

MARS: Memory Attention-Aware Recommender System

Lei Zheng, Chun-Ta Lu
Department of Computer Science
University of Illinois at Chicago
Chicago, IL, U. S. A.
lzheng21,clu29@uic.edu

Sihong Xie
Computer Science and Engineering Department
Lehigh University
Bethlehem, PA, U. S. A.
sxie@cse.lehigh.edu

Lifang He
Weill Cornell Department of Healthcare Policy & Research
Cornell University
New York, NY, U. S. A.
lifanghescut@gmail.com

Vahid Noroozi, He Huang and Philip S. Yu
Department of Computer Science
University of Illinois at Chicago
Chicago, IL, U. S. A.
vnoroo2,hehuang,psyu@uic.edu

ABSTRACT

In this paper, we study the problem of modeling users' diverse interests. Previous methods usually learn a fixed user representation, which has a limited ability to represent distinct interests of a user. In order to model users' various interests, we propose a Memory Attention-aware Recommender System (MARS). MARS utilizes a memory component and a novel attentional mechanism to learn deep *adaptive user representations*. Trained in an end-to-end fashion, MARS adaptively summarizes users' interests. In the experiments, MARS outperforms seven state-of-the-art methods on three real-world datasets in terms of recall and mean average precision. We also demonstrate that MARS has a great interpretability to explain its recommendation results, which is important in many recommendation scenarios.

KEYWORDS

Recommender Systems, Deep Learning, Attention

ACM Reference Format:

Lei Zheng, Chun-Ta Lu, Lifang He, Sihong Xie, and Vahid Noroozi, He Huang and Philip S. Yu. 2018. MARS: Memory Attention-Aware Recommender System. In *Proceedings of ACM Conference (Conference'17)*, Jennifer B. Sartor, Theo D'Hondt, and Wolfgang De Meuter (Eds.). ACM, New York, NY, USA, 10 pages. https://doi.org/10.475/123_4

1 INTRODUCTION

How can we accurately model users' interests? It is a fundamental question for building recommender systems (RS). To answer this question, we observed one essential characteristic of users' interests: *diversity*. Users are interested in different kinds of items and interests of users are diverse. For example, a user purchased several books while she or he also bought an electrical gadget. Furthermore, owing to the diversity of a user's interests, only a subset of the user's purchased products can reveal if the user is interested in another product. For instance, a user may purchase an iPad case

Permission to make digital or hard copies of part or all 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 third-party components of this work must be honored. For all other uses, contact the owner/author(s).
Conference'17, July 2017, Washington, DC, USA
© 2018 Copyright held by the owner/author(s).
ACM ISBN 123-4567-24-567/08/06.
https://doi.org/10.475/123_4

because the user bought an iPad rather than a book or a pair of shoes in her last week's shopping list. Hence, it is a nontrivial task to model the diversities of users' interests.

However, how to devise representation vectors to express users' diverse interests is challenging. Existing methods often project users and items into fixed low-dimensional representation vectors in a user-item joint space. We argue that a fixed user representation largely restrains models from accurately modeling users' diverse interests. In the space, similar items are close to each other and the distance between a user and an item implies how much the user is interested in the item. As shown in Figure 1, a user i liked a list of different kinds of items $\{j_1, j_2, j_3, j_4\}$. Because of the diversity of items liked by user i , their representation vectors ($v_{j_1}, v_{j_2}, v_{j_3}$ and v_{j_4}) form two clusters in the space. As a result, the fixed user representation u_i resides between the two clusters in the space. In this case, a system employing fixed user representations will recommend item s to user i instead of item j or item k . Although v_s is closer to u_i than v_j or v_k in the space, it is obvious that either item j or item k is a much more reliable recommendation than item s . One may increase the representation dimensionality to overcome the restriction but will make item representations more scattered in the space and cause a huge increase of model parameters.

Moreover, for most of current deep RS, interpreting its recommendations is also difficult and demanding [12, 26]. Despite the effectiveness of representation vectors for predicting interactions, these vectors reside in latent spaces and are not understandable. In order to enhance transparency and trust in the recommendation process, an increasing need of RS is to provide not only accurate but also interpretable recommendations. This is particularly important in a business-to-business setting, where recommendations are generated for experienced sales staff and not directly for the end-client.

In this paper, in order to model representations of users to express their diverse interests and build a recommendation model capable of interpreting its recommendations, we develop MARS: Memory Attention-Aware Recommender System. First of all, MARS exploits the power of deep learning to learn representations of items from the item content. More importantly, motivated by the observation of user behaviors, we facilitate a memory component and a novel item-level attention mechanism to devise a deep *adaptive user representation*. Unlike fixed user representations, *adaptive user*

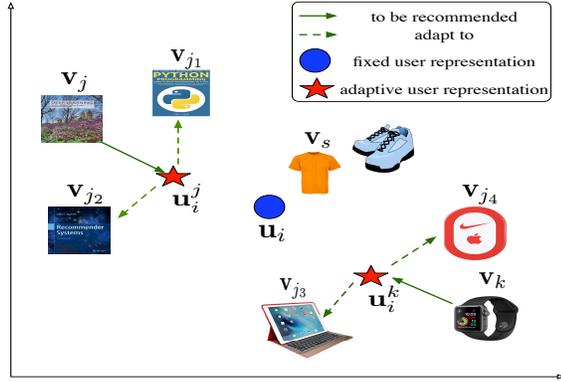


Figure 1: A fixed user representation u_i fails to express user i 's diverse interests. Adaptive user representations dynamically adapt to relevant items.

representations dynamically adapt to locally activated items. As shown in Figure 1, for a candidate item j , an *adaptive user representation* u_i^j dynamically adapts to locally activated items: j_1 and j_2 . To recommend item k , another *adaptive user representation* u_i^k adapts to relevant items like j_3 and j_4 . The concept of *adaptive user representation* can be defined as :

DEFINITION 1. (Adaptive User Representation). For a user i and an item list: $\mathcal{L} = \{j_1, j_2, \dots, j_{n_i}\}$ with n_i items liked by the user, in order to recommend a candidate item j , an *adaptive user representation* dynamically adapts to items in \mathcal{L} which are highly relevant to item j .

The contributions of this work can be summarized as follows:

- **Adaptive User Representations:** To the best of our knowledge, we are the first to introduce a deep end-to-end recommendation model to learn *adaptive user representations*. Especially, a memory component is utilized to capture users' interests in an end-to-end fashion. An attentional mechanism is facilitated to handle the diversities of users' interests.
- **An Interpretable Model:** Benefiting from the item-level attention mechanism, MARS has a good interpretability to explain why an item gets recommended to a user by showing relevant items liked by the user.
- **Strong Performance:** In the experiments, we demonstrate that MARS can significantly outperform strong baselines on three real-world datasets.

2 BACKGROUND AND PRELIMINARIES

In this section, we present the background and preliminaries of this study. Throughout the paper, we denote scalars by either lowercase or uppercase letters, vectors by boldfaced lowercase letters, and matrices by boldfaced uppercase letters.

We consider the most common scenario with implicit feedbacks. For implicit feedbacks, only positive observations, such as clicks, purchase or likes, are available. The non-observed user-item pairs, e.g. a user has not bought an item yet, are a mixture of real negative

feedback (the user is not interested in buying the item) and missing values (the user might want to buy the item in the future).

Although, in different recommendation scenarios, the item content can be in distinct formats, such as text or images, we limit our study to recommend items associated with the textual content, which usually describe important characteristics of items. However, note that the proposed model can be easily generalized to recommend different types of items, such as songs or videos.

Let us denote a set of users \mathcal{U} and an item set \mathcal{I} . Each item $j \in \mathcal{I}$ is associated with a document $Y^j = \{w_1, w_2, \dots, w_{n_j}\}$ containing n_j words. For a user $i \in \mathcal{U}$, let $\mathcal{I}_i^+ = \{I_1, I_2, \dots, I_{n_i}\}$ denote the set of n_i items liked by user i and \mathcal{I}_i^- denote the remaining items. Furthermore, all remaining items liked by user i except item j , denoted as \mathcal{I}_i^+ / j , are utilized to model user representations.

Finally, we define the recommendation problem which we study in this paper as follows:

DEFINITION 2. (Problem Definition). Given user and item sets: \mathcal{U} and \mathcal{I} , for each user $i \in \mathcal{U}$, we aim to recommend a ranked list of items from \mathcal{I}_i^- that are of interests to the user.

3 PROPOSED MODEL

The overall architecture of MARS is shown as in Figure 3. Inspired by the recent success of Convolutional Neural Network (CNN) [15] in various tasks [13, 14], we first facilitate two CNN models to serve for different purposes. One CNN model, denoted as $f_{user}(\cdot; \Psi)$ and parameterized by Ψ , is designed to learn a memory component C from documents associated with items liked by user i . The other CNN model, denoted as $f_{item}(\cdot; \Omega)$ and parameterized by Ω , is responsible for learning an representation v_j of item j . Later, a novel item-level attention mechanism is proposed to build on the top of the memory component to learn a deep *adaptive user representation* u_i^j . At last, with u_i^j and v_j , user i 's preference score of item j can be derived.

3.0.1 CNN Models. In this subsection, since $f_{user}(\cdot; \Psi)$ and $f_{item}(\cdot; \Omega)$ only differ in their inputs, we focus on illustrating the process of $f_{item}(\cdot; \Omega)$ in detail. The same process is applied for $f_{user}(\cdot; \Psi)$ with similar layers. The architecture of CNN used in this paper is shown in Figure 2.

Given an item j and its associated document $Y^j = \{w_1, w_2, \dots, w_{n_j}\}$, in order to make use of complex lexical semantics within Y^j , we first utilize a word embedding layer to transform the document into a dense numeric matrix Π . Formally, the word embedding layer is defined as $\mathcal{H} : \mathcal{V} \rightarrow \mathcal{R}^e$, where \mathcal{V} represents the dictionary of words. We map each word in Y^j into an e dimensional vector as:

$$\Pi = \mathcal{H}(w_1) \oplus \mathcal{H}(w_2) \oplus \dots \oplus \mathcal{H}(w_{n_j}), \quad (1)$$

where $\Pi \in \mathcal{R}^{e \times n_j}$, each column corresponds to a vector for a word in the document Y^j . Note that Π is randomly initialized and is further trained through the optimization process.

Following the word embedding layer, a convolutional layer is built to extract important contextual features that can represent items. A convolution layer consists of g neurons in total, each of which applies convolution operator on Π as:

$$\mathbf{z} = \text{ReLU}(\Pi * \mathcal{K} + b). \quad (2)$$

Here symbol $*$ is the convolution operator, $\mathcal{K} \in \mathcal{R}^{e \times c}$ and $b \in \mathcal{R}$ are a convolutional filter and a bias term, respectively, and c denotes the window size of the convolutional filter. We use *ReLU* [20] as our activation function.

After the convolutional layer, $\mathbf{z} \in \mathcal{R}^{(n_j-c+1) \times 1}$ becomes a vector consisting of $n_j - c + 1$ contextual features captured by the convolutional kernel \mathcal{K} . To extract the most important feature value from \mathbf{z} , we apply a max-pooling operation on \mathbf{z} as shown in Figure 2. With g neurons in the convolutional layer, we obtain a column vector \mathbf{s}_j consisting of g important contextual features as:

$$\mathbf{s}_j = \{\max(\mathbf{z}_1), \max(\mathbf{z}_2), \dots, \max(\mathbf{z}_g)\}. \quad (3)$$

At last, contextual features \mathbf{s}_j of item j is projected to a K -dimensional space through a fully connected layer as:

$$\mathbf{v}_j = \tanh(\hat{\mathbf{W}}\mathbf{s}_j + \hat{\mathbf{b}}), \quad (4)$$

where $\hat{\mathbf{W}} \in \mathcal{R}^{K \times g}$ and $\hat{\mathbf{b}} \in \mathcal{R}$. Since the document \mathbf{Y}^j is associated with item j and characterizes item j , $\mathbf{v}_j \in \mathcal{R}^{K \times 1}$ can be regarded as a deep representation of item j .

3.1 Memory Component

A set of items liked by a user naturally reflects the user's interests. However, it is not easy to summarize diverse items liked by user i into a user representation. As discussed in the introduction section, a fixed user representation fails to convey diverse interests of a user.

Different from previous methods [4, 36], in order to learn a deep and accurate user representation, as shown in Figure 3, we first propose to first utilize $f_{user}(\cdot; \Psi)$ to learn a memory component for items in \mathcal{I}_i^+ / j as:

$$\begin{aligned} \mathbf{C} &= \{f_{user}(\mathbf{Y}^1; \Psi), f_{user}(\mathbf{Y}^2; \Psi), \dots, f_{user}(\mathbf{Y}^{n_i-1}; \Psi)\} \\ &= \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{n_i-1}\}, \end{aligned} \quad (5)$$

Each column of $\mathbf{C} \in \mathcal{R}^{K \times (n_i-1)}$ corresponds to a representation of an item in \mathcal{I}_i^+ / j . Hence, the matrix \mathbf{C} characterizes the user's interests.

3.1.1 Item-level Attention. Although the memory component \mathbf{C} stores all item information of user i , our goal is to learn a deep representation of user i based on \mathbf{C} . As discussed in the Introduction section, items in \mathcal{I}_i^+ / j are diverse and only a subset of \mathcal{I}_i^+ / j is relevant to item j . Motivated by the intuition above, we introduce an item-level attention mechanism on the memory component to capture items in \mathcal{I}_i^+ / j relevant to item j . To do so, as shown in Figure 3, for a given item j , we first feed its corresponding document \mathbf{Y}^j into $f_{item}(\cdot; \Omega)$ to derive an item representation as:

$$\mathbf{v}_j = f_{item}(\mathbf{Y}^j; \Omega), \quad (6)$$

where $\mathbf{v}_j \in \mathcal{R}^{K \times 1}$ is a column vector serving as a representation of item j .

More items in \mathcal{I}_i^+ / j with high relevance scores or locally activated to item j mean that the user i is more likely interested in j . Thus, with the memory component \mathbf{C} and \mathbf{v}_j , we measure the relevance score or activation degree between \mathbf{v}_j and each item representation in $\mathbf{C} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{n_i-1}\} \in \mathcal{R}^{K \times (n_i-1)}$ by taking the inner product followed by a softmax as:

$$\boldsymbol{\alpha}^{ij} = \text{softmax}(\mathbf{C}^\top \mathbf{v}_j). \quad (7)$$

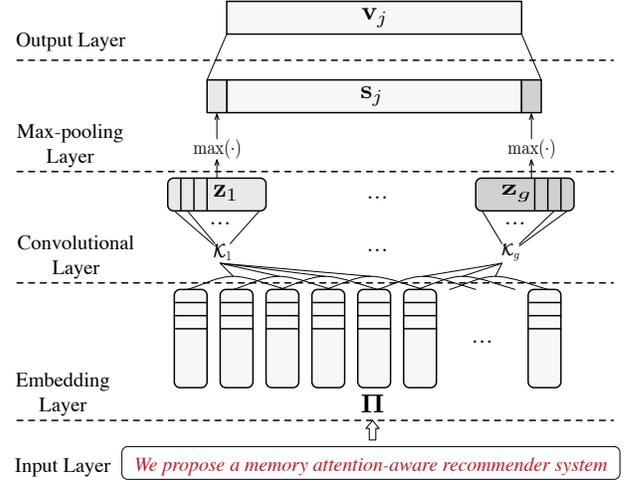


Figure 2: Our CNN architecture used for learning item representations.

Defined in this way, $\boldsymbol{\alpha}^{ij} \in \mathcal{R}^{(n_i-1) \times 1}$ is an attention column vector of user i for item j . A larger value in $\boldsymbol{\alpha}^{ij}$ indicates higher relevance between item j and a corresponding item in \mathcal{I}_i^+ / j . A larger value means a higher weight for deriving interests of user i for item j .

3.1.2 Deep Adaptive User Representations. A fixed user representation is unable to express diverse interests of the user. In order to uncover a user's interests in a candidate item j , we introduce a deep *adaptive user representation*. The proposed user representation is dependent on the candidate item j and non-fixed. By focusing on items in \mathcal{I}_i^+ / j which are highly activated, it is able to model users' diverse interests.

With the matrix \mathbf{C} and the attention vector $\boldsymbol{\alpha}^{ij}$, a deep attention-aware *adaptive user representation* is formed by a weighted sum of $\mathbf{C} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{n_i-1}\}$ as:

$$\mathbf{u}_i^j = \mathbf{C} \boldsymbol{\alpha}^{ij}. \quad (8)$$

In Eq. (8), the attention vector $\boldsymbol{\alpha}^{ij}$ weights each item representation in \mathbf{C} . An *adaptive user representation* \mathbf{u}_i^j is obtained by adaptively focusing on items in \mathcal{I}_i^+ / j activated by item j .

3.2 Pair-wise Learning and Prediction

Recall that an *adaptive user representation* \mathbf{u}_i^j is obtained by placing an attentional vector on our memory component. Given the representation of a candidate item \mathbf{v}_j , one can compute the user i 's preference score over item j as:

$$r_{ij} = \mathbf{u}_i^{j\top} \mathbf{v}_j. \quad (9)$$

Here, we derive a preference score r_{ij} that reflects the preference of the user for the item. Usual approaches for item recommendations are to rank items by sorting them according to the scores. However, these approaches treat recommendation tasks as regression problems and are not optimized for ranking. To train MARS optimized for ranking, inspired by BPR [22], we formalize the training data

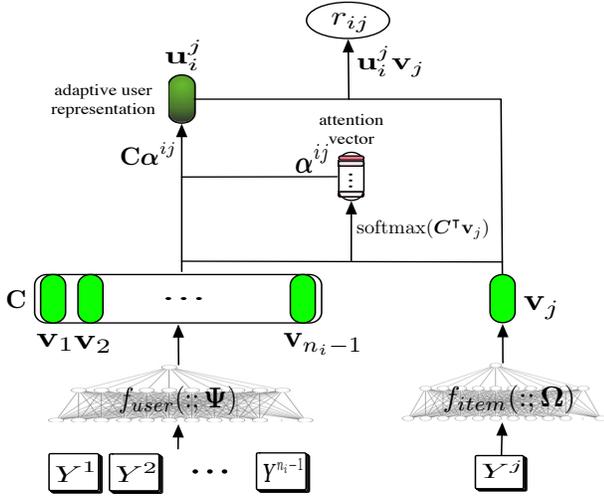


Figure 3: The architecture of MARS.

\mathcal{D} as:

$$\mathcal{D} = \{(i, I_i^+ / j, j, j') | i \in \mathcal{U} \wedge j \in I_i^+ \wedge j' \in I_i^-\}, \quad (10)$$

where i, j and j' are uniformly sampled from \mathcal{U}, I_i^+ and I_i^- , respectively. Note that for each quadruple $(i, I_i^+ / j, j, j')$ in \mathcal{D} , I_i^+ / j is a different sampled subset of I_i^+ . The sampling strategy is a common practice for pair-wise recommendation learning. [22] also shows that training on all item pairs can result in slow and poor convergence.

Similar to BPR, instead of scoring single items, we use item pairs to model a user u 's preference of item j over item j' as:

$$\begin{aligned} \mathcal{L} = & -\frac{1}{|\mathcal{D}|} \sum_{(i, I_i^+ / j, j, j') \in \mathcal{D}} \{\ln(\sigma(\mathbf{u}_i^{jT} \mathbf{v}_j - \mathbf{u}_i^{j'T} \mathbf{v}_{j'})) \\ & + \lambda_u \mathbf{u}_i^{jT} \mathbf{u}_i^j + \lambda_u \mathbf{u}_i^{j'T} \mathbf{u}_i^{j'} + \lambda_v \mathbf{v}_j^T \mathbf{v}_j + \lambda_v \mathbf{v}_{j'}^T \mathbf{v}_{j'}\}, \end{aligned} \quad (11)$$

where \mathbf{u}_i^j and $\mathbf{u}_i^{j'}$ denote *adaptive user representations* of user i for item j and item j' ; \mathbf{v}_j and $\mathbf{v}_{j'}$ stand for representations of item j and item j' ; λ_u and λ_v are regularization terms. The sigmoid function σ maps user i 's preference score of item j over item j' into probabilities.

MARS is trained by minimizing Eq. (11). The derivatives of parameters in different layers can be computed by applying differentiation chain rule [23]. We optimize the model through RMSprop [31] over a batch of tuple $\{i, I_i^+ / j, j, j'\}$.

Overall, the training procedure of MARS is summarized in Algorithm 1. During the test, given a user i 's past likes I_i^+ , the final item recommendation list for user i is given according to the following ranking criterion:

$$i : j_1 \geq j_2 \geq \dots \geq j_n \Rightarrow r_{i,j_1} > r_{i,j_2} > \dots > r_{i,j_n}. \quad (12)$$

4 EXPERIMENTS

4.1 Datasets

We evaluate the proposed method on three real-world datasets with different kinds of items as the following:

Algorithm 1: Training procedure of MARS

Input: Training set:

$\mathcal{D} := \{(i, I_i^+ / j, j, j') | i \in \mathcal{U} \wedge j \in I_i^+ \wedge j' \in I_i^-\}$, number of epochs T , batch size m , window size c , number of convolutional neurons g

Output: Model's parameter set: $\Theta = \{\mathbf{u}, \mathbf{v}, \Omega\}$

for $t = 1, 2, \dots, T$ **do**

Generate the t^{th} batch of size m by uniformly sampling from \mathcal{U} ,

I_i^+ and I_i^- ;

Calculate the memory component C according to Eq. (5);

Calculate \mathbf{v}_j and $\mathbf{v}_{j'}$ according to Eq. (6);

Calculate α^{ij} and $\alpha^{ij'}$ according to Eq. (7);

Calculate \mathbf{u}_i^j and $\mathbf{u}_i^{j'}$ according to Eq. (8);

Calculate \mathcal{L} according to Eq. (11);

Estimate gradients $\frac{\partial \mathcal{L}}{\partial \Theta_t}$ by back propagation;

Calculate Θ_{t+1} with RMSprop [31];

end

return Θ_T ;

- **Yahoo! Movies:** it consists of users rating movies on a scale of 1-5 with a short synopsis. To be consistent with the implicit feedback setting, as in [8], we extract only positive ratings (rating 5) for training and testing. After removing movies without a synopsis and users with less than 3 ratings, we obtain a dataset that contains 7,642 users, 11,915 items, and 221,367 positive ratings with 0.24% density.
- **Amazon Video Games:** it is collected by [6, 19] and contains game descriptions and ratings from 1 to 5. Similarly, we transformed it into implicit data, where each entry is marked as 0 or 1 indicating whether the user has rated the games. After removed games with descriptions and users with less than 10 ratings, the resulting dataset contains 47,063 video games and 2,670 users with 0.037% density.
- **Amazon Movies and TV:** this dataset is also collected by [6, 19]. Similarly, items in this dataset are movies and TV shows available on *Amazon.com*. Users rate them from 1 to 5. In order to transform explicit ratings into implicit feedbacks, we only preserve ratings with the value of 5 and treat them as positive feedbacks. Other entries are marked as negative ones. By removing users with less than 10 ratings, we obtain a dataset with 22,147 users, 178,086 items, and 0.0128% density. Note that, compared with the previous two datasets, this is a much harder dataset due to its sparsity.

The synopses and descriptions of items serve as item contents. After removing stop words and frequencies of the word less than 5, the vocabulary sizes of *Yahoo! Movies*, *Amazon Video Games* and *Amazon Movies and TV* are 33,195, 25,035 and 68,919, respectively. All three datasets are publicly available ¹.

4.2 Baselines

To validate the effectiveness of MARS, we compare MARS with seven state-of-the-art baseline models. Among them, BPR and NCF completely ignore the textual content associated with items and

¹ *Yahoo! Movies* can be downloaded at <https://webscope.sandbox.yahoo.com/>; *Amazon Video Games* and *Amazon Movies and TV* are available at <https://snap.stanford.edu/data/web-Amazon.html>

all other baselines utilize the content information to boost their performances.

- **BPR** [22]²: We use **B**ayesian **P**ersonalized **R**anking based Matrix Factorization, which is based on users' pair-wise preference, as the single collaborative filtering method. BPR completely ignores the usage of item content.
- **NCF** [8]³: **N**eural **C**ollaborative **F**iltering fuses matrix factorization and Multi-Layer Perceptron (MLP) to learn from user-item interactions. The MLP endows NCF with the ability of modelling non-linearities.
- **CMF** [29]⁴: **C**ollective **M**atrix **F**actorization simultaneously factorizes multiple matrices to incorporate different sources of information. We factorize the rating matrix and a matrix consisting of *bag-of-words* features. Note that we follow the work of [11] to adapt CMF for implicit feedbacks.
- **CTR** [34]⁵: **C**ollaborative **T**opic **R**egression is based on topic modeling techniques and shows very good performance on recommending articles.
- **DeepFM** [28]⁶: This baseline combines the power of factorization machine [21] for recommendations and deep learning for feature learning in a neural network. We concatenate a user id and *bag-of-words* features of each item for training and predictions.
- **CDL** [35]⁷: **C**ollaborative **D**eep **L**earning is a recently proposed deep recommender system. It tightly couples a Bayesian formulation of the stacked denoising auto-encoders and Probabilistic Matrix Factorization (PMF) [24]. The middle layer of auto-encoders serves as a bridge between auto-encoders and probabilistic matrix factorization.
- **Wide & Deep** [4]⁸: This model is proposed to jointly train wide linear models and deep neural networks to combine the benefits of memorization and generalization for recommender systems. Similar to DeepFM, a user id and *bag-of-words* features of each item are combined to be fed into both the "wide" and "deep" parts of the model.

All baseline models can be categorized into three groups: (1) **BPR** and **NCF**: these two models learn only from users' implicit feedback and ignore the textual content; (2) **CMF** and **CTR**: this group includes two "shallow" recommendation models; (3) **DeepFM**, **CDL** and **Wide & Deep**: three state-of-the-art deep learning based models are included to be compared with MARS.

4.3 Experimental Setup

Following [36], we evaluate the models on held-out user-item likes. For each dataset, to evaluate MARS and baselines in a sparse setting, we randomly select only 30% items associated with each user to form the training set. All the remaining are split evenly to serve as the validation and test set. We repeat the evaluation five times with different randomly selected training sets. The average performances are reported in the following sections. For each dataset, to achieve

the best performance for each individual model, we conducted carefully parameter study in Section 4.4.

Since we expect RS to not only be able to retrieve relevant items out of all available items but also provide a ranking where items of users' interests are ranked in the top. Therefore, for evaluation, we use two metrics to evaluate the proposed model and the baselines. The first metric we use is recall@N. The recall is often used to measure how well a model can retrieve relevant items out of all available items. The recall@N for each user is then defined as:

$$\text{recall@N} = \frac{\# \text{ items the user likes among the top N}}{\text{total number of items the user likes}}. \quad (13)$$

Another evaluation metric we use is Mean Average Precision (MAP). The MAP is employed to measure the ranking performances of MARS and baselines. The definition of MAP is given as:

$$\text{MAP} = \frac{1}{|\mathcal{U}|} \sum_{i \in \mathcal{U}} \text{AveP}(i),$$

$$\text{AveP}(i) = \frac{1}{|K'|} \sum_{k=1}^{K'} P_i(k) \text{rel}_i(k), \quad (14)$$

where $P_i(k)$ represents the precision of the top k products recommended to user i ; $\text{rel}_i(k)$ denotes whether the k_{th} item has interacted with user i in the test set. Similar to [32], we set the cut-off point K' as 500 for each user.

4.4 Hyper-Parameter Study

Although different models have different hyper-parameters, some hyper-parameters are common and play important roles on model performances. In this section, we optimize the performances of MARS and the baselines by studying the impacts of several important hyper-parameters on the validation sets of *Yahoo! Movies* and *Amazon Video Games*. The learning rate and batch size are empirically set as 0.001 and 512 for MARS. All the other hyper-parameters for baselines follow the original papers.

4.4.1 Dimension of Latent Vectors. A low-dimensional latent vector of users and items has a limitation of modeling complex user-item interactions. However, a high-dimensional vector may harm the generalization of the model and increases the number of parameters. In order to optimize the performances of MARS as well as the baselines, we conduct an experiment to investigate the impacts of the dimension of latent vectors. The dimension of latent vectors is searched from [5, 10, 20, 50, 70, 100] via validation sets from *Yahoo! Movies* and *Amazon Video Games*. As shown in Figure 4, different models reach their own best performances at different dimensions of latent vectors. MARS achieves its best performances when its dimensions are set to 50 and 70 on the dataset of *Yahoo! Movies* and *Amazon Video Games*, respectively.

4.4.2 Regularization Terms. In order to combat the over-fitting problem, many models place \mathcal{L}_2 regularization terms (λ_u and λ_v) on the representation vectors of users and items. We search the two hyper-parameters from [0.001, 0.002, 0.005, 0.01, 0.02, 0.05] to optimize performances of MARS and the baselines. In Figure 5 and 6, we can see that MARS reaches its best performances when both λ_u and λ_v are set to 0.002 on the two datasets. Note that instead

²Code: <https://github.com/zenogantner/MyMediaLite>

³Code: https://github.com/hexiangnan/neural_collaborative_filtering

⁴Code: <https://github.com/david-cortes/cmfrac>

⁵Code: <https://github.com/blei-lab/ctr>

⁶Code: <https://github.com/ChenglongChen/tensorflow-DeepFM>

⁷Code: <http://www.wanghao.in/code/cdl-release.rar>

⁸Code: https://www.tensorflow.org/tutorials/wide_and_deep

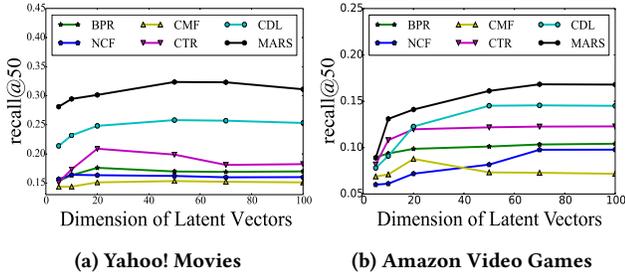


Figure 4: The dimension of latent vectors of users and items is varied from 5 to 100 to investigate its impacts on the performances.

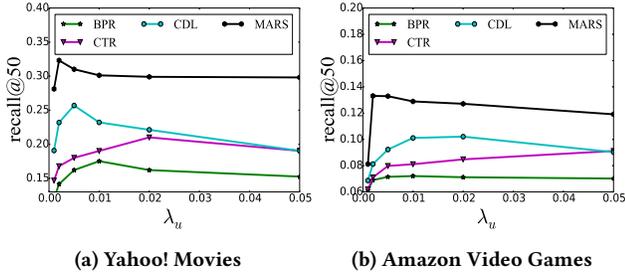


Figure 5: The \mathcal{L}_2 regularization term λ_u is varied from 0.001 to 0.05 to investigate its impacts on the performances.

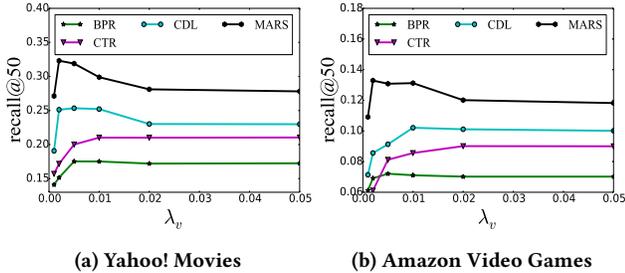


Figure 6: The \mathcal{L}_2 regularization term λ_v is varied from 0.001 to 0.05 to examine its impacts on the performances.

of \mathcal{L}_2 regularization, NCF and DeepFM employ Dropout [30] to prevent over-fitting. And, Wide & Deep uses the \mathcal{L}_1 regularization.

4.4.3 Network Architecture. To find a good network shape for MARS, we also investigate two hyper-parameters: e and g , each of which denotes the dimension of word embeddings and the number of convolutional filters, respectively. We perform a grid search on the two hyper-parameters on the validation set of *Yahoo! Movies*. We found MARS can achieve good performances when we set $e = 300$ and $g = 64$. For NCF, as suggested in the original paper [8], we employ a three-layer MLP with the shape of (32, 16, 8). In Wide & Deep, a three-layer MLP with the shape of (1024, 512, 256) is used in the "deep" part. In DeepFM, as in [28], we build a MLP network with the shape of (200, 200, 200).

4.4.4 Other Hyper-Parameters. For CMF, *bag-of-words* vectors as in [35] for each item forms a matrix to be simultaneously factorized with the rating matrix. By performing the grid search, we optimize the performances of CMF by setting the weights for the

rating matrix and content matrix to 5 and 1. α of CTR is set as 0.1. To optimize the performance of CDL, we perform a grid search on the hyper-parameters: λ_n , λ_w and L .

4.5 Experimental Results

Table 1 illustrates the experimental results for MARS and seven state-of-the-art baselines on all three datasets. Overall, MARS improves the best baseline by **20.8%** and **13.0%** in terms of recall@50 and MAP, respectively, averaging on all three datasets.

These experiments reveal a number of interesting points:

- Regardless of the data sets and the evaluation metrics, our MARS always achieves the best performance. This shows that by leveraging the power of *adaptive user representations*, MARS can better model users' diverse interests, resulting in better recommendations.
- Besides CMF, models considering the textual content generally give better results than models ignoring the textual content. It validates the usefulness of the textual content.
- Deep models perform better than "shallow" ones in our experiments. This observation indicates the advantages of deep features extracted by various deep models.
- Among all baseline models, Wide & Deep shows itself as a strong performer and defeats all other baselines in all three datasets in terms of recall@50 and MAP. It calls us to combine a deep model with a "shallow" one for improvements.
- Because of the sparsity differences, we also observe that the performances of all models degrade as the dataset become sparser.

5 MODEL ANALYSIS

In this section, we conduct experiments to answer the following questions:

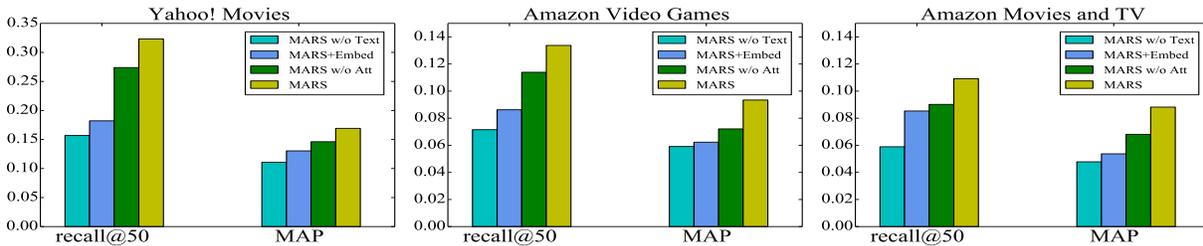
- Q1:** How much does MARS benefit from taking the textual content associated with items into consideration?
Q2: How much do the incorporated CNN models help MARS?
Q3: Do the proposed *adaptive user representations* and attention mechanism assist MARS in achieving better performances?

In order to answer questions above, we include three variants of MARS: MARS w/o Text, MARS+Embed and MARS w/o Att as the following:

- **MARS w/o Text:** To answer **Q1**, instead of learning item embeddings from the textual content, we ignore the content information and randomly initialize item embeddings based on a Gaussian distribution. After the initialization, embeddings of items are optimized during the training.
- **MARS+Embed:** In order to answer **Q2**, in this variant, we replace the CNN models by simply averaging on embeddings of words contained each document associated with every item.
- **MARS w/o Att:** For **Q3**, we include a variant of MARS without the proposed attention mechanism. To do so, we fix all values of α as one. In this way, for a user i , the attention vector is canceled and representation of user i becomes fixed and is not adaptive to relevant items liked by user i .

Table 1: Performance comparison with baselines. The best performance is indicated in bold (higher is better). The standard deviation is shown in parentheses.

Dataset		BPR	NCF	CMF	CTR	DeepFM	CDL	Wide & Deep	MARS	MARS vs. best
Yahoo! Movies	recall@50	0.1756 (0.001)	0.1649 (0.008)	0.1532 (0.002)	0.2067 (0.002)	0.2416 (0.002)	0.2534 (0.001)	0.2611 (0.001)	0.3230 (0.001)	23.7%
	MAP	0.1045 (0.001)	0.0973 (0.002)	0.0894 (0.003)	0.1223 (0.002)	0.1389 (0.001)	0.1452 (0.001)	0.1523 (0.002)	0.1692 (0.002)	11.1%
Amazon Video Games	recall@50	0.0716 (0.001)	0.0849 (0.002)	0.0528 (0.002)	0.0827 (0.003)	0.1086 (0.001)	0.1034 (0.004)	0.1123 (0.001)	0.1337 (0.002)	19.1%
	MAP	0.0625 (0.002)	0.0654 (0.001)	0.0517 (0.002)	0.0658 (0.001)	0.0815 (0.004)	0.0746 (0.005)	0.0821 (0.003)	0.0934 (0.003)	13.8%
Amazon Movies and TV	recall@50	0.0643 (0.001)	0.0604 (0.002)	0.0496 (0.002)	0.0711 (0.002)	0.0935 (0.001)	0.0901 (0.001)	0.1001 (0.003)	0.1196 (0.002)	19.5%
	MAP	0.0543 (0.001)	0.0551 (0.001)	0.0493 (0.002)	0.0659 (0.004)	0.0771 (0.002)	0.0746 (0.004)	0.0785 (0.001)	0.0895 (0.002)	14.0%

**Figure 7: Performance comparison with variants of MARS on the three datasets in terms of recall@50 and MAP.**

As shown in Figure 7, MARS w/o Text performs the worst, regardless of datasets and evaluation metrics. Compared with MARS w/o Text, MARS+Embed gains improvements in all three datasets. It further validates the advantages of incorporating the textual content for RS. However, averaging on embeddings of words associated with items is not an ideal method to make use of information existing in the textual content. Powered by deep features extracted by the CNN models from the textual content, MARS w/o Att achieves additional improvements over MARS+Embed. In terms of MAP, MARS w/o Att improves MARS+Embed by 12.0%, 15.73% and 26.82% on the dataset of *Yahoo! Movies*, *Amazon Video Games* and *Amazon Movies and TV*, respectively. However, MARS defeats MARS w/o Att in terms of recall@50 and MAP on all three datasets. This demonstrates that, with the proposed attention mechanism and the incorporated CNN models, MARS is able to learn an effective *adaptive user representation* and leverage the rich information existing in the textual content.

6 A CASE STUDY ON THE INTERPRETABILITY OF MARS

To gain a better insight into the interpretation ability of MARS, we conduct a qualitative experiment in this section. We run MARS on the dataset of *Yahoo! Movies* and examine two example users. For each of them, we show two recommended movies in the second column of Table 2. And the top three movies with the highest attention values shown in the third column explain why corresponding

movies in the second column are recommended by MARS. For example, MARS recommends *Sleepless in Seattle* to *User I* because movies, such as *Where the Heart Is* and *When Harry met Sally...*, are locally activated by high attention values. Besides *romance* movies, MARS also discovers that *User I* is interested in *action* movies because of his or her past interactions with *X2: X-Men United* and *Top Gun*. As a result, MARS recommends *Enemy at the Gates*. For *User II*, MARS recommends two movies of different genres: *The Lord of the Rings: The Two Towers* and *Kill Bill Vol. 1 (2003)*. It is because movies watched by *User II*, such as *Pirates of the Caribbean: The Curse of the Black Pearl* and *Donnie Brasco (1997)*, indicate the user is a fan of *fantasy* and *action* movies. Overall, this case study shows that MARS is not only able to capture users' diverse interests but also has a great interpretability to explain its recommendation results.

7 RELATED WORK

Our work is closely related to two areas: deep learning based RS and attentional mechanisms employed for RS. We will first give a brief review on deep learning based RS. Then, we covers works utilizing attentional mechanisms for RS.

7.1 Deep Recommender Systems

Several studies [18, 37, 47] recently propose deep learning based models for the recommendation tasks. One pioneer work [25] in this area uses a Restricted Boltzmann Machines (RBM) [9] based method to model users using their rating preferences. Followed by this trend, [41] utilize denoising Auto-encoders to learn latent vectors of users

Table 2: A case study on the interpretability of MARS. The second column displays the top 2 movies recommended by MARS to *User I* and *User II*. The third column includes the top three movies with highest attention values. The actual attention value is shown in parentheses.

	Recommended Movies	Because you watched
<i>User I</i>	<i>Sleepless in Seattle</i>	1. <i>Where the Heart Is</i> (0.03) 2. <i>When Harry Met Sally...</i> (0.02) 3. <i>Ronnie & Julie</i> (0.02)
	<i>Enemy at the Gates</i>	1. <i>X2: X-Men United</i> (0.02) 2. <i>Top Gun</i> (0.015) 3. <i>The Matrix Reloaded</i> (0.014)
<i>User II</i>	<i>The Lord of the Rings: The Two Towers</i>	1. <i>Pirates of the Caribbean: The Curse of the Black Pearl</i> (0.03) 2. <i>Harry Potter And the Chamber of Secrets</i> (0.02) 3. <i>Harry Potter and the Sorcerer's Stone</i> (0.02)
	<i>Kill Bill Vol. 1</i>	1. <i>Donnie Brasco</i> (0.03) 2. <i>The Italian Job</i> (0.021) 3. <i>S.W.A.T.</i> (0.02)

and items from the rating matrix. [8] and [7] leverage Multilayer Perceptron (MLP) to learn from user-item interactions. In [47], a CF-based Neural Autoregressive Distribution Estimator (CF-NADE) model is proposed for collaborative filtering tasks. However, all works above differ from MARS because they ignore the rich content associated with items.

Some studies also propose to utilize deep learning techniques to build a content-based recommender system. In [32, 39], authors introduce Convolutional Neural Networks (CNN) [15] and Deep Belief Network (DBN) [10] to learn users' preferences from music data. [1, 35, 36] propose a deep recommendation model learning from the textual content associated with items. [5, 38] propose deep recommender systems for video and point-of-interest recommendation. [43, 45, 46] investigate how to leverage the multi-view information to improve the quality of recommender systems. In [16], they propose a model which is able to simultaneously predict ratings and generate abstractive tips. For more works on the deep learning based RS, readers can refer to a survey paper [44].

7.2 Attentional Mechanisms for RS

Recently, attentional mechanisms attract considerable interests of researchers, owing to their ability of modeling users' attention. A number of works [33] employing attentional mechanisms to build RS have been proposed. [2] introduces an attentional mechanism utilizing reviews. In [42], authors propose a neural attention network to model the importance of each feature interaction from data. [17] proposed an attentional mechanism based model to learn dynamics of user interests from a sequence of items. To build an interpretable recommendation model, [27] proposes to utilize an attentional method to extract text from reviews. For the purpose of capturing editors' dynamic criteria for selecting news articles, [40] proposes a dynamic attentional method.

To the best of our knowledge, two methods close to ours are presented in [3] and [48]. Although the two works proposed to place attentional mechanisms on items, none of them are end-to-end models. It means that they use either hand-crafted or pre-trained features. As a result, item features are not optimized for the task of recommendation. It leads to an inaccurate attentional mechanism. In contrast, MARS learns item features and users' attention in

an end-to-end fashion. Therefore, compared with them, MARS achieves better item features and more accurate attention of users.

In terms of interpretability, a deep recommender system with interpretability is presented in [27]. However, they interpret recommendations based on reviews written by users. In contrast, we interpret recommendations based on items purchased/liked by users. In fact, interpreting recommendations based on purchased items is more popular in real life (as *Amazon* or *Netflix* does).

Although all neural network based approaches above are deep content-based recommender systems, they differ from MARS because all of works above are unable to learn an *adaptive user representation* in an end-to-end fashion.

8 CONCLUSIONS

Although deep learning based RS have shown promising results in a variety of recommendation tasks, previous methods often focus on utilizing deep models for modeling item representations and learning a fixed user representation. However, a fixed user representation restrains models from modeling diverse interests of users. Furthermore, while interpretability is demanded in many recommendation scenarios, most of existing deep learning based RS are unable to interpret their recommendation results.

In this paper, to tackle the problems and challenges above, we present a Memory Attention-aware Recommender System (MARS) model. With a proposed memory component and an item-level attention mechanism, instead of modeling fixed deep user representations, MARS learns a deep *adaptive user representation*. For an item j and a set of items liked by user i , a deep *adaptive user representation* can dynamically adapt to those items in the set which are relevant item j . Owing to its adaptability, a deep *adaptive user representation* can overcome the difficulty of modeling users' diverse interests. Moreover, with the help of the proposed attention mechanism, MARS is able to interpret its recommendation results based on purchased items of users.

In the experiments, we demonstrate that MARS achieves superior performances by comparing with seven state-of-the-art methods on three real-world datasets. Also, we demonstrate that MARS can not only overcome the difficulty of modeling diverse interests of users but also has a great interpretability.

REFERENCES

- [1] Trapit Bansal, David Belanger, and Andrew McCallum. 2016. Ask the gru: Multi-task learning for deep text recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems*. ACM, 107–114.
- [2] 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 on World Wide Web*. International World Wide Web Conferences Steering Committee, 1583–1592.
- [3] Jingyuan Chen, Hanwang Zhang, Xiangnan He, Liqiang Nie, Wei Liu, and Tat-Seng Chua. 2017. Attentive Collaborative Filtering: Multimedia Recommendation with Item-and Component-Level Attention. In *Proceedings of the 40th International ACM SIGIR conference on Research and Development in Information Retrieval*.
- [4] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishikesh Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispr, et al. 2016. Wide & deep learning for recommender systems. In *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems*. ACM, 7–10.
- [5] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems*. ACM, 191–198.
- [6] Ruining He and Julian McAuley. 2016. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *Proceedings of the 25th International Conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 507–517.
- [7] Xiangnan He and Tat-Seng Chua. 2017. Neural Factorization Machines for Sparse Predictive Analytics. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval, Shinjuku, Tokyo, Japan, August 7-11, 2017*. 355–364.
- [8] 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.
- [9] Geoffrey Hinton. 2010. A practical guide to training restricted Boltzmann machines. *Momentum* 9, 1 (2010), 926.
- [10] Geoffrey Hinton, Simon Osindero, and Yee-Whye Teh. 2006. A fast learning algorithm for deep belief nets. *Neural computation* 18, 7 (2006), 1527–1554.
- [11] Yifan Hu, Yehuda Koren, and Chris Volinsky. 2008. Collaborative filtering for implicit feedback datasets. In *2008 Eighth IEEE International Conference on Data Mining*. Ieee, 263–272.
- [12] Dietmar Jannach, Markus Zanker, Alexander Felfernig, and Gerhard Friedrich. 2010. *Recommender systems: an introduction*. Cambridge University Press.
- [13] 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.
- [14] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*. 1097–1105.
- [15] Yann LeCun, Bernhard E Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne E Hubbard, and Lawrence D Jackel. 1990. Handwritten digit recognition with a back-propagation network. In *Advances in neural information processing systems*. 396–404.
- [16] Piji Li, Zihao Wang, Zhaochun Ren, Lidong Bing, and Wai Lam. 2017. Neural Rating Regression with Abstractive Tips Generation for Recommendation. (2017).
- [17] Pablo Loyola, Chen Liu, and Yu Hirate. 2017. Modeling User Session and Intent with an Attention-based Encoder-Decoder Architecture. In *Proceedings of the Eleventh ACM Conference on Recommender Systems*. ACM, 147–151.
- [18] Chun-Ta Lu, Lifang He, Hao Ding, Bokai Cao, and S Yu Philip. 2018. Learning from Multi-View Multi-Way Data via Structural Factorization Machines. In *Proceedings of the 2018 World Wide Web Conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 1593–1602.
- [19] Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. 2015. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 43–52.
- [20] Vinod Nair and Geoffrey E Hinton. 2010. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*. 807–814.
- [21] Steffen Rendle. 2012. Factorization Machines with libFM. *ACM Trans. Intell. Syst. Technol.* 3, 3, Article 57 (May 2012), 22 pages.
- [22] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian personalized ranking from implicit feedback. In *Proceedings of the twenty-fifth conference on uncertainty in artificial intelligence*. AUAI Press, 452–461.
- [23] David E Rumelhart, Geoffrey E Hinton, Ronald J Williams, et al. 1988. Learning representations by back-propagating errors. *Cognitive modeling* 5, 3 (1988), 1.
- [24] Ruslan Salakhutdinov and Andriy Mnih. 2011. Probabilistic matrix factorization. In *NIPS*, Vol. 20. 1–8.
- [25] 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.
- [26] J Ben Schafer, Joseph Konstan, and John Riedl. 1999. Recommender systems in e-commerce. In *Proceedings of the 1st ACM conference on Electronic commerce*. ACM, 158–166.
- [27] 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 Eleventh ACM Conference on Recommender Systems*. ACM, 297–305.
- [28] Carles Sierra (Ed.). 2017. *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*. ijcai.org. <http://www.ijcai.org/Proceedings/2017/>
- [29] Ajit P Singh and Geoffrey J Gordon. 2008. Relational learning via collective matrix factorization. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 650–658.
- [30] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research* 15, 1 (2014), 1929–1958.
- [31] Tijmen Tieleman and Geoffrey Hinton. 2012. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning* 4, 2 (2012), 26–31.
- [32] Aaron Van den Oord, Sander Dieleman, and Benjamin Schrauwen. 2013. Deep content-based music recommendation. In *Advances in Neural Information Processing Systems*. 2643–2651.
- [33] Tran Dang Quang Vinh, Tuan-Anh Nguyen Pham, Gao Cong, and Xiao-Li Li. 2018. Attention-based Group Recommendation. *arXiv preprint arXiv:1804.04327* (2018).
- [34] Chong Wang and David M Blei. 2011. Collaborative topic modeling for recommending scientific articles. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 448–456.
- [35] Hao Wang, Naiyan Wang, and Dit-Yan Yeung. 2015. Collaborative deep learning for recommender systems. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 1235–1244.
- [36] Hao Wang, SHI Xingjian, and Dit-Yan Yeung. 2016. Collaborative recurrent autoencoder: Recommend while learning to fill in the blanks. In *Advances in Neural Information Processing Systems*. 415–423.
- [37] Jun Wang, Lantao Yu, Weinan Zhang, Yu Gong, Yinghui Xu, Benyou Wang, Peng Zhang, and Dell Zhang. 2017. Irgan: A minimax game for unifying generative and discriminative information retrieval models. In *Proceedings of the 40th International ACM SIGIR conference on Research and Development in Information Retrieval*. ACM, 515–524.
- [38] Suhang Wang, Yilin Wang, Jiliang Tang, Kai Shu, Suhas Ranganath, and Huan Liu. 2017. What your images reveal: Exploiting visual contents for point-of-interest recommendation. In *Proceedings of the 26th International Conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 391–400.
- [39] Xinxin Wang and Ye Wang. 2014. Improving content-based and hybrid music recommendation using deep learning. In *Proceedings of the ACM International Conference on Multimedia*. ACM, 627–636.
- [40] Xuejian Wang, Lantao Yu, Kan Ren, Guanyu Tao, Weinan Zhang, Yong Yu, and Jun Wang. 2017. Dynamic attention deep model for article recommendation by learning human editors' demonstration. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2051–2059.
- [41] Yao Wu, Christopher DuBois, Alice X Zheng, and Martin Ester. 2016. Collaborative denoising auto-encoders for top-n recommender systems. In *Proceedings of the Ninth ACM International Conference on Web Search and Data Mining*. ACM, 153–162.
- [42] Jun Xiao, Hao Ye, Xiangnan He, Hanwang Zhang, Fei Wu, and Tat-Seng Chua. 2017. Attentional factorization machines: Learning the weight of feature interactions via attention networks. *arXiv preprint arXiv:1708.04617* (2017).
- [43] Fuzheng Zhang, Nicholas Jing Yuan, Defu Lian, Xing Xie, and Wei-Ying Ma. 2016. Collaborative knowledge base embedding for recommender systems. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 353–362.
- [44] Shuai Zhang, Lina Yao, and Aixin Sun. 2017. Deep Learning based Recommender System: A Survey and New Perspectives. *arXiv preprint arXiv:1707.07435* (2017).
- [45] Yongfeng Zhang, Qingyao Ai, Xu Chen, and W Croft. 2017. Joint representation learning for top-n recommendation with heterogeneous information sources. *CIKM*. ACM (2017).
- [46] Lei Zheng, Vahid Noroozi, and Philip S Yu. 2017. Joint Deep Modeling of Users and Items Using Reviews for Recommendation. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*. ACM, 425–434.
- [47] Yin Zheng, Bangsheng Tang, Wenkui Ding, and Hanning Zhou. 2016. A neural autoregressive approach to collaborative filtering. *arXiv preprint arXiv:1605.09477* (2016).

[48] Guorui Zhou, Chengru Song, Xiaoqiang Zhu, Xiao Ma, Yanghui Yan, Xingya Dai, Han Zhu, Junqi Jin, Han Li, and Kun Gai. 2017. Deep Interest Network for

Click-Through Rate Prediction. *arXiv preprint arXiv:1706.06978* (2017).