

A Convolutional Attention Network for Unifying General and Sequential Recommenders

Shahpar Yakhchi^a, Amin Behehsti^a, Seyed-mohssen Ghafari^a, Imran Razzak^b,
Mehmet Orgun^a, Mehdi Elahi^c

^a*Macquarie University, Sydney, Australia*

^b*Deakin University, Geelong, Australia*

^c*Bergen University, Bergen, Norway*

Abstract

General recommenders and sequential recommenders are two modeling paradigms of recommender. The main focus of a general recommender is to identify long-term user preferences, while the user's sequential behaviours are ignored and sequential recommenders try to capture short-term user preferences by exploring item-to-item relations, failing to consider general user preferences. Recently, better performance improvement is reported by combining these two types of recommenders. However, most of the previous works typically treat each item separately and assume that each user-item interaction in a sequence is independent. This may be a too simplistic assumption, since there may be a particular purpose behind buying the successive item in a sequence. In fact, a user makes a decision through two sequential processes, i.e., start shopping with a particular intention and then select a specific item which satisfies her/his preferences under this intention. Moreover, different users usually have different purposes and preferences, and the same user may have various intentions. Thus, different users may click on the same items with an attention on a different purpose. Therefore, a user's behavior pattern is not completely exploited in most of the current methods and they neglect the distinction between users' purposes

*Shahpar Yakhchi is the corresponding author.

Email addresses: `Shahpar.Yakhchi@hdr.mq.edu.au` (Shahpar Yakhchi),
`SAmin.Behehsti@mq.edu.au` (Amin Behehsti), `seyed-mohssen.ghafari@hdr.mq.edu.au` (
Seyed-mohssen Ghafari), `imran.razzak@deakin.edu.au` (Imran Razzak),
`Mehmet.Orgun@mq.edu.au` (Mehmet Orgun), `Mehdi.Elahi@uib.no` (Mehdi Elahi)

and their preferences. To alleviate those problems, we propose a novel method named, CAN, which takes both users' purposes and preferences into account for the next-item recommendation. We propose to use Purpose-Specific Attention Unit (PSAU) in order to discriminately learn the representations of user purpose and preference. The experimental results on real-world datasets demonstrate the advantages of our approach over the state-of-the-art methods.

Keywords: General recommenders, Sequential recommenders, User purpose modeling, personal preference modeling, Attention mechanism, Convolutional neural network.

1. Introduction

Due to the information explosion, people are surrounded by too many options and services. Therefore, there is a need for a tool to help customers with their decision-making process, find their interested items and alleviate the information overload problem. Recommendation systems have emerged as a platform which automatically recommends a small set of items in order to help users find their desired items in online services. Based on how the users' preferences are modelled, there are two types of recommenders: general recommenders and sequential recommenders[1][2][3].

General recommenders aim to learn what items a user is typically interested in. Matrix factorization is one of the most widely used methods in this setting, which learns user-item interactions in a latent vector space to model the general user preferences [4]. While sequential recommenders try to capture sequential patterns from previously visited items. Markov Chains-based classic sequential recommenders assume that the next visited item highly depends on the only most recent visited items [5].

Soon after, convolutional neural networks (CNNs) and recurrent neural networks (RNNs) have become dominant paradigms in modeling complex relations over user-item interaction sequences [6] [7] [8] [9] [10]. Lately, an attention-based approaches such as SHAN [11] can surpass the traditional methods due to the

strong capability of attention mechanism in highlighting the selective parts in a user and item interaction sequence [12].

Both the aforementioned classes of approaches have their strengths and shortcomings [3]. Although general recommenders have been widely adopted
25 to capture long-term user preferences, their performance is limited due to ignoring short-term user preferences. A major advantage of the sequential recommenders is their capability to model sequential dependencies, e.g. a customer who has recently purchased an iPhone is more likely to buy an iWatch next. However, sequential recommenders discard prior user-item interactions within
30 user behaviors, and thus failing to capture general user preferences [1].

Based on the above observation, it is better to build a recommender system which benefits from the advantages of both general- and sequential recommenders. FPMC as an example, is a combination of MC and MF, in which instead of using the same transition matrix for all users, an individual transition matrix is used for each user [2]. FPMC can well capture both sequential
35 behavior and general taste of the users and then linearly combine them [2]. HRM takes one step forward to make progress by using different types of aggregation operations, especially non-linearity into its model [3]. However, users decision-making pattern is not exploited thoroughly by the existing models as
40 they mainly take each user-item interaction independently and consider each item in a sequence as a separated entity. Hence, the current studies may fail to capture local contexts in a session and ignore a user’s purpose which is reflected by a set of clicked successive items in a session. The same user may have various purposes and different users may have different purposes by clicking on the
45 same items. Furthermore, different items within a session may also have different informativeness for revealing purposes and preferences of different users. Therefore, the previous works neglect the hierarchical distinction between user purposes and user preferences, which in turn makes it a challenging task to fully exploit users’ decision-making patterns.

50 Usually, a user’s decision-making process is a combination of two sequential steps; a user’s main purpose and his/her preference. Taking the shopping

event of a user as an example, she/he starts shopping with a specific purpose and then keeps looking into different items until she/he finds items that satisfy her/his preference. Suppose Alice is a PhD student and her previous actions are
55 mostly related to her field of study such as looking for a workshop, and finding an article. Alice has a plan to travel overseas for presenting her work in an international conference. She starts booking her flight and hotel and her next action may be visiting some universities or institutions. While current systems may recommend tourist attractions or car rental companies to her because many
60 users may look for them after booking a hotel and a flight, ignoring her educational purpose of this travel which is hidden inside her long-term interacted item set. Based on this observation, we can see that the user’s main purpose may be hidden inside her/his very previous actions, while analysing her/his very current actions can show her/his preferences on particular items.

65 The above illustrations reveal the difficulty of capturing collective dependency in a session. In the other words, the next choice of item may not be only affected by a part of current session, but all items need to be taken into consideration as a collective of interacted items may have a particular purpose. Moreover, most of these works have taken user-item relationships into consider-
70 ation from the static views and the dynamic property of users’ preferences are ignored. More importantly, the users’ main purposes are not only forgotten, but also there is no difference between the contributions of the same items in modeling preferences of different users. Therefore, how to fully exploit users decision-making process and completely take both the users’ motivations along
75 with their current interests are still largely unexplored.

To address the above issues, we propose a novel model called CAN, A convolutional attention network for unifying general and sequential recommenders, which unifies the benefits of both general- and sequential recommenders. CAN consists of two main modules: purpose encoder and preference encoder. In the
80 purpose encoder we first embed users and items into low-dimensional vectors and then use the CNN network to identify user purposes by capturing the local and high-level information of the long-term interacted item set. Then, we

propose to use a Purpose-Specific Attention Unit (PSAU) to differently attend to different items and fully exploit different informativeness of different items.

85 Next, at preference encoder we also utilize PSAU in order to learn the items' informativeness in the short-term interacted item set to better understand users' preferences. Lastly, the final user representation is learned through coupling user long-term and short-term preferences. The model's parameters are learned by employing the Bayesian personalized ranking optimization criterion to generate a pair-wise loss function [13]. From the experiments, we can observe the
90 superiority of our model over the state-of-the-art algorithms on two datasets. The **key contributions** of the paper are summarized as follows:

- We introduce a unified framework, named CAN, integrating a CNN network and attention-based PSAU module to model the users' purposes and
95 personal preferences.
- We propose a Purpose-Specific Attention Unit, PSAU, which takes user embedding as the query vector of the purpose- and personal preference-level attention networks to differentially attend to important items according to user purposes and preferences.
- 100 • We use the PSAU in both the long- and short-term interacted item set to generate a high-level hybrid user representation.
- We conduct extensive experiments on two real-world datasets. The experimental results demonstrate the superiority of our proposed model compared to the state-of-the-art methods.

105 The rest of the paper is organised as follows: we discuss the related works in Section 2. The proposed methodology and our experiments are presented in Section 3 and Section 4, respectively, before we conclude the paper in Section 5.

2. Related Work

110 Based on different aspects of user behavior, there are two types of paradigms that are applied to recommendation tasks: general recommender and sequential recommender. Both paradigms have strengths and weaknesses, which in the following discussion, we will analyze each paradigms.

2.1. General Recommender

115 The main goal of general recommenders is to discover the users' long-term preferences by exploiting their past items interactions. Early works on this kind of recommenders mostly use Collaborative Filtering (CF) to model users' preferences [14] [15]. Matrix factorization (MF) is one of the widely adopted techniques in CF, which aims to learn user and item latent vectors in order
120 to compute a user's preference on an item [4] [16]. Basically there are two different types of data with which MF-based approaches deal: explicit feedback, e.g., given ratings, and implicit feedback, e.g., mouse clicking. The first one treats making a recommendation as a rating prediction problem, referring to the approaches that try to predict users' preference scores by utilizing their
125 rating patterns [4]. Unlike approaches belonging to the first class, implicit feedback oriented methods formulate making a recommendation as a ranking problem based on the idea of the Learning to-Rank technique [17]. Although general recommenders may better model the long-term user preferences, their performance is limited due to ignoring short-term user preferences.

130 2.2. Sequential Recommender

Different from general recommenders, sequential recommenders try to understand the sequential user behaviors and model the short-term user preferences [18]. Markov chain (MC) has been known as a typical solution in this setting. For instance, SPMC exploits both sequential and social information to
135 make a more personalized recommendation model [19]. In the past few years,

deep learning methods have shown their great capability in modeling the complex interactions between users and items. Among deep neural networks techniques, Recurrent Neural Network (RNN) is one of the widely adopted methods in sequential recommenders due to its capability in sequence modeling. Apart
140 from using basic RNN [6] [20], improved architectures like long short-term-memory (LSTM) [21] and gated recurrent unit (GRU) [22] have also been introduced to better model dependencies in a longer sequence. Different from RNN, Convolutional Neural Network (CNN) stores the embedding of the user-item interaction sequences in a matrix and then treats this matrix as an image [7] [8].
145 Although the basic deep neural networks (i.e., RNN, CNN) have shown a great success in modeling sequential dependencies, they may have some shortcomings in modeling complex relations between users and items. Thus, three advanced models have been introduced to overcome this problem: (i) *attention mechanism*: by more focusing on relevant and important interactions in a sequence [23] [11];
150 (ii) *memory networks*: by incorporating an external memory matrix [24] [25]; and (iii) *mixture models*: by combining the strength of the current deep neural models [26].

Inspired by the outperformance of Transformer [27] [28] [29] [9] in NLP tasks, SRs have motivated to use self-attention technique to better capture sequential
155 dependency. BERT4Rec [30] for instance, has used the deep bidirectional self-attention algorithm to model the sequences of users' behaviors. Except these methods, Graph Neural Network (GNN) has emerged as a solid structure With the strong capability of modelling complex transition patterns of items [31]. SURGE is an example of GNN-based model, in which different types of users' preferences are modelled. The authors have also used graph network to model
160 users' dynamic behaviour.

While sequential recommender models are good at capturing the sequential dependency, they mostly recommend items similar to those that a user currently visited and the general user preference is ignored.

165 2.3. Unified Recommender

There are some recent attempts to combine both general- and sequential recommenders in a unified system. For instance, FPMC is one of the pioneering works in the literature which fuse MF and MC into one model in order to learn the both users' long-and short-term preferences [2]. Soon after, Hierarchical Representation Model (HRM) is proposed by Wang et.al [3] which
170 non-linearly models both sequential behaviors and users' general taste to make a better recommendation. While FPMC and HRM have exploited user long-term preferences to improve the performance of sequential recommenders, Co-Factor benefits from integrating a co-occurrence item-to-item matrix into an MF model [32]. BINN which is proposed by Li et al. [33], is another attempt in
175 unifying both types of users' preferences. The authors have stated that different types of users' actions (e.g., browse, click, collect, cart, and purchase) need to be treated differently. Their proposed model consists of two main components: Neural Item Embedding and Discriminative Behaviors Learning. At first component, BINN tries to find the items' similarities by analysing users' sequential behaviors. While at second component, two alignments Session Behaviors Learning (SBL) and Preference Behaviors Learning (PBL) are introduced to learn discriminative behaviors [33]. Although BINN can record a significant improvement over several state-of-the-art models, it uses LSTM for discriminative
180 behaviors learning part, which may limit the performance of their recommender system as it may not be able to capture the dynamic property of users' preferences. Moreover, BINN only considers purchase behavior for modelling users' historical preferences. This may not only cause in losing some useful information by exploiting other types of users' behaviours (e.g., click, add to cart, and etc),
185 but also may fail to learn latent users' purposes which is hidden in a collection of successive user-item interaction
190

Our model falls under this category and the difference of our method over the existing works can be seen in three different aspects. First, the main purpose of a user's shopping behaviour is ignored in most of the current unified recommenders, which in turn may lead to performance degradation as it plays an
195

important role in the user's decision-making. Second, current methods mostly consider the same informativeness for clicked items in the sequence of user-item interaction, which may result in uncompleted exploited short-term users' preferences. Third, we propose to use a PSAU component to apply in both long-and short term interacted item set in order to dynamically recognize important items for recommendation based on user preferences.

3. Proposed Methodology: Convolutional Attention Network

Before introducing the details of our proposed model, we first define and formulate the research problem and basic concepts and then we present the optimization procedures.

3.1. Notations and Problem Formulation

In this section, we investigate the next-item recommendation problem with implicit feedback data. Let us consider $U = \{u_1, u_2, \dots, u_{|u|}\}$ as the user set and $V = \{v_1, v_2, \dots, v_{|v|}\}$ as the item set, where $|u|$ and $|v|$ are the total number of users and items, respectively. For each user u , we define $G^u = \{S_1^u, S_2^u, \dots, S_T^u\}$

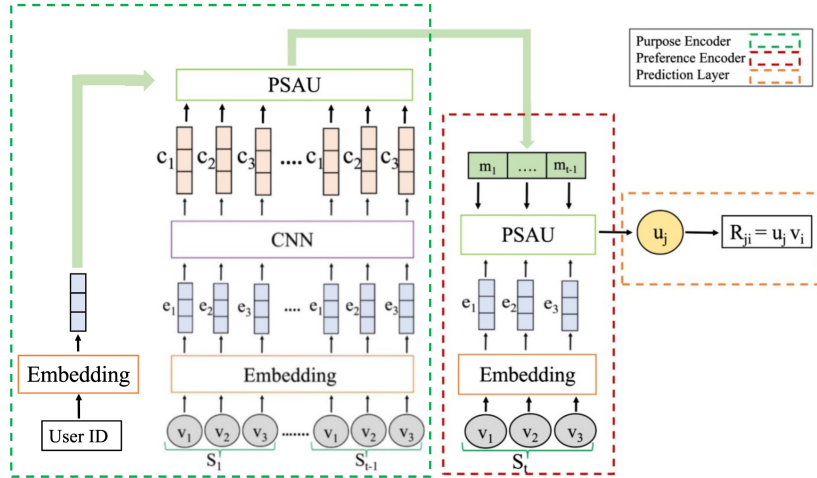


Figure 1: The architecture of CAN, which consists of two main modules purpose encoder and preference encoder.

as her/his transaction history, where T is the total number of sessions and each session $S_t^u \subseteq V (t \in [1, T])$, where S_t^u represents a set of interacted items for users u at time step t . We denote S_t^u as the short-term preference of user u (i.e., her/his sequential behaviour) at specific time step t . In addition to short-term preference, long-term preference of user u is also important for identifying items that users will interact in the near future. Therefore, we consider $G_{t-1}^u = \bigcup_{t=1}^{t-1} S_t^u$ to reflect the long-term preference of user u (i.e., general preference), where G_{t-1}^u is a set of interacted item sets before time step t . For the rest of this paper, we call G_{t-1}^u and S_t^u as the long- and short-term interacted item sets regarding time step t , respectively. Given user u transaction history G^u , we aim to predict the next items which the user will likely purchase by learning her/his long- and short-term preferences.

3.2. Modeling and Learning

The framework of CAN is illustrated in Figure 1. As shown in Figure 1, our proposed model consists of two main modules: (1) *the purpose encoder* and (2) *the preference encoder*. The first module aims to learn the main purpose of the long-term interacted item set for the users. It takes a set of user-item interactions in the long-term item set and embeds them into low-dimensional vector representations, and then these vectors are passed to a CNN network to effectively capture the local contextual information of the sequence in order to identify a user’s main purpose. Then, we propose to use a Purpose-Specific Attention Unit (PSAU) to differentially attend to the users’ main purposes. The reason behind applying PSAU is that different users may have different purposes of buying the same items. For instance, both users a and b buy item i , while user a buys this item as a souvenir for her friend, but user b is interested in this item for herself. Then, we propose to use a Purpose-Specific Attention Unit (PSAU) to differentially attend to the users’ main purposes. The reason behind applying PSAU is that different users may have different purposes of buying the same items. For instance, both users a and b buy item i , while user a buys this item as a souvenir for her friend, but user b is interested in

this item for herself. Thus, we propose to use PSAU in order to incorporate the informativeness of purchasing the same items for different users. The next module is (2) *the preference encoder*, which aims to learn the users’ current preferences. The same user may have different preferences and each item may
245 be more or less informative for that specific preference. Hence, PSAU is also applied here to discriminate each item informativeness.

3.3. Purpose Encoder

Our purpose encoder module has three core components: (i) embedding look-up, (ii) convolutional neural network and (iii) Purpose-Specific Attention Unit (PSAU). Usually users’ decision-making process consists of two sequential
250 vital steps, namely, users’ main purposes and users’ preferences. Normally, people start shopping with an intention and then view different items until they find interesting items that satisfy their preferences. In this block, we aim to first convert a session of items into a sequence of low-dimensional dense vectors.
255 Then, we use a convolutional neural network for capturing local information. Since local contexts within a set of interacted items may imply a user’s purpose. For instance, Julia wants to have a Halloween party. She goes to a shopping and puts a set of *{hanging ghost, pumpkin, lollipop, plastic blood bag}* together. In this collection of items, the local combination of the “hanging ghost”, and
260 “plastic blood bag” may be more important to show the user’s main intention of this shopping. Therefore, we use a CNN network here to learn contextual representations of a set of items. Finally, at this block, PSAU unit is applied to distinguish the level of informativeness of different items in revealing the users’ motivations of purchasing a set of items together. The reason behind
265 using PSAU unit in purpose encoder is that different items may have different level of contributions in presenting a user’s main purpose, and the same words may have different informativeness for the recommendation of different users. Based on this observation, we need to identify important items in demonstrating shopping’s purpose of different users, and thus the personalized attention-based
270 network is proposed to apply in this block.

Embedding Look-up. First, we use embedding look-up to embed user and item IDs (i.e., one-hot representations) into two continuous low-dimensional spaces, where e_i represents the item embedding vector of item i , and u_j denotes the user embedding vector of user j . The embedding matrix is denoted by $E = [e_1, e_2, \dots, e_i]$, $E \in R^{|V| \times D}$, where D and $|V|$ represent the embedding dimension and the total number of items, respectively. The matrix $U \in R^{D \times |U|}$ is the user embedding matrix, where u_j denotes the user embedding vector of user j .

Convolutional Neural Network (CNN). Second, we employ CNN to learn contextual information of user-item interactions [34]. CNN is one of the deep learning techniques with a great capability in capturing local information [35]. Therefore, we use CNN to capture the user’s main purpose in the long-term item set. Next, we perform a convolution operator on the matrix E as the concatenation of the items’ embedding vectors. Let $K_w \in R^{N_f \times (2K+1)D}$, and $b_w \in R^{N_f}$ denote the parameters of CNN network, in which K_w is the kernel and b_w represents the bias parameters. N_f is the number of CNN filters, and $2K+1$ is the window size of CNN. Then, c_i illustrates the contextual representation of item i :

$$c_i = \text{ReLU}(K_w \times e_{[i-k]:[i+k]} + b_w) \quad (1)$$

, where $e_{[i-k]:[i+k]} \in G_{t-1}^u$ is the combination of the embedding vectors of items from position $[i-k]$ to position $[i+k]$. We use ReLU as our non-linear activation function.

Purpose-Specific Attention Unit (PSAU). The last component in the *purpose encoder* is the Purpose-Specific Attention Unit (PSAU), to differentially attend to important items according to user purposes. In a sequence of user-item interactions, each item may be more or less informative for learning users’ purpose representation. For instance, imagine $\{\text{pizza bread}, \text{pepperoni}, \text{cheese}\}$ as a set of purchased items together for making a pizza. In this shopping basket, *pizza bread* is more informative to represent the users’ purposes than *cheese*.

Furthermore, different users may purchase the same items for a different purpose. Therefore, based on these observations, identifying the contributions of different items for different users play an important role in personalized recommendation. However, most of the current approaches use a classic attention network which computes attention score as a weighted sum over the embeddings of items and a fixed attention query vector, ignoring users' main purposes. To learn the informativeness of each item for different users, we propose to employ the PSAU cell to identify the most informative items related to the users' main purpose within a user-item interaction sequence. PSAU first takes the embedded user-ID vector $u'_j \in R^{D_u}$, where D_u is the user embedding dimension. Then, we use a dense non-linear layer to transform the embedding vector u'_j to the purpose-level user preference vector p_j , which is formulated as:

$$p_j = ReLU(W_1 \times u'_j + b_1), \quad (2)$$

where $W_1 \in R^{D_u \times D_p}$ and $b_1 \in R^{D_p \times 1}$ are model parameters, and D_p is the preference vector dimension. Next, we denote α_j as the attention score of item j , which can extract the level of informativeness of each item according to the users' main purpose. The attention score α_j , is calculated based on the interaction between the user preference vector and the contextual item representations, which is shown as :

$$a_i = c_i^T \tanh(W_2 \times p_j + b_2), \quad (3)$$

$$\alpha_i = \frac{\exp(a_i)}{\sum_{i \in G_{t-1}^u} \exp(a_i)} \quad (4)$$

, where $W_2 \in R^{D_p \times N_f}$ and $b_2 \in R^{N_f \times 1}$ are model parameters. Next, the user's main purpose representation m_i is modeled as a weighted sum of the contextual representation of item i with their attention scores. Formally, this representation can be formulated as follows:

$$m_i = \sum_{i \in G_{t-1}^u} \alpha_i c_i \quad (5)$$

3.4. Preference Encoder

As it is clear from Figure 1, PSAU is also employed in the preference encoder module in order to learn an informative user short-term preference representation. Different users may have different preferences by clicking on the same items and different items are more or less informative for modeling user preferences. Hence, we use PSAU here as well to model the different informativeness of the same items for different users. Hence, we first take the item embedding $e_i \in S_t^u$ in a short-term interacted item set to model a user preference vector p_d , which is shown as:

$$p_d = ReLU(W_3 \times e_i + b_3), \quad (6)$$

where $W_3 \in R^{D_u \times D_q}$ and $b_3 \in R^{D_q \times 1}$, and D_q is the preference query size. Next, the attention weight α'_i represents the level of informativeness of item i in the short-term user preference, which can be computed by the interactions between the user's purpose representation and user preference vector. Then, the softmax function is used to normalize the attention weight, which is calculated as follows:

$$a'_i = m_i^T \tanh(W_4 \times p_d + b_4), \quad (7)$$

$$\alpha'_i = \frac{\exp(a_i)}{\sum_{i \in S_t^u} \exp(a_i)} \quad (8)$$

where $W_4 \in R^{D_q \times N_f}$ and $b_4 \in R^{N_f \times 1}$ are model parameters. Finally, the contextual user representation u_j is computed as follows:

$$u_j = \sum_{i \in S_t^u} \alpha'_i m_i \quad (9)$$

3.5. Prediction Layer

After the final user representation u_j has been learned, we calculate the inner product of it and item representation v_i in order to compute the user preference score R_{ij} as follows:

$$R_{ij} = u_j v_i \quad (10)$$

Next, followed by [13], we utilize a pair-wise loss function in order to train our
 345 model. We aim to provide a ranked list of the next items to be recommended,
 where observed items should have higher score than unobserved ones. Let $D =$
 $\{(u, v_i, v_j) : u \in U, v_i \in G^u, v_j \in V/G^u\}$ denote the set of pair-wise training
 instances. Then we train our model by maximizing a posterior (MAP) as follows:

$$\arg \min_{\Theta} \sum_{(u, v_i, v_j) \in D} -\ln \sigma(R_i^u - R_j^u) + \lambda_{uv} \|\theta_{uv}\|^2 + \lambda_a \|\theta_a\|^2 \quad (11)$$

where $\theta_{uv} = \{U, V\}$ is the set of user and item embedding parameters, $\theta_a =$
 350 $\{W_1, W_2, W_3, W_4\}$ is the set of weights of attention networks, λ_{uv} and λ_a are
 the regularization parameters, and σ is a logistic function.

4. Experiments

In this section, we present experimental evaluation of proposed recomender
 and compare the performance with state-of-the-art baseline methods such as
 355 BPR [13], FOSSIL [36], Caser [7], FPMC [2], HRM [3], GRU4Rec [6], NARM [37],
 SHAN [11], and MEANS [25].

4.1. Datasets and Experimental Setting

We conduct our experiments on two widely used datasets Tmall ¹ and
 Gowalla ². The Tmall dataset records the user's consumption and browsing
 360 behavior during the user's shopping process. It has too many interactions of
 424,170 users on 1,090,390 items within six months. In this dataset there are
 four kinds of activities: click, collect, add-to-cart and purchase. Following the
 settings in [11] and [38] we only consider the users' purchase activities in our
 experiment. The Gowalla aggregates the users' check-in information from the

¹<https://tianchi.aliyun.com/dataset/dataDetail?dataId=53>

²<https://snap.stanford.edu/data/loc-gowalla.html>

Table 1: Statistics of our datasets

Dataset	Users	Items	Sessions length	Training sessions	Testing sessions	Interactions
Tmall	20,716	25,143	2.81	71,998	3565	85,432
Gowalla	15,254	13,052	2.99	128,115	3611	94,654

location-based social networking website, Gowalla from February 2009 to October 2010. Gowalla consists of 6,442,890 number of total check-ins, where each record consists of user id, timestamp, GPS location and POI id. We follow the same preprocessing procedure as in SHAN [11] and we treat user transactions or check-ins in one day as a session. Sessions with only one item and items with less than 20 time observations are removed from datasets. We randomly select the sessions in the last week as a test set, and the rest are used for training. In addition, we randomly keep one item in each session as the next item to be predicted. The statistics of the datasets after the preprocessing stage are illustrated in Table 1.

Baselines: To demonstrate the effectiveness of our method, we compare it with the following representative state-of-the-art recommender systems built on various frameworks including RNN, CNN, attention models and memory networks:

- TOP: This method identifies the top popular items based on the number of occurrences in each session in the training data, and then recommends those items in test data.
- BPR [13]: This is a state-of-the-art baseline for binary implicit feedback through pairwise learning to rank.
- FOSSIL [36]: This method integrates factored item similarity with a Markov chain to model the user’s long- and short-term preferences.
- Caser [7]: This is a state-of-the-art model, which uses CNN for sequence embedding.

- FPMC [2]: This is a combination of MF and MC model in order to learn user preferences.
- 390 • HRM [3]: This model non-linearly learns both sequential behavior and users' general taste to make a better recommendation.
- GRU4Rec [6]: This is a state-of-the-art sequential recommender, which applies modern recurrent neural network (GRU) to be able to model the whole session.
- 395 • NARM [37]: This is a sequential recommender which combines a recurrent neural network with an attention network.
- SHAN [11]: This is a state-of-the-art sequential recommender, which employs a two-layer hierarchical attention network to learn long- and short-term preference.
- 400 • MEANS [25]: This model first operates a max-pooling technique on the most recent sessions and the results are stored into an external memory. Then the attention mechanism is applied to learn long-term user preference. Finally, at prediction layer a recommendation is made by learning a mixture of long- and short-term preference.

405 **Evaluation Metrics.** Similar to the previous work [11], we also adopt several widely used evaluation metrics AUC, Recall@N, and Precision@N to evaluate the performance of our model, where $N \in \{5, 10, 20\}$. Recall measures the proportion of the right ranked items overall top-k recommendation items in a list, while Precision measures the proportion of results which are relevant. Different from both above metrics, AUC computes how highly predicted items are ranked over all items. The larger metric scores show better model performance. Due to the space limitation, we name Recall and Precision as Re and Pre in the rest of the paper, respectively.

Parameter Settings. We set the item embedding and user embedding 415 dimensions, D , to 100, which is a trade-off between the performance of recommendation and the computation cost for both datasets. Similar to the [39],

Table 2: Impact of different regularization at Recall@20

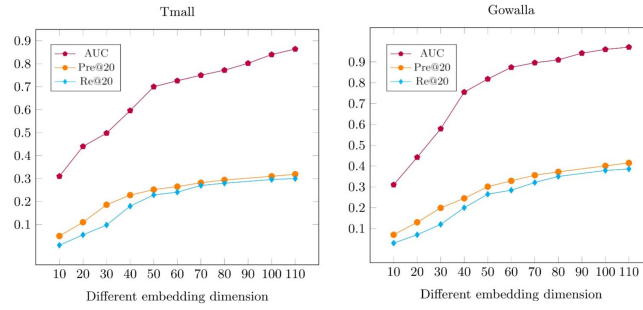
Dataset	λ_{uv} λ_{α}	0	1	10	50
Tmall	0.01	0.085	0.126	0.143	0.146
	0.001	0.079	0.124	0.138	0.139
	0.0001	0.073	0.111	0.129	0.133
Gowalla	0.01	0.250	0.344	0.355	0.372
	0.001	0.321	0.397	0.423	0.432
	0.0001	0.342	0.421	0.452	0.461

we set the number of CNN kernels N_f and the window size to 400 and 3, respectively. We apply dropout strategy [40] to each layer of CNN in order to avoid overfitting. The dropout rate is set to 0.2, the batch size is empirically set to 50, the sizes of both the user purpose query D_p and preference query D_q are set to 200. The learning rate η is 0.01. Items and users dimensions are randomly initialized with normal distribution $N(0, 0.01)$ and then learned during the training process. The attention parameters are initialized with the $U(-\sqrt{\frac{3}{k}}, \sqrt{\frac{3}{k}})$.

4.2. Impact of Hyper-parameters

In this subsection, we investigate the impact of hyper-parameters on the performance of CAN. We consider $\lambda_{uv} = \{0.01, 0.001, 0.0001\}$ as our user and

Figure 2: Impact of different embedding dimension on Gowalla and Tmall datasets. In each figure, we have shown the impact of different embedding sizes on three evaluation metrics AUC, Precision and Recall.



item embedding regularization, and $\lambda_a = \{0, 1, 10, 50\}$ as our attention network regularization. Based on the Table 2, the performance of CAN is gradually
430 increased when $\lambda_a > 0$ in both Tmall and Gowalla datasets, which indicates the effectiveness of applying attention mechanism in our model. We also test the impact of different embedding dimensions, D , related to the user, item and hidden layer parameters in attention network. As it is clear from Figure 2, the higher embedding dimension can result in better AUC, Recall@20, and
435 Precision@20 as it can learn more latent features form user and item as well as their interactions through attention mechanism. From this figure, a slight improvement is recorded while the embedding dimension is increased from 100, and thus we set the embedding size to 100.

4.3. Impact of Different Sessions Lengths

440 We examine the performance of CAN under different sequence lengths as the local features captured by CNN network may be different. Table 3 demonstrates the results of our investigation. We consider sessions with less than 3 items as a short session and treat sessions with more than 3 items as a long session. The percentage of short and long sessions are 90%, 10% and 83%, 17% in both Tmall
445 and Gowalla datasets, respectively. In Table 3, CAN-S refers to a situation

Table 3: Impact of different session lengths
Tmall Gowalla

Methods	AUC	Re@20	Pre@20	Methods	AUC	Re@20	Pre@20
CAN-S	0.745	0.196	0.213	CAN-S	0.814	0.219	0.263
CAN-L	0.889	0.221	0.282	CAN-L	0.916	0.298	0.342

Table 4: Impcat of CAN modules
Tmall width=0.95 Gowalla width=0.95

Methods	AUC	Re@20	Pre@20	Methods	AUC	Re@20	Pre@20
CAN-PurEn	0.817	0.256	0.278	CAN-PurEn	0.924	0.284	0.312
CAN-PreEn	0.781	0.210	0.264	CAN-PreEn	0.899	0.256	0.299
CAN	0.915	0.317	0.322	CAN	0.989	0.392	0.401

where short sessions are modelled, while only long sessions are considered in CAN-L. From this table, we can have several observation. First, the performance of both CAN-L and CAN-S are too close. Second, CAN-L performs slightly better than CAN-S with respect to AUC, Pre@20, and Re@20 in both Tmall and Gowalla datasets. This is probably because of capturing the more contextual features through long sessions. Third, the performance of CAN-L is still too close to the overall performance of our model.

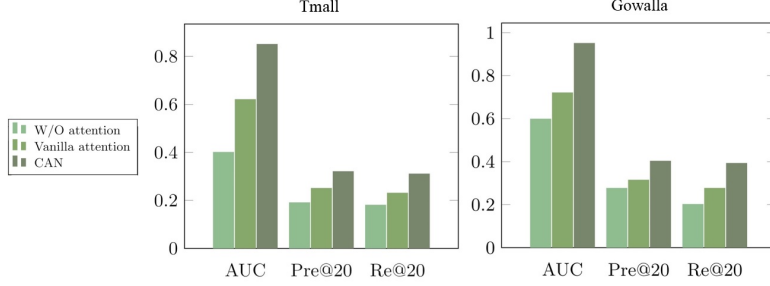
4.4. Impact of CAN Modules

In this experiment, we aim to test the performance of two modules, i.e., purpose encoder and preference encoder in Table 4. CAN-PurEn means only user purpose module is used, while CAN-PreEn only considers a user’s preference. According to the Table 4, we can have several observations. First, the CAN-PurEn can effectively improve the performance of our approach, as it can help our model CAN to achieve the higher performance compared to the state-of-the-art models. This may be due to the capturing the local patterns in a long-term interacted item set through CNN and highlighting the important items according to user preferences by PSAU cell. Second, the CAN-PreEn is also another effective module in our model, which indicates a significant improvement in the performance of CAN. This is probably because items in a short-term interacted item set usually have different informativeness and recognizing the important items can help better modeling user representations. Third, generally CAN performs better than two single modules. It demonstrates that combining these two modules is helpful in learning user representation and predicting next items.

4.5. Impact of PSAU component

In order to verify the effectiveness of the PSAU component in our model, we compare the performance of our model in the presence and absence of the PSAU cell. As it is clear from Figure 3, we have different findings: (1) applying attention mechanism can show better performance compared to the model without attention. The reason behind this observation may be because of assigning

Figure 3: Impact of PSAU on Gowalla and Tmall datasets. W/O attention means no attention mechanism is used.



different weight to different items, and attention mechanism can discover the important items in a user-item interaction; (2) our model CAN consistently outperforms the model without attention mechanism and vanilla attention. The reason behind this observation may be because of assigning different score to the same items for modeling different users, while vanilla attention assigns a fixed score and thus is not able to differentiate the importance of the same items in modeling the different user preferences. Attention mechanism pays same attention to each item by computing the attention weights only based on the input representation sequence via a fixed vector, and thus the user preferences are not incorporated. While in contrast to vanilla attention, the attention scores in PSAU are computed based on the interaction between the user preference vector and the contextual item representations. Therefore, our model can highlight important items in user's purpose according to her/his personal preference, which in turn can help in better user representation learning. Based on these results, we can validate the effectiveness of the PSAU cell in our approach.

4.6. Overall Performance Comparison

In this subsection, we compare the results of our model with the other state-of-the-art approaches in both Tmall and Gowalla datasets, which is summarized in Tables 5 and 6. This table illustrates that:

1. According to Tables 5 and 6, where the best result in each row is highlighted in boldface, our proposed model significantly and consistently out-

Table 5: The performance of different methods regarding the evaluation metrics in Tmall dataset.

Datasets	Tmall						
Metrics	Re5	Re10	Re20	Pre5	Pre10	Pre20	AUC
Top	0.021	0.052	0.084	0.051	0.062	0.074	0.392
BPR	0.024	0.090	0.122	0.062	0.069	0.074	0.481
Fossil	0.110	0.120	0.125	0.083	0.088	0.092	0.691
Caser	0.041	0.049	0.052	0.100	0.108	0.115	0.701
FPMC	0.050	0.055	0.061	0.118	0.125	0.130	0.742
HRM	0.060	0.065	0.070	0.121	0.129	0.133	0.751
GRU4Rec	0.062	0.065	0.069	0.138	0.145	0.149	0.762
NARM	0.063	0.068	0.073	0.141	0.149	0.159	0.781
SHAN	0.071	0.076	0.079	0.155	0.160	0.166	0.789
MEANS	0.074	0.079	0.082	0.163	0.172	0.177	0.790
CAN	0.201	0.278	0.317	0.200	0.260	0.322	0.915

performs all state-of-the-art models in terms of Precision@N, Recall@N and AUC in different N s in both Tmall and Gowalla datasets. Specifically, compared to MEANS which is the best baseline in terms of all evaluation metrics, CAN has shown 14% and 16% improvements with respect to the AUC on Tmall and Gowalla datasets, respectively. This indicates the effectiveness of CAN, which can recognize important items in users' purposes according to their preferences through CNN network and PSAU component.

2. Deep learning methods using attention network (CAN, MEANS, SHAN, and NARM) show better performance compared with the methods without attention mechanism. The reason may be due to the capability of attention mechanism in recognizing the most important items in user and item interaction.
3. Overall, all unified approaches (CAN, MEANS, SHAN, NARM, HRM,

Table 6: The performance of different methods regarding the evaluation metrics in Gowalla dataset.

Datasets	Gowalla						
Metrics	Re5	Re10	Re20	Pre5	Pre10	Pre20	AUC
Top	0.038	0.048	0.059	0.061	0.066	0.071	0.711
BPR	0.069	0.074	0.081	0.077	0.082	0.089	0.800
Fossil	0.215	0.298	0.312	0.091	0.095	0.099	0.810
Caser	0.075	0.083	0.089	0.114	0.119	0.124	0.815
FPMC	0.115	0.129	0.138	0.127	0.133	0.142	0.820
HRM	0.119	0.125	0.145	0.150	0.157	0.161	0.824
GRU4Rec	0.121	0.135	0.141	0.155	0.160	0.165	0.828
NARM	0.130	0.136	0.140	0.156	0.159	0.163	0.830
SHAN	0.135	0.140	0.144	0.163	0.169	0.175	0.832
MEANS	0.142	0.150	0.158	0.170	0.175	0.180	0.840
CAN	0.250	0.312	0.392	0.360	0.399	0.401	0.989

510 FPMC, and Fossil) outperform the best general- and sequential recom-
menders such as BPR and GRU4Rec, respectively.

4. Among all unified approaches, after CAN, MEANS outperforms others
like SHAN, NARM, HRM, FPMC, and Fossil. While the performance of
MEANS and SHAN are too close, MEANS can achieve around 5% and
515 9% improvement compared to SHAN at Recall@20 in Tmall and Gowalla
datasets, respectively. This indicates the effect of using external memory
to store long-term user and item interaction after a max-pooling oper-
ation. However, MEANS cannot effectively model the local contexts in
the long term user preference, and is not able to find important items for
revealing purposes and preferences of different users. Moreover, although
520 MEANS uses attention mechanism, it can not model the informativeness
of different items. Different from all mentioned approaches, our proposed
model can dynamically find important items according to user purposes

and preferences.

525 5. Conclusion

In this paper, we propose a novel unified recommendation approach which consists of a Purpose-Specific Attention Unit (PSAU). In our approach, CAN, we learn the users' purposes in long-term interacted item set by using CNN. We use PSAU cell to recognize important items in users' purposes according to their
530 preferences. Since same items may have different informativeness for different users, we use PSAU in short-term interacted item set as well to model users' preferences. The extensive experimental results on the real-world datasets validate the effectiveness of our approach compared to other state-of-the-art methods. As our future work, we aim to take contextual information into sequential rec-
535 ommenders in order to make a more accurate recommendation. Furthermore, modelling different heterogeneous actions can be another direction for our future work.

References

- [1] D. Dong, X. Zheng, R. Zhang, Y. Wang, Recurrent collaborative filtering
540 for unifying general and sequential recommender, in: Proceedings of the 27th Int. Joint Conf. on Artificial Intelligence, IJCAI, ijcai.org, 2018, pp. 3350–3356.
- [2] S. Rendle, C. Freudenthaler, L. Schmidt-Thieme, Factorizing personalized markov chains for next-basket recommendation, in: Proc. of the 19th Int.
545 Conf. on World Wide Web, WWW 2010, ACM, 2010, pp. 811–820.
- [3] P. Wang, et al., Learning hierarchical representation model for nextbasket recommendation, in: Proc. of the 38th Int. SIGIR Conf. on Research and Development in Information RetrievalAug, 2015, pp. 403–412.
- [4] Y. Koren, R. M. Bell, C. Volinsky, Matrix factorization techniques for
550 recommender systems, IEEE Journal of Computer 42 (8) (2009) 30–37.

- [5] M. Grbovic, V. Radosavljevic, N. Djuric, N. Bhamidipati, J. Savla, V. Bhagwan, D. Sharp, E-commerce in your inbox: Product recommendations at scale, in: Proc. of the 21th SIGKDD Int. Conf. on Knowledge Discovery and Data Mining, Australia, ACM, 2015, pp. 1809–1818.
- 555 [6] B. Hidasi, A. Karatzoglou, L. Baltrunas, D. Tikk, Session-based recommendations with recurrent neural networks, in: 4th Int. Conf. on Learning Representations, ICLR, 2016.
- [7] J. Tang, K. Wang, Personalized top-n sequential recommendation via convolutional sequence embedding, in: Proc. of the 11th Int. Conf. on Web Search and Data Mining, WSDM, USA, ACM, 2018, pp. 565–573.
- 560 [8] F. Yuan, A. Karatzoglou, I. Arapakis, J. M. Jose, X. He, A simple convolutional generative network for next item recommendation, in: Proc. of the 12th Int. Conf. on Web Search and Data Mining, WSDM, Australia, ACM, 2019, pp. 582–590.
- 565 [9] H. Zogan, I. Razzak, S. Jameel, G. Xu, Depressionnet: Learning multi-modalities with user post summarization for depression detection on social media, in: Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, 2021, pp. 133–142.
- [10] H. Zogan, I. Razzak, X. Wang, S. Jameel, G. Xu, Explainable depression detection with multi-modalities using a hybrid deep learning model on social media, arXiv preprint arXiv:2007.02847 (2020).
- 570 [11] H. Ying, F. Zhuang, F. Zhang, Y. Liu, G. Xu, X. Xie, H. Xiong, J. Wu, Sequential recommender system based on hierarchical attention networks, in: Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI, Sweden, ijcai.org, 2018, pp. 3926–3932.
- 575 [12] D. Bahdanau, K. Cho, Y. Bengio, Neural machine translation by jointly learning to align and translate, in: Third Int. Conf. on Learning Representations, ICLR, USA, 2015.

- [13] S. Rendle, C. Freudenthaler, Z. Gantner, L. Schmidt-Thieme, BPR: bayesian personalized ranking from implicit feedback, in: Proc. of the 25th Conf. on Uncertainty in Artificial Intelligence (UAI), Canada, AUAI Press, 2009, pp. 452–461.
- [14] Y. Koren, R. M. Bell, in: Recommender Systems Handbook, Springer, 2011, pp. 145–186.
- [15] B. M. Sarwar, G. Karypis, J. A. Konstan, J. Riedl, Item-based collaborative filtering recommendation algorithms, in: Proc. of the 10th Int. World Wide Web Conf. WWW, China, ACM, 2001, pp. 285–295.
- [16] R. Salakhutdinov, A. Mnih, Probabilistic matrix factorization, in: Proc. of the 21th Annual Conf. on Neural Information Processing Systems, Canada, Curran Associates, Inc., 2007, pp. 1257–1264.
- [17] A. Karatzoglou, L. Baltrunas, Y. Shi, Learning to rank for recommender systems, in: Seventh ACM Conf. on Recommender Systems, RecSys '13, China, ACM, 2013, pp. 493–494.
- [18] S. Wang, L. Hu, Y. Wang, L. Cao, Q. Z. Sheng, M. A. Orgun, Sequential recommender systems: Challenges, progress and prospects, in: Proc. of the 28th Int. Joint Conf. on Artificial Intelligence, IJCAI, China, ijcai.org, 2019, pp. 6332–6338.
- [19] C. Cai, R. He, J. J. McAuley, SPMC: socially-aware personalized markov chains for sparse sequential recommendation, CoRR (2017).
- [20] Y. Zhang, H. Dai, C. Xu, J. Feng, T. Wang, J. Bian, B. Wang, T. Liu, Sequential click prediction for sponsored search with recurrent neural networks, CoRR (2014).
- [21] C. Wu, A. Ahmed, A. Beutel, A. J. Smola, H. Jing, Recurrent recommender networks, in: Proc. of the 10th Int. Conf. on Web Search and Data Mining, WSDM, United Kingdom, ACM, 2017, pp. 495–503.

- [22] B. Hidasi, M. Quadrana, A. Karatzoglou, D. Tikk, Parallel recurrent neural network architectures for feature-rich session-based recommendations, in: Proc. of the 10th ACM Conf. on Recommender Systems, USA, ACM, 2016, pp. 241–248.
- 610 [23] W. Kang, M. Wan, J. J. McAuley, Recommendation through mixtures of heterogeneous item relationships, CoRR abs/1808.10031 (2018).
- [24] X. Chen, H. Xu, Y. Zhang, J. Tang, Y. Cao, Z. Qin, H. Zha, Sequential recommendation with user memory networks, in: Proc. of the 11th Int. Conf. on Web Search and Data Mining, WSDM, ACM, 2018, pp. 108–116.
- 615 [25] C. Hu, P. He, C. Sha, J. Niu, Memory-augmented attention network for sequential recommendation, in: R. Cheng, N. Mamoulis, Y. Sun, X. Huang (Eds.), Int. Conf. on the Web Information Systems Engineering - WISE, China, Vol. 11881, Springer, 2019, pp. 228–242.
- [26] J. Tang, F. Belletti, S. Jain, M. Chen, A. Beutel, C. Xu, E. H. Chi, Towards
620 neural mixture recommender for long range dependent user sequences, in: The World Wide Web Conf., WWW, USA, ACM, 2019, pp. 1782–1793.
- [27] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. L. Scao, S. Gugger, M. Drame,
625 Q. Lhoest, A. M. Rush, Transformers: State-of-the-art natural language processing, in: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, EMNLP 2020 - Demos, Online, November 16-20, 2020, 2020, pp. 38–45.
- [28] U. Naseem, I. Razzak, K. Musial, M. Imran, Transformer based deep intelligent contextual embedding for twitter sentiment analysis, Future Generation Computer Systems 113 (2020) 58–69.
630
- [29] U. Naseem, I. Razzak, M. Khushi, P. W. Eklund, J. Kim, Covidsent:

A large-scale benchmark twitter data set for covid-19 sentiment analysis,
IEEE Transactions on Computational Social Systems (2021).

- 635 [30] F. Sun, J. Liu, J. Wu, C. Pei, X. Lin, W. Ou, P. Jiang, Bert4rec: Sequential recommendation with bidirectional encoder representations from transformer, in: Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM 2019, Beijing, China, November 3-7, 2019, 2019, pp. 1441–1450.
- 640 [31] S. Wu, W. Zhang, F. Sun, B. Cui, Graph neural networks in recommender systems: A survey, CoRR abs/2011.02260 (2020).
- [32] D. Liang, J. Alotaibi, L. Charlin, D. M. Blei, Factorization meets the item embedding: Regularizing matrix factorization with item co-occurrence, in: Proc. of the 10th ACM Conf. on Recommender Systems, USA, ACM, 2016, pp. 59–66.
- 645 [33] Z. Li, H. Zhao, Q. Liu, Z. Huang, T. Mei, E. Chen, Learning from history and present: Next-item recommendation via discriminatively exploiting user behaviors, CoRR abs/1808.01075 (2018).
- [34] Y. Kim, Convolutional neural networks for sentence classification, in: A. Moschitti, B. Pang, W. Daelemans (Eds.), Proc. of the 2014 Conf. on Empirical Methods in Natural Language Processing, EMNLP, Qatar, ACL, 2014, pp. 1746–1751.
- 650 [35] C. Wu, F. Wu, J. Liu, S. He, Y. Huang, X. Xie, Neural demographic prediction using search query, in: 12th ACM Int. WSDM Conf., Australia, ACM, 2019, pp. 654–662.
- 655 [36] R. He, J. J. McAuley, Fusing similarity models with markov chains for sparse sequential recommendation, CoRR abs/1609.09152 (2016).
- [37] J. Li, P. Ren, Z. Chen, Z. Ren, T. Lian, J. Ma, Neural attentive session-based recommendation, in: Proc. of the Conf. on Information and Knowledge Management, CIKM, Singapore, ACM, 2017, pp. 1419–1428.
- 660

- [38] L. Hu, L. Cao, S. Wang, G. Xu, J. Cao, Z. Gu, Diversifying personalized recommendation with user-session context, in: Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017, ijcai.org, 2017, pp. 1858–1864.
- 665 [39] C. Wu, F. Wu, M. An, J. Huang, Y. Huang, X. Xie, NPA: neural news recommendation with personalized attention, in: Proc. of the 25th SIGKDD Int. Conf. on Knowledge Discovery and Data Mining, KDD, USA, ACM, 2019, pp. 2576–2584.
- [40] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov,
670 Dropout: a simple way to prevent neural networks from overfitting, The Journal of Machine Learning Research 15 (1) (2014) 1929–1958.