edu/data/amazon/>.

Table 3. Statistical Details of the Datasets

Datasets	# users	# items	# ratings	# avg.review length	# words per user	# words per item	density (%)
Musical Instruments	1,429	900	10,261	32.45	141.32	200.12	0.798
Digital Music	5,540	3,568	64,666	69.57	216.21	266.51	0.327
Video Games	24,303	10,672	231,577	72.13	188.79	260.60	0.089
Office Products	4,905	2,420	53,228	48.15	197.93	229.52	0.448
Tools Improvement	16,638	10,217	134,345	38.75	162.53	212.48	0.079
Yelp16-17	167,106	100,229	1,217,208	38.86	133.60	155.18	0.007

domains, named *Amazon*. Note that all datasets consist of users explicit ratings on items (from 1 to 5) and contain textual reviews with respect to user-item pairs. We follow the preprocessing method in previous work [4, 21] and keep the length and the number of reviews covering p percentage users and items respectively to address the long tail effect of the length and quantity of reviews. We set the percentage $p = 0.9$ in our experiments.

The detailed characteristics of the datasets are shown in Table 3.

For evaluation, we randomly split each dataset into a training set and testing set, with the ratio 90%/10% respectively. We further split 10% of the training set as the validation set. The validation set is to tune hyper-parameters in our method. At least one interaction per user/item is included in the training set. Considering that the target reviews are not available in the real-world scenarios, therefore the target reviews are excluded from validation and testing sets.

4.1.2 Hyper-parameters Settings. In our experiment, we adopt the validation set to tune the hyper-parameters in our model and report hyper-parameters setting here. We set the dimension of word embedding d_w to 300, and we adopt the pre-trained word embedding that is trained on more than 100 billion words from Google News [29]. The dimension of user or item ID embedding (i.e., d) and the number of latent factors (i.e., dimension of \mathbf{P}_u or \mathbf{Q}_i) are set to 32 (tuning in [8, 16, 32, 64]), the number of convolution filters (i.e., K) is 100 (tuning in [50, 100, 150, 200]) and the window size of CNN (i.e., l) is 3. The batch size of our model is set to 256. We utilize the dropout technology [39] to alleviate the overfitting problem and set the dropout ratio is 0.6. We adopt Adam [17] optimization strategy to optimize our model and set the learning rate to 0.001, and the weight decay is 0.0001. The student t -test is applying in our experiments to conduct the statistical significance. And the comparison difference of our proposed method over all baselines are significant with $p < 0.05$.

4.1.3 Evaluation Metric. To verify the effectiveness of our proposed model, Mean Square Error (MSE) is utilized as the evaluation metric. It is noticed that the lower MSE indicates the better performance. Given the ground truth rating $R_{u,i}$ and the rating score $\hat{R}_{u,i}$ estimated by model, we compute the MSE as follows:

$$MSE = \frac{\sum_{u,i \in \Omega} (\hat{R}_{u,i} - R_{u,i})^2}{|\Omega|}, \quad (28)$$

where Ω indicates the set of the user-item pairs in the testing set. Note that in our experiments, we independently repeat each experiment for 5 times and report the average results in terms of MSE for fair comparison.

4.1.4 Baselines. To evaluate the performance of our method MRCP, we compare MRCP with the recent competitive state-of-art baseline methods including:

Table 4. MSE Comparisons of Our Model MRCP (in Bold) and Baseline Methods

Method	Musical Instruments	Digital Music	Office Products	Tools Improvement	Video Games	Amazon	Yelp16-17
PMF	1.398	1.206	1.092	1.566	1.672	1.504	2.574
RBLT	0.815	0.870	0.759	0.983	1.143	1.122	1.569
DeepCoNN	0.814	1.056	0.860	1.061	1.238	1.192	1.593
D-Attn	0.982	0.911	0.825	1.043	1.145	1.136	1.573
TARMF	0.943	0.853	0.789	1.069	1.195	1.168	1.914
NARRE	0.910	<u>0.812</u>	0.732	0.957	1.112	1.087	1.522
MPCN	0.824	<u>0.903</u>	0.769	1.017	1.201	1.171	1.617
ANR	0.795	0.867	0.742	0.975	1.182	1.166	1.553
DAML	0.871	0.813	<u>0.705</u>	<u>0.945</u>	1.165	1.113	1.601
CARP	<u>0.773</u>	0.820	0.719	0.960	<u>1.084</u>	<u>1.064</u>	<u>1.508</u>
MRCP	0.760	0.801	0.702	0.928	1.069	1.043	1.489

The best performance of the baselines are underlined. The improvements of our proposed method over all baselines are significant with $p\text{-value} < 0.05$.

- **PMF** [30]: Probabilistic matrix factorization with Gaussian distribution is introduced to model the latent factors for users and items.
- **RBLT** [41]: This method utilizes a rating-boosted approach with rating-boosted reviews and rating scores for rating prediction.
- **DeepCoNN** [47]: The model utilizes two parallel CNNs to learn the latent feature vectors of user and item from reviews respectively and FM mechanism as prediction layer.
- **D-Attn** [38]: The dual local and global attentions are leveraged to enable an interpretable embedding of users and items and capture the precise semantic features.
- **TARMF** [27]: A novel recommendation model utilizes attention-based recurrent neural networks to extract topical information from review documents.
- **NARRE** [4]: A newly proposed method that introduces neural attention mechanism to complete the rating prediction task and select highly-useful reviews simultaneously.
- **MPCN** [42]: A review-by-review pointer-based learning scheme is exploited to enable informative reviews and deeper word-level interaction between user and item.
- **ANR** [8]: The state-of-art aspect-based neural recommendation systems performs aspect-based representation learning of users and items for rating prediction.
- **DAML** [24]: DAML adopts local attention layer and mutual attention layer for selecting informative words and realizing the dynamic interaction of users and items.
- **CARP** [21]: CARP proposes a novel deep learning model that exploits capsule network and reviews for rating prediction and explains what the users like and dislike.

Note that there are many more other models, such as HFT [28], JMARS [10], ConvMF [16], ALFM [7], A³NCF [6]. These works have been outperformed by one or more baselines compared here. Therefore, we omit to compare our method with these approaches.

4.2 Performance Comparison

Table 4 reports the results of all the methods over all the datasets. We compare our method MRCP with the above baselines in terms of MSE and we have the following observations. First, it is not surprising that the review-based methods consistently outperform the rating-based method (i.e., PMF) over all all datasets. This is because reviews contain rich information of user preference and

item features that would be useful for rating prediction. This observation is consistent with many review-based works [4, 21].

Second, the neural network methods (i.e., DeepCoNN, D-Attn, TARMF, NARRE, MPCN, ANR, CARP, MRCP) usually perform better than the topic-based model (i.e., RBLT). The reason is that neural networks have powerful representation capabilities and can learn more semantic information from reviews to model users and items, while the topic-based modeling approaches may not fully capture deep review features. In addition, deep learning can model user and item in a non-linear way, which is the limitation of traditional models. Furthermore, we find that the methods utilizing attention mechanism (e.g., D-Attn, TARMF, NARRE, MPCN, ANR, MRCP) generally perform better than those without attention (e.g., DeepCoNN). We argue that this phenomenon is because there is noise information in text reviews, and the attention mechanism can focus on the more important reviews and words, which will help models to learn the representations of users and items more precisely. This observation is also in line with the experimental results reported in References [4, 38] and validates the effectiveness of the attention mechanism.

Third, as shown in Table 4, our method MRCP achieves the best MSE score over all the datasets. And compared with all the baselines, our method achieves the minimum 1.3% and the maximum 24.1% improvement on all different datasets. We can find out the reason and clue from Table 1. Our method MRCP utilizes multi-aspect word- and review-level attention to select those important words and reviews under different aspects. In addition, MRCP can learn dynamic representations of users and items under the context-aware scheme instead of single static representations (e.g., NARRE, TARMF). Furthermore, compared with other methods, our MRCP could learn the personalized representations for users and items from personalized aspect-level attention. Our method outperforms two very recent work DAML [24] and CARP [21]. The reason may be that DAML adopts a Euclidean distance approach to capture context-aware interaction information between user reviews and item reviews, which may not take diversity and personalization into account, and CARP omits the review-level information. And both of them are incapable of exploiting the personalization of users and items. Overall, the experimental results meet our motivation in Section 1 and demonstrate the effectiveness of our method in recommendation.

4.3 Analysis of MRCP

In this section, we investigate the effectiveness of different components in our model. The three-tier attentions are the most important modules in our method, and hence we try to design elaborate experiments to analyze their effects. In addition to the three-tier attentions, other basic technologies in our method including word embeddings (word2vec) and text feature extraction methods (CNN) are not the focus in this article, and hence we conduct no experiments to study their impacts.

4.3.1 Effect of Different Level Attention. In our model, we use word- and review-level attention to indicate the usefulness of different words and reviews, and then a personalized aspect-level attention is utilized for learning a unique representation of the user or item. In this section, we conduct ablation experiments to investigate the contributions of different levels attention (i.e., word-level, review-level, aspect-level) in our MRCP across three datasets. We design three different model variants by removing different attention as follows:

- MRCP-WA: the model without word-level attention (WA), and averages the features of all the words in a review.
- MRCP-RA: the model without review-level attention (RA), and averages the review representation as the user representation.

Table 5. MSE Comparison of Different Attention Components

Variant	Tools Improvement	Digital Music	Office Products
MRCP-WA	0.931	0.808	0.709
MRCP-RA	0.933	0.807	0.708
MRCP-PA	0.958	0.815	0.728
MRCP	0.928	0.801	0.702

- MRCP-PA: the model without personalized aspect-level attention (PA), and averages all aspect representations of users and items rather than utilizing personalized aspect-level attention.

The experimental results are reported in Table 5. We can find that removing any attention would degrade the performance, which indicates the effectiveness of our three attention components. In word-level, it is not surprising that utilizing word-level attention is better than averaging each word representation. This is because while learning the representation of a review, the multi-aspect attention mechanism can learn the aspect-specific importance for each word of the review, nevertheless the average method treats all words equally. Similarly in review-level, different reviews may be different informativeness for user and item representations, hence simply averaging the review representation would be insufficient and incapable of fully capturing the review usefulness.

In addition, the model MRCP-PA achieves the worst MSE scores, which indicates the important effect of personalized aspect-level attention. This is because users have their individual inherent preference and items have individual inherent attributes, which could influence the weight of each aspect for representing users and items. We utilize the ID embedding of users or items to guide the aspect-level personalized attention procedure to capture personalized user preference and item features. The combination of three modules can further improve the performance in recommendation, which validates the effectiveness of our model.

4.3.2 Effect of Context-Aware Attention. Users will show specific preference and the same word/review may reflect different implication while facing different target items. Hence, our method introduces context-aware mechanism to capture the user-item pair-based relevant information (i.e., context-aware information). The attention components of our model MRCP utilizes context-aware mechanism (i.e., the attention query vectors are learned from target item ID while modelling user) to learn dynamic representations of users and items instead of static representations. In this section, we conduct experiments over three datasets to verify the validity of the context-aware mechanism. Instead of using the target ID embedding as attention query vector, we randomly initialize the \mathbf{q}_w and \mathbf{q}_r to guide the word- and review-level attention procedure respectively, denoted as

- MRCP-CWA: the model randomly initializes the user/item word-level attention query vectors instead of deriving from target item/user ID embedding.
- MRCP-CRA: the model randomly initializes the user/item review-level attention query vectors instead of deriving from target item/user ID embedding.
- MRCP-CWRA: the model randomly initializes the user/item word- and review-level attention query vectors instead of deriving from target item/user ID embedding.

The experimental results in Table 6 show that randomly initializing user/item word- or review-level attention vectors would degrade the performance in recommendation. It is worth noted that simultaneously randomly initializing the user/item word- and review-level (i.e., MRCP-CWRA)

Table 6. Effect of Context-aware Mechanism

Variant	Tools Improvement	Digital Music	Office Products
MRCP-CWA	0.949	0.829	0.717
MRCP-CRA	0.944	0.835	0.719
MRCP-CWRA	0.951	0.838	0.722
MRCP	0.928	0.801	0.702

Table 7. Effect of Attention in User-/Item-Net

Variant	Tools Improvement	Digital Music	Office Products
MRCP-UN	0.938	0.821	0.715
MRCP-IN	0.940	0.826	0.712
MRCP-UIN	0.945	0.832	0.721
MRCP	0.928	0.801	0.702

obtains the worst performance. These phenomena indicate that the effectiveness of context-aware mechanism. In this way, the model learns that the user/item representation is not static, but a dynamic representation based on the target item/user.

4.3.3 Effect of Attention Mechanism in User-/Item-Net. In our model, we adopt attention mechanism on User-/Item-Net to help the User-/Item-Net select informative words and reviews and learn more precise representations of user/item. In this section, we conduct ablation experiments to evaluate the effectiveness of attention mechanism in User- and Item-Net over three datasets respectively. We design three different model variants by removing attention mechanisms from our model as follows:

- MRCP-UN: The UN averages the semantic features of all words and reviews without utilizing multi-aspect word-/review-level attention and personalized attention, while the Item-Net is the same with the main experiment setting.
- MRCP-IN: The IN averages the semantic features of all words and reviews without utilizing multi-aspect word-/review-level attention and personalized attention, while the User-Net is the same with the main experiment setting.
- MRCP-UIN: The model without any attention mechanism in UN and IN averages all semantic features to learn representations of users and items.

The experimental results are published in Table 7. We can find that eliminating attention mechanisms from User-Net or Item-Net would degrade predictive performance of our model. Furthermore, it is not surprising that the model with removing the attention mechanisms from both User-Net and Item-Net obtains the worst rating prediction performance. These results indicate the effectiveness of our multi-aspects context-aware and personalized attention mechanism in User-/Item-Net. Because different users or items always have their different multiple preference and characteristics while facing different items and users. Hence, the multi-aspect context-aware attention mechanisms are utilized in User-Net and Item-Net for capturing multiple context-aware user preference and item features, respectively. Meanwhile, considering that users and items have their own unique preference and characteristics, the personalized attention mechanism is utilized in User-Net and Item-Net for capturing the personal preference and characteristics of users and items, which are conducive for modelling users and items.

Table 8. The Results of Different Aspect Numbers in MRCP

Number of Aspects	Musical Instruments	Digital Music	Office Products	Tools Improvement	Video Games	Amazon	Yelp16-17
a=1	0.771	0.814	0.719	0.937	1.079	1.070	1.533
a=3	0.765	0.811	0.718	0.948	1.077	1.066	1.520
a=5	0.760	0.805	0.702	0.928	1.073	1.043	1.489
a=7	0.768	0.801	0.716	0.956	1.069	1.051	1.521

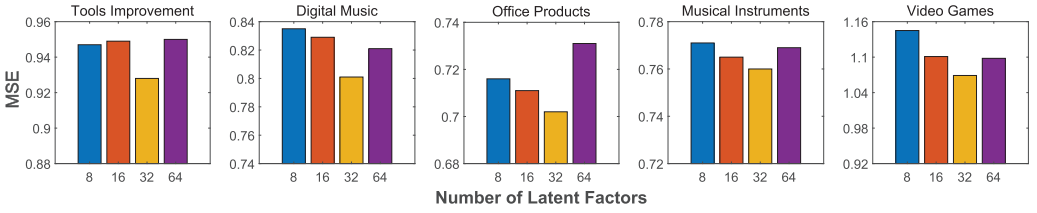


Fig. 3. The impact of Number of Latent Factors.

4.4 Hyper-Parameters Analyses

In this section, we conduct hyper-parameter analysis experiments to explore the effects of hyper-parameters on the validity of the model in this article. We mainly analyze three important hyper-parameters: the number of aspects, the number of latent factors, and the number of convolution filters in CNN.

4.4.1 Effect of Number of Aspects. Considering that the number of aspects plays a critical role in the maintenance of the overall performance of our model MRCP, we conduct experiments to analysis the impact of different number of aspects across all the datasets. Table 8 illustrates the effect of varying the number of aspects (i.e., a) arranged from $\{1, 3, 5, 7\}$. From the experimental results, we can find that (1) the optimal number of aspects varies from different datasets. (2) multi-aspect model (i.e., $a > 1$) can achieve better performance than that with only one aspect (i.e., $a = 1$), which indicates the effectiveness of multi-aspect diverse representations of users and items. We argue that the optimal aspect number is related to the review length in datasets, since the longer the review, the more aspects about users or items would be covered (e.g., Digital Music and Video Games datasets). In addition, we can find that in the large Amazon dataset, $a = 5$ is the optimal choice. From the results, we can observe that in most datasets of our experiments the model with 5 as aspect number perform best. Hence in real applications, setting the aspect number to 5 would be a good choice for recommender platforms.

4.4.2 Number of Latent Factors. Since the latent factors of users or items contain rich user preference and product feature information, the dimension of latent factors is important for representing users and items. We explore the effect of the number of user/item latent factors, i.e., the dimension of $\mathbf{P}_u, \mathbf{Q}_i$. As shown in Figure 3, we can find that as the dimension of latent factors increases, the MSE first decreases and then reaches the lowest (i.e, the best performance) when the dimension of latent factors is 32 and increases afterwards. We argue that when the dimension of latent factors is too small, the latent factors may not be able to capture the inherent preference of users and potential attributes of items. When the dimension is too large, the model may suffer from overfitting and the computational complexity would increase. Hence we set up the optimal dimension of latent factors is 32.

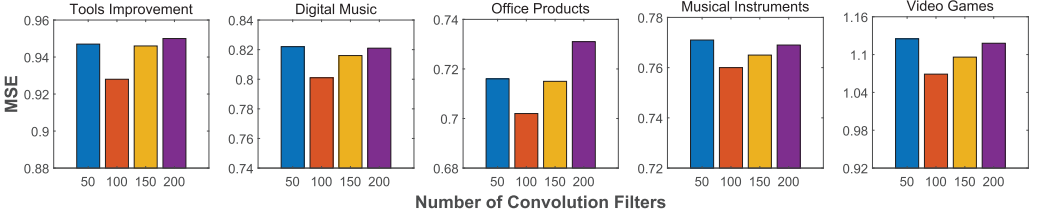


Fig. 4. The impact of Number of Convolution Filters.

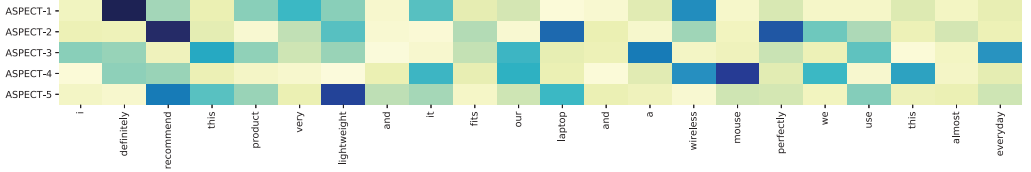


Fig. 5. Visualization of the weights of words in multi-aspect word-level attention.

4.4.3 Number of Convolution Filters. Considering that the number of convolution filters might affect the ability of the convolutional neural network to extract semantic features from review. We analyze the influence of the number of convolution filters on the model performance. From Figure 4, with the increase of the convolution filters number, the MSE declines at first, and reaches the minimum value (i.e., the best model performance) when the convolution filters number is 100, and then increases. When convolution filters number is too small, the convolution operator may not extract rich information of users and items from reviews. However, if it becomes too large, then the model may suffer from overfitting and degrade performance. Therefore we set the optimal number of convolution filters is 100 regardless of different datasets.

4.5 Case Study

In this section, we visualize multi-aspect attention weights in word- and review-level, respectively, to qualitatively study the behavior of our multi-aspect attention mechanism intuitively. In word-level attention, we save the multiple weights of words in a review under different aspects by Equation (6), and in review-level attention, we save the multiple weights of all the reviews of a user under different aspects by Equation (12). We sample reviews from *Office-Products* datasets for case study. Note that the deeper the color, the larger the attention weight.

As shown in Figure 5, we can see that the weight of each word in the review is different based on different word-level aspect. For the Aspect-1, the words “definitely” and “wireless” are most informative for representation of the review while the word “recommend” and “laptop” contribute great deals to expression of the review with the Aspect-4. This indicates the effectiveness of the word-level multi-aspect attention mechanism, i.e., it can support our model to learn multiple representations of reviews in view of multi-aspect from words.

We select all the five reviews of a user from the dataset, and visualize the attention weights under different aspects shown in Figure 6. it is apparent that the weight of each review written by the user is disparate established in multiple review-level aspects. The Review-2 is most informative for the user representation with Aspect-5. Nevertheless, under the Aspect-2, the Review-3 contributes the most to representing user. This manifests the availability of the review-level multi-aspect attention mechanism, i.e., it can support our model to learn multiple representations of the user based on multi-aspect from reviews. In a word, utilizing the word- and review-level multi-aspect

Review-1	i bought these to use with my small household shredder works great with the shredder no problems with the product
Review-2	seemed to be full yield i didn t have any leaks like some other users had they seemed to work as described
Review-3	i definitely recommend this product very lightweight and it fits our laptop and a wireless mouse perfectly we use this almost everyday
Review-4	bought this on black friday as part of an amazon lightning deal i like using scotch tape products and this product is great nice big pack will last a long time
Review-5	'product works great and i have no issues large pack should last a long time would definitely buy this again

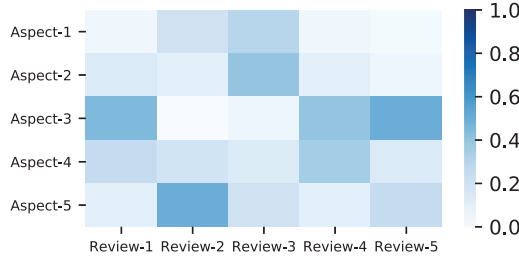


Fig. 6. Visualization of the weights of reviews in multi-aspect review-level attention.

attention mechanisms could help to obtain the multifaceted review representations and user/item representations.

5 CONCLUSION AND FUTURE WORKS

In this article, we propose a multi-aspect neural recommendation model with context-aware personalized attention, named MRCP. We design a word- and review-level context-aware attention network to select the informative words and reviews for users and items under a dynamic manner, since different words and reviews are of different importance. Note that both the two level attentions are multi-aspect to mine more diverse features of reviews and users/items. In addition, we use an aspect-level personalized attention to describe the unique characteristic of users and items. In this way, the comprehensive and accurate representations of users and items can be learned from the rich reviews. Extensive experiments are conducted on six public recommendation datasets from Amazon and Yelp. The results show that our model can effectively improve the performance of recommendation, and consistently outperforms the competitive state-of-the-art baselines. The ablation study indicates the effectiveness of the three-tier attentions, and the case study further provides an intuitive result of the multi-aspect attention weights.

In future, we would like to explore the following directions: (1) considering the different characteristics between users and items during learning their representations instead of using the same modeling architectures, and (2) extending our framework to explore other paradigms in recommender systems beyond rating prediction, and considering more user behaviors in addition to ratings toward better recommendation performance.

REFERENCES

- [1] Yang Bao, Hui Fang, and Jie Zhang. 2014. Topicmf: Simultaneously exploiting ratings and reviews for recommendation. In *Proceedings of the 28th AAAI Conference on Artificial Intelligence*.
- [2] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent dirichlet allocation. *J. Mach. Learn. Res.* 3 (Jan. 2003), 993–1022.
- [3] Rose Catherine and William Cohen. 2017. Transnets: Learning to transform for recommendation. In *Proceedings of the 11th ACM Conference on Recommender Systems*. 288–296.

- [4] Chong Chen, Min Zhang, Yiqun Liu, and Shaoping Ma. 2018. Neural attentional rating regression with review-level explanations. In *Proceedings of the 2018 World Wide Web Conference* (2018), 1583–1592.
- [5] Zhiyong Cheng, Xiaojun Chang, Lei Zhu, Rose C. Kanjirathinkal, and Mohan Kankanhalli. 2019. MMALFM: Explainable recommendation by leveraging reviews and images. *ACM Trans. Inf. Syst.* 37, 2 (2019), 1–28.
- [6] Zhiyong Cheng, Ying Ding, Xiangnan He, Lei Zhu, Xuemeng Song, and Mohan S. Kankanhalli. 2018. A³NCF: An adaptive aspect attention model for rating prediction. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI'18)*. 3748–3754.
- [7] Zhiyong Cheng, Ying Ding, Lei Zhu, and Mohan Kankanhalli. 2018. Aspect-aware latent factor model: Rating prediction with ratings and reviews. In *Proceedings of the 2018 World Wide Web Conference*. International World Wide Web Conferences Steering Committee, 639–648.
- [8] Jin Yao Chin, Kaiqi Zhao, Shafiq Joty, and Gao Cong. 2018. ANR: Aspect-based neural recommender. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. ACM, 147–156.
- [9] Mukund Deshpande and George Karypis. 2004. Item-based top-n recommendation algorithms. *ACM Trans. Inf. Syst.* 22, 1 (2004), 143–177.
- [10] Qiming Diao, Minghui Qiu, Chao-Yuan Wu, Alexander J. Smola, Jing Jiang, and Chong Wang. 2014. Jointly modeling aspects, ratings and sentiments for movie recommendation (JMARS). In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 193–202.
- [11] Kostadin Georgiev and Preslav Nakov. 2013. A non-iid framework for collaborative filtering with restricted boltzmann machines. In *Proceedings of the International Conference on Machine Learning*. 1148–1156.
- [12] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th International Conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 173–182.
- [13] Xiangnan He, Hanwang Zhang, Min-Yen Kan, and Tat-Seng Chua. 2016. Fast matrix factorization for online recommendation with implicit feedback. In *Proceedings of the 39th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 549–558.
- [14] Dietmar Jannach, Paul Resnick, Alexander Tuzhilin, and Markus Zanker. 2016. Recommender systems—Beyond matrix completion. *Commun. ACM* 59, 11 (2016), 94–102.
- [15] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *Proceedings of the 2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 197–206.
- [16] Donghyun Kim, Chanyoung Park, Jinoh Oh, Sungyoung Lee, and Hwanjo Yu. 2016. Convolutional matrix factorization for document context-aware recommendation. In *Proceedings of the 10th ACM Conference on Recommender Systems*. ACM, 233–240.
- [17] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations* (2015).
- [18] Yehuda Koren. 2008. Factorization meets the neighborhood: A multifaceted collaborative filtering model. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 426–434.
- [19] Yehuda Koren, Robert Bell, and Chris Volinsky. 2009. Matrix factorization techniques for recommender systems. *Computer* 42, 8 (2009), 30–37.
- [20] Daniel D. Lee and H. Sebastian Seung. 2001. Algorithms for non-negative matrix factorization. In *Advances in Neural Information Processing Systems*. 556–562.
- [21] Chenliang Li, Cong Quan, Li Peng, Yunwei Qi, Yuming Deng, and Libing Wu. 2019. A capsule network for recommendation and explaining what you like and dislike. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 275–284.
- [22] Greg Linden, Brent Smith, and Jeremy York. 2003. Amazon. com recommendations: Item-to-item collaborative filtering. *IEEE Internet Comput.* 7, 1 (2003), 76–80.
- [23] Guang Ling, Michael R. Lyu, and Irwin King. 2014. Ratings meet reviews, a combined approach to recommend. In *Proceedings of the 8th ACM Conference on Recommender systems*. ACM, 105–112.
- [24] Donghua Liu, Jing Li, Bo Du, Jun Chang, and Rong Gao. 2019. DAML: Dual attention mutual learning between ratings and reviews for item recommendation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 344–352.
- [25] Hongtao Liu, Yian Wang, Qi Yao Peng, Fangzhao Wu, Lin Gan, Lin Pan, and Pengfei Jiao. 2020. Hybrid neural recommendation with joint deep representation learning of ratings and reviews. *Neurocomputing* 374 (2020), 77–85.
- [26] Hongtao Liu, Fangzhao Wu, Wenjun Wang, Xianchen Wang, Pengfei Jiao, Chuhan Wu, and Xing Xie. 2019. NRPA: Neural recommendation with personalized attention. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM.
- [27] Yichao Lu, Ruihai Dong, and Barry Smyth. 2018. Coevolutionary recommendation model: Mutual learning between ratings and reviews. In *Proceedings of the 2018 World Wide Web Conference*. 773–782.

- [28] Julian McAuley and Jure Leskovec. 2013. Hidden factors and hidden topics: Understanding rating dimensions with review text. In *Proceedings of the 7th ACM Conference on Recommender Systems*. ACM, 165–172.
- [29] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*. 3111–3119.
- [30] Andriy Mnih and Ruslan R Salakhutdinov. 2008. Probabilistic matrix factorization. In *Proceedings of the Neural Information Processing Systems (NIPS'08)*. 1257–1264.
- [31] Rong Pan, Yunhong Zhou, Bin Cao, Nathan N. Liu, Rajan Lukose, Martin Scholz, and Qiang Yang. 2008. One-class collaborative filtering. In *Proceedings of the 2008 8th IEEE International Conference on Data Mining*. IEEE, 502–511.
- [32] Steffen Rendle. 2010. Factorization machines. In *Proceedings of the 2010 IEEE International Conference on Data Mining*. 995–1000.
- [33] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian personalized ranking from implicit feedback. In *Proceedings of the 25th Conference on Uncertainty in Artificial Intelligence*. AUAI Press, 452–461.
- [34] Sara Sabour, Nicholas Frosst, and Geoffrey E. Hinton. 2017. Dynamic routing between capsules. In *Advances in Neural Information Processing Systems*. 3856–3866.
- [35] Ruslan Salakhutdinov, Andriy Mnih, and Geoffrey Hinton. 2007. Restricted Boltzmann machines for collaborative filtering. In *Proceedings of the 24th International Conference on Machine Learning*. ACM, 791–798.
- [36] Badrul Munir Sarwar, George Karypis, Joseph A. Konstan, John Riedl, et al. 2001. Item-based collaborative filtering recommendation algorithms. In *Proceedings of the 10th International Conference on World Wide Web* 1, 285–295.
- [37] J. Ben Schafer, Dan Frankowski, Jon Herlocker, and Shilad Sen. 2007. Collaborative filtering recommender systems. In *The Adaptive Web*. 291–324.
- [38] Sungyong Seo, Jing Huang, Hao Yang, and Yan Liu. 2017. Interpretable convolutional neural networks with dual local and global attention for review rating prediction. In *Proceedings of the 11th ACM Conference on Recommender Systems*. ACM, 297–305.
- [39] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* 15, 1 (2014), 1929–1958.
- [40] Xiaoyuan Su and Taghi M. Khoshgoftaar. 2009. A survey of collaborative filtering techniques. *Adv. Artif. Intell.* 2009 (2009).
- [41] Yunzhi Tan, Min Zhang, Yiqun Liu, and Shaoping Ma. 2016. Rating-boosted latent topics: Understanding users and items with ratings and reviews. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI'16)*, Vol. 16. 2640–2646.
- [42] Yi Tay, Anh Tuan Luu, and Siu Cheung Hui. 2018. Multi-pointer co-attention networks for recommendation. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 2309–2318.
- [43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*. 5998–6008.
- [44] Xianchen Wang, Hongtao Liu, Peiyi Wang, Fangzhao Wu, Hongyan Xu, Wenjun Wang, and Xing Xie. 2019. Neural review rating prediction with hierarchical attentions and latent factors. In *Proceedings of the International Conference on Database Systems for Advanced Applications*. Springer, 363–367.
- [45] Libing Wu, Cong Quan, Chenliang Li, Qian Wang, Bolong Zheng, and Xiangyang Luo. 2019. A context-aware user-item representation learning for item recommendation. *ACM Trans. Inf. Syst.* 37, 2 (2019), 1–29.
- [46] Wei Zhang, Quan Yuan, Jiawei Han, and Jianyong Wang. 2016. Collaborative multi-Level embedding learning from reviews for rating prediction. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI'16)*. 2986–2992.
- [47] Lei Zheng, Vahid Noroozi, and Philip S. Yu. 2017. Joint deep modeling of users and items using reviews for recommendation. In *Proceedings of the 10th ACM International Conference on Web Search and Data Mining*. 425–434.
- [48] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1059–1068.

Received July 2020; revised December 2020; accepted December 2020