

# Deep Reinforcement Learning for Page-wise Recommendations

Xiangyu Zhao  
Data Science and Engineering Lab  
Michigan State University  
zhaoxi35@msu.edu

Long Xia  
JD.com  
xialong@jd.com

Liang Zhang  
JD.com  
zhangliang16@jd.com

Zhuoye Ding  
JD.com  
dingzhuoye@jd.com

Dawei Yin  
JD.com  
yindawei@acm.org

Jiliang Tang  
Data Science and Engineering Lab  
Michigan State University  
tangjili@msu.edu

## ABSTRACT

Recommender systems can mitigate the information overload problem by suggesting users' personalized items. In real-world recommendations such as e-commerce, a typical interaction between the system and its users is – users are recommended a page of items and provide feedback; and then the system recommends a new page of items. To effectively capture such interaction for recommendations, we need to solve two key problems – (1) how to update recommending strategy according to user's *real-time feedback*, and (2) how to generate a page of items with proper display, which pose tremendous challenges to traditional recommender systems. In this paper, we study the problem of page-wise recommendations aiming to address aforementioned two challenges simultaneously. In particular, we propose a principled approach to jointly generate a set of complementary items and the corresponding strategy to display them in a 2-D page; and propose a novel page-wise recommendation framework based on deep reinforcement learning, DeepPage, which can optimize a page of items with proper display based on real-time feedback from users. The experimental results based on a real-world e-commerce dataset demonstrate the effectiveness of the proposed framework.

## KEYWORDS

Recommender Systems; Deep Reinforcement Learning; Actor-Critic; Item Display Strategy; Sequential Preference

### ACM Reference Format:

Xiangyu Zhao, Long Xia, Liang Zhang, Zhuoye Ding, Dawei Yin, and Jiliang Tang. 2018. Deep Reinforcement Learning for Page-wise Recommendations. In *Twelfth ACM Conference on Recommender Systems (RecSys '18)*, October 2–7, 2018, Vancouver, BC, Canada. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3240323.3240374>

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

RecSys '18, October 2–7, 2018, Vancouver, BC, Canada

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-5901-6/18/10...\$15.00

<https://doi.org/10.1145/3240323.3240374>



Figure 1: An example to illustrate the interactions between recommender systems and users.

## 1 INTRODUCTION

Recommender systems are intelligent E-commerce applications [4, 21, 26]. They assist users in their information-seeking tasks by suggesting items (products, services, or information) that best fit their needs and preferences. Recommender systems have become increasingly popular in recent years, and have been utilized in a variety of domains including movies, music, locations, and social events [2, 11, 13, 30, 31, 39, 41, 45]. Figure 1 illustrates a typical example of the interactions between an e-commerce recommender system and a user – each time the system recommends a page of items to the user; next the user browses these items and provides real-time feedback and then the system recommends a new page of items. This example suggests two key challenges to effectively take advantage of these interactions for e-commerce recommender systems – 1) how to efficiently capture user's preference and update recommending strategy according to user's real-time feedback; and 2) how to generate a page of items with proper display based on user's preferences.

### 1.1 Real-time Feedback

Most existing recommender systems consider the recommendation procedure as a static process and make recommendations following a fixed greedy strategy. However, these approaches may fail in capturing the dynamic nature of the users' preferences, and they become infeasible to efficiently and continuously update their recommending strategies according to user's real-time feedback. Thus, in this work, we consider the recommendation procedure as sequential interactions between users and the recommender agent; and leverage Reinforcement Learning (RL) to automatically learn the optimal recommendation strategies. Recommender systems based

on reinforcement learning have two major advantages. First, they are able to continuously update their strategies based on user's real-time feedback during the interactions, until the system converges to the optimal strategy that generates recommendations best fitting users' dynamic preferences. Second, the optimal strategy is made by maximizing the expected long-term cumulative reward from users; while the majority of traditional recommender systems are designed to maximize the immediate (short-term) reward of recommendations [32]. Therefore, the system can identify items with small immediate rewards but making big contributions to the rewards for future recommendations.

## 1.2 Page-wise Recommendations

As mentioned in the example, users are typically recommended a page of items. To achieve this goal, we introduce a page-wise recommender system, which is able to jointly (1) generate a set of diverse and complementary items and (2) form an item display strategy to arrange the items in a 2-D page that can lead to maximal reward. Conventional RL methods could recommend a set of items each time, for instance, DQN can recommend a set of items with highest Q-values according to the current state[24]. However, these approaches recommend items based on the same state, which leads to the recommended items to be similar. In practice, a bundling of complementary items may receive higher rewards than recommending all similar items. For instance, in real-time news feed recommendations, a user may want to read diverse topics of interest[40]. In addition, page-wise recommendations need to properly display a set of generated items in a 2-D page. Traditional approaches treat it as a ranking problem, i.e., ranking items into a 1-D list according to the importance of items. In other words, user's most preferred item is posited in the top of list. However, in e-commerce recommender systems, a recommendation page is a 2-D grid rather than a 1-D list. Also eye-tracking studies [33] show that rather than scanning a page in a linear fashion, users do page chunking, i.e., they partition the 2-D page into chunks, and browse the chunk they prefer more. In addition, the set of items and the display strategy are generated separately; hence they may be not optimal to each other. Therefore, page-wise recommendations need principled approaches to simultaneously generate a set of complementary items and the display strategy in a 2-D page.

## 1.3 Contributions

In this paper, we tackle the two aforementioned challenges simultaneously by introducing a novel page-wise recommender system based on deep reinforcement learning. We summarize our major contributions as follows – (1) we introduce a principled approach to generate a set of complementary items and properly display them in one 2-D recommendation page simultaneously; (2) we propose a page-wise recommendation framework DeepPage, which can jointly optimize a page of items by incorporating real-time feedback from users; and (3) we demonstrate the effectiveness of the proposed framework in a real-world e-commerce dataset and validate the effectiveness of the components in DeepPage for accurate recommendations.

## 2 THE PROPOSED FRAMEWORK

In this section, we first give an overview of the proposed Actor-Critic based reinforcement learning recommendation framework

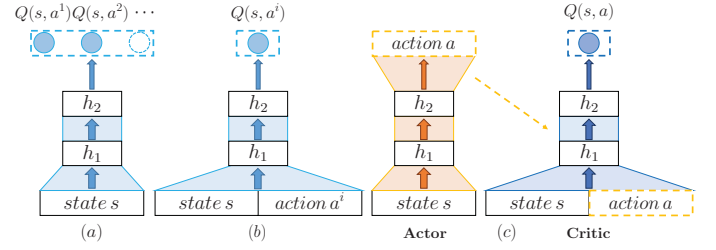


Figure 2: Framework architecture selection.

with notations. Then we present the technical details of components in Actor and Critic, respectively.

### 2.1 Framework Overview

As mentioned in Section 1.1, we model the recommendation task as a Markov Decision Process (MDP) and leverage Reinforcement Learning (RL) to automatically learn the optimal recommendation strategies, which can continuously update recommendation strategies during the interactions and the optimal strategy is made by maximizing the expected long-term cumulative reward from users. With the above intuitions, we formally define the tuple of five elements  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$  of MDP – (a) **State space**  $\mathcal{S}$ : A state  $s \in \mathcal{S}$  is defined as user's current preference, which is generated based on user's browsing history, i.e., the items that a user browsed and her corresponding feedback; (b) **Action space**  $\mathcal{A}$ : An action  $a = \{a^1, \dots, a^M\} \in \mathcal{A}$  is to recommend a page of  $M$  items to a user based on current state  $s$ ; (c) **Reward**  $\mathcal{R}$ : After the RA takes an action  $a$  at the state  $s$ , i.e., recommending a page of items to a user, the user browses these items and provides her feedback. She can skip (not click), click, or purchase these items, and the agent receives immediate reward  $r(s, a)$  according to the user's feedback; (d) **Transition**  $\mathcal{P}$ : Transition  $p(s'|s, a)$  defines the state transition from  $s$  to  $s'$  when RA takes action  $a$ ; and (e) **Discount factor**  $\gamma$ :  $\gamma \in [0, 1]$  defines the discount factor when we measure the present value of future reward. In particular, when  $\gamma = 0$ , RA only considers the immediate reward. In other words, when  $\gamma = 1$ , all future rewards can be counted fully into that of the current action.

Specifically, we model the recommendation task as a MDP in which a recommender agent (RA) interacts with environment  $\mathcal{E}$  (or users) over a sequence of time steps. At each time step, the RA takes an action  $a \in \mathcal{A}$  according to  $\mathcal{E}$ 's state  $s \in \mathcal{S}$ , and receives a reward  $r(s, a)$  (i.e. the RA recommends a page of items according to user's current preference, and receives user's feedback). As the consequence of action  $a$ , the environment  $\mathcal{E}$  updates its state to  $s'$  with transition  $p(s'|s, a)$ . The goal of reinforcement learning is to find a recommendation policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ , which can maximize the cumulative reward for the recommender system.

In practice, conventional RL methods like Q-learning[36, 37] and POMDP[14, 17, 27, 28, 32] become infeasible with the increasing number of items for recommendations. Thus, we leverage Deep Reinforcement Learning[20] with (adapted) artificial neural networks as the non-linear approximators to estimate the action-value function in RL. This model-free reinforcement learning method does not estimate the transition probability and not store the Q-value table. Hence it can support huge amount of items in recommender systems.

There are two major challenges when we apply deep reinforcement learning to the studied problem – (a) the large (or even continuous) and dynamic action space (item space), and (b) the computational cost to select an optimal action (a page of items). In practice, only using discrete indices to denote items is not sufficient since we cannot know the relations between different items only from indices. One common way is to use extra information to represent items with continuous embeddings[19]. Besides, the action space of recommender systems is dynamic as items are arriving and leaving. Moreover, computing Q-value for all state-action pairs is time-consuming because of the enormous state and action spaces.

To tackle these challenges, in this paper, our recommending policy builds upon the Actor-Critic framework [35], shown in Figure 2 (c). The Actor-Critic architecture is preferred from the studied problem since it is suitable for large and dynamic action space, and can also reduce redundant computation simultaneously compared to alternative architectures as shown in Figures 2 (a) and (b). The conventional Deep Q-learning architectures shown in Figure 2 (a) inputs only the state space and outputs Q-values of all actions. This architecture is suitable for the scenario with high state space and small/fixed action space like Atari[24], but cannot handle large and dynamic action space scenario, like recommender systems. Also, we cannot leverage the second conventional deep Q-learning architecture as shown in Fig.2(b) because of its temporal complexity. This architecture inputs a state-action pair and outputs the Q-value correspondingly, and makes use of the optimal action-value function  $Q^*(s, a)$ . It is the maximum expected return achievable by the optimal policy, and should follow the Bellman equation [3] as:

$$Q^*(s, a) = \mathbb{E}_{s'} [r + \gamma \max_{a'} Q^*(s', a') | s, a]. \quad (1)$$

In practice, selecting an optimal  $a'$ ,  $|\mathcal{A}|$  evaluations is necessary for the inner operation “ $\max_{a'}$ ”. In other words, this architecture computes Q-value for all  $a' \in \mathcal{A}$  separately, and then selects the maximal one. This prevents Equation (1) from being adopted in practical recommender systems.

In the Actor-Critic framework, the Actor architecture inputs the current state  $s$  and aims to output a deterministic action (or recommending a deterministic page of  $M$  items), i.e.,  $s \rightarrow a = \{a^1, \dots, a^M\}$ . The Critic inputs only this state-action pair rather than all potential state-action pairs, which avoids the aforementioned computational cost as follows:

$$Q(s, a) = \mathbb{E}_{s'} [r + \gamma Q(s', a') | s, a], \quad (2)$$

where the Q-value function  $Q(s, a)$  is a judgment of whether the selected action matches the current state, i.e., whether the recommendations match user’s preference. Finally, according to the judgment from Critic, the Actor updates its’ parameters in a direction of boosting recommendation performance so as to output proper actions in the next iteration. Next we will elaborate the Actor and Critic architectures.

## 2.2 Architecture of Actor Framework

The Actor is designed to generate a page of recommendations according to user’s preference, which needs to tackle three challenges – 1) setting an initial preference at the beginning of a new recommendation session, 2) learning the real-time preference in the current session, which should capture the dynamic nature of user’s preference in current session and user’s preferable item display patterns in a page, and 3) jointly generating a set of recommendations and displaying them in a 2-D page. To address these challenges, we

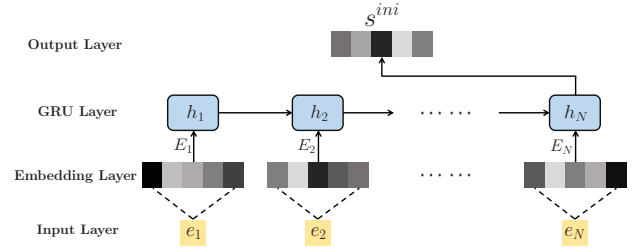


Figure 3: Encoder to generate initial state  $s^{ini}$ .

propose an Actor framework with the Encoder-Decoder architecture.

**2.2.1 Encoder for Initial State Generation Process.** Figure 3 illustrates the model for generating initial preference. We introduce a RNN with Gated Recurrent Units (GRU) to capture users’ sequential behaviors as user’s initial preference. The inputs of GRU are user’s last clicked/purchased items  $\{e_1, \dots, e_N\}$  (sorted in chronological order) before the current session, while the output is the representation of users’ initial preference by a vector. The input  $\{e_1, \dots, e_N\}$  is dense and low-dimensional vector representations of items<sup>1</sup>. We add an item-embedding layer to transform  $e_i$  into a low-dimensional dense vector via  $E_i = \tanh(W_E e_i + b_E) \in \mathbb{R}^{|E|}$  where we use “tanh” activate function since  $e_i \in (-1, +1)$ .

We leverage GRU rather than Long Short-Term Memory (LSTM) because that GRU outperforms LSTM for capturing users’ sequential preference in some recommendation tasks [15]. Unlike LSTM using input gate  $i_t$  and forget gate  $f_t$  to generate a new state, GRU utilizes an update gate  $z_t$ :

$$z_t = \sigma(W_z E_t + U_z h_{t-1}). \quad (3)$$

GRU leverages a reset gate  $r_t$  to control the input of the former state  $h_{t-1}$ :

$$r_t = \sigma(W_r E_t + U_r h_{t-1}). \quad (4)$$

Then the activation of GRU is a linear interpolation between the previous activation  $h_{t-1}$  and the candidate activation  $\hat{h}_t$ :

$$h_t = (1 - z_t)h_{t-1} + z_t \hat{h}_t, \quad (5)$$

where candidate activation function  $\hat{h}_t$  is computed as:

$$\hat{h}_t = \tanh[W E_t + U(r_t \cdot h_{t-1})]. \quad (6)$$

We use the final hidden state  $h_t$  as the representation of the user’s initial state  $s^{ini}$  at the beginning of current recommendation session, i.e.,  $s^{ini} = h_t$ .

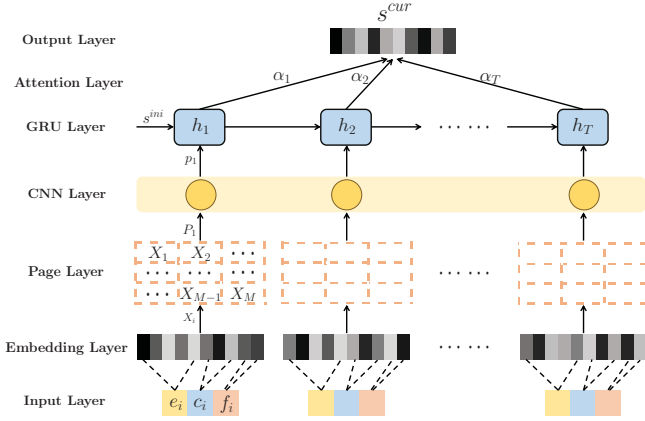
**2.2.2 Encoder for Real-time State Generation Process.** Figure 4 illustrates the model to generate real-time preference in current session. In the page-wise recommender system, the inputs  $\{x_1, \dots, x_M\}$  for each recommendation page are the representations of the items in the page and user’s corresponding feedback, where  $M$  is the size of a recommendation page and  $x_i$  is a tuple as:

$$x_i = (e_i, c_i, f_i), \quad (7)$$

where  $e_i$  is the aforementioned item representation. To assist the RA in capturing user’s preference among different categories of items and generating complementary recommendations, we incorporate item’s category  $c_i$ . The item’s category  $c_i$  is an one-hot

<sup>1</sup>These item representations are pre-trained using users’ browsing history by a company, i.e. each item is treated as a word and the clicked items in one recommendation session as a sentence, and item representations are trained via word embedding[19]. The effectiveness of these item representations is validated by their business such as searching, ranking, bidding and recommendations.





**Figure 4: Encoder to generate real-time state  $s^{cur}$ .**

indicator vector where  $c_i(i) = 1$  if this item belongs to the  $i^{th}$  category and other entities are zero. The one-hot indicator vector is extremely sparse and high-dimensional; hence we add a category-embedding layer transforming  $c_i$  into a low-dimensional dense vector  $C_i = \tanh(W_C c_i + b_C) \in \mathbb{R}^{|C|}$ .

In addition to information from items,  $e_i$  and  $c_i$ , we also want to capture user's interests or feedback in current recommendation page. Thus, we introduce user's feedback vector  $f_i$ , which is an one-hot vector to indicate user's feedback for item  $i$ , i.e., skip/click/purchase. Similarly, we transform  $f_i$  into a dense vector  $F_i = \tanh(W_F f_i + b_F) \in \mathbb{R}^{|F|}$  via the embedding layer. Finally, we get a low-dimensional dense vector  $X_i \in \mathbb{R}^{|X|}$  ( $|X| = |E| + |C| + |F|$ ) by concatenating  $E_i$ ,  $C_i$  and  $F_i$  as:

$$X_i = \text{concat}(E_i, C_i, F_i) \\ = \tanh(\text{concat}(W_E e_i + b_E, W_C c_i + b_C, W_F f_i + b_F)). \quad (8)$$

Note that all item-embedding layers share the same parameters  $W_E$  and  $b_E$ , which reduces the number of parameters and achieves better generalization. We apply the same constraints for category and feedback embedding layers.

Then, we reshape the transformed item representations  $\{X_1, \dots, X_M\}$  as the original arrangement in the page. In other words, we arrange the item representations in one page  $P_t$  as 2D grids similar to one image. For instance, if one recommendation page has  $h$  rows and  $w$  columns ( $M = h \times w$ ), we will get a  $h \times w|X|$  matrix  $P_t$ . To learn item spatial display strategy in one page that leads to maximal reward, we introduce a Convolutional Neural Network (CNN). CNN is a successful architecture in computer vision applications because of its capability to apply various learnable kernel filters on image to discover complex spatial correlations [18]. Hence, we utilize 2D-CNN followed by fully connected layers to learn the optimal item display strategy as:

$$p_t = \text{conv2d}(P_t), \quad (9)$$

where  $p_t$  is a low-dimensional dense vector representing the information from the items and user's feedback in page  $P_t$  as well as the spatial patterns of the item display strategy of page  $P_t$ .

Next, we feed  $\{p_1, \dots, p_T\}$  into another RNN with Gated Recurrent Units (GRU) to capture user's real-time preference in the current session. The architecture of this GRU is similar to the one in Section 2.2.1, but we utilize the final hidden state  $s^{ini}$  in Section 2.2.1 as the initial state in current GRU. Furthermore, to capture the user's real-time preference in the current session, we employ attention mechanism [1], which allows the RA to adaptively focus

on and linearly combine different parts of the input sequence:

$$s^{cur} = \sum_{t=1}^T \alpha_t h_t, \quad (10)$$

where the weighted factors  $\alpha_t$  determine which parts of the input sequence should be emphasized or ignored when making predictions. Here we leverage location-based attention mechanism [22] where the weighted factors are computed from the target hidden state  $h_t$  as follows:

$$\alpha_t = \frac{\exp(W_\alpha h_t + b_\alpha)}{\sum_j \exp(W_\alpha h_j + b_\alpha)}. \quad (11)$$

This GRU with attention mechanism is able to dynamically select more important inputs, which is helpful to capture the user's real-time preference in the current session. Note that - 1) the length of this GRU is flexible according to that of the current recommendation session. After each user-agent interaction, i.e., user browse one page of generated recommendations and give feedback to RA, we can add one more GRU unit, and use this page of items, corresponding categories and feedback as the input of the new GRU unit; 2) in fact, the two RNNs in Section 2.2.1 and Section 2.2.2 can be integrated into one RNN, we describe them separately to clearly illustrate their architecture, and clearly validate their effectiveness in Section 4.4.

**2.2.3 Decoder for Action Generation Process.** In this subsection, we will propose the action  $a^{cur}$  generation process, which generates (recommends) a new page of items to users. In other words, given user's current preference  $s^{cur}$ , we aim to recommend a page of items and displays them properly to maximize the reward. It is the inverse process of what the convolutional layer does. Hence, we use deconvolution neural network (DeCNN) to restore one page from the low-dimensional representation  $s^{cur}$ . It provides a sophisticated and automatic way to generate a page of recommendations with the corresponding display as:

$$a^{cur} = \text{deconv2d}(s^{cur}). \quad (12)$$

Note that - 1) the size of  $a^{cur}$  and  $P$  are different, since  $a^{cur}$  only contains item-embedding  $E_i$ , while  $P$  also contains item's category embedding  $C_i$  and feedback-embedding  $F_i$ . For instance, if one recommendation page has  $h$  rows and  $w$  columns ( $M = h \times w$ ),  $P$  is a  $h \times w|X|$  matrix, while  $a^{cur}$  is a  $h \times w|E|$  matrix; and 2) the generated item embeddings in  $a^{cur}$  may be not in the real item embedding set, thus we need to map them to valid item embeddings, which will be provided in later sections.

### 2.3 The Architecture of Critic Framework

The Critic is designed to leverage an approximator to learn an action-value function  $Q(s, a)$ , which is a judgment of whether the action  $a$  (or a recommendation page) generated by Actor matches the current state  $s$ . Note that we use " $s$ " as  $s^{cur}$  in the last subsection for simplicity. Then, according  $Q(s, a)$ , the Actor updates its' parameters in a direction of improving performance to generate proper actions (or recommendations) in the following iterations.

Thus we need to feed user's current state  $s$  and action  $a$  (or a recommendation page) into the critic. To generate user's current state  $s$ , the RA follows the same strategy from Equation (3) to Equation (10), which uses embedding layers, 2D-CNN and GRU with attention mechanism to capture user's current preference. For action  $a$ , because  $a^{cur}$  generated in Equation (12) is a 2D matrix similar to an image, we utilize the same strategy in Equation (9), a 2D-CNN, to degrade  $a^{cur}$  into a low-dimensional dense vector  $a$

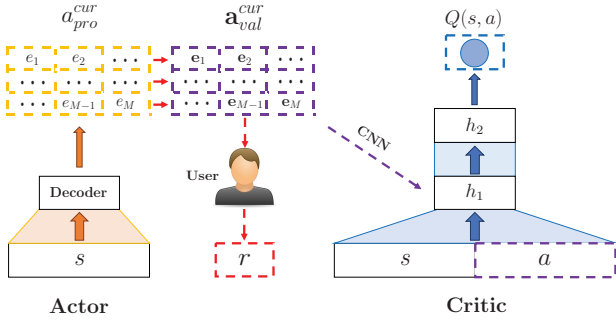


Figure 5: An illustration of the proposed framework.

as:

$$a = \text{conv2d}(a^{\text{cur}}). \quad (13)$$

Then the RA concatenates current state  $s$  and action  $a$ , and feeds them into a Q-value function  $Q(s, a)$ . In real recommender systems, the state and action spaces are enormous, thus estimating the action-value function  $Q(s, a)$  for each state-action pair is infeasible. In addition, many state-action pairs may not appear in the real trace such that it is hard to update their values. Therefore, it is more flexible and practical to use an approximator function to estimate the action-value function. In practice, the action-value function is usually highly nonlinear. Thus we choose Deep neural networks as approximators. In this work, we refer to a neural network approximator as deep Q-value function (DQN).

### 3 TRAINING AND TEST PROCEDURE

In this section, we discuss the training and test procedures. We propose two policies, i.e., online-policy and off-policy, to train and test the proposed framework based on online environment and offline historical data, respectively. Off-policy is necessary because the proposed framework should be pre-trained offline and be evaluated before launching them online to ensure the quality of the recommendations and mitigate possible negative impacts on user experience. After the offline stage, we can apply the framework online, and then the framework can continuously improve its strategies during the interactions with users.

#### 3.1 The Training Procedure

As aforementioned in Section 2.2, we map user's preference  $s^{\text{cur}}$  to a new page of recommendations ( $a^{\text{cur}}$ ). In a page of  $M$  items,  $a^{\text{cur}}$  contains item-embeddings of  $M$  items, i.e.,  $\{e_1, \dots, e_M\}$ . However,  $a^{\text{cur}}$  is a *proto-action*, because the generated item embedding  $e_i \in a^{\text{cur}}$  may be not in the existing item-embedding space  $\mathcal{I}$ . Therefore, we need to map from proto-action  $a_{\text{pro}}^{\text{cur}}$  to *valid-action*  $a_{\text{val}}^{\text{cur}}$  where we have  $\{e_i \in \mathcal{I} | \forall e_i \in a_{\text{val}}^{\text{cur}}\}$ . With this modification, an illustration of the proposed Actor-Critic recommending framework is demonstrated in Figure 5 where we omit Encoders for state generation part.

**3.1.1 Online Training Procedure.** When we train the proposed framework in online environment, RA can interact with users by sequentially choosing recommendation items over a sequence of time steps. Thus, in online environment, the RA is able to receive real-time feedback for the recommended items from users. In this setting, for each  $e_i \in a_{\text{pro}}^{\text{cur}}$ , we select the most similar  $e_i \in \mathcal{I}$  as the *valid item-embedding* in  $a_{\text{val}}^{\text{cur}}$ . In this work, we select *cosine similarity* as:

$$e_i = \arg \max_{e \in \mathcal{I}} \frac{e_i^\top \cdot e}{\|e_i\| \|e\|} = \arg \max_{e \in \mathcal{I}} e_i^\top \cdot \frac{e}{\|e\|}. \quad (14)$$

To decrease the computational cost, we pre-compute  $\frac{e}{\|e\|}$  and also adopt item recalling mechanism to reduce the number of relevant items<sup>2</sup>. Note that Equation (14) is of the same complexity as the  $\max_{a'} Q^*(s', a')$  in Equation (1),  $|\mathcal{A}|$ , but each step of evaluation is a cosine similarity instead of a full value-function evaluation [9].

#### Algorithm 1 Mapping Algorithm.

**Input:** User's browsing history, item-embedding space  $\mathcal{I}$ , the size of recommendation page  $M$ .

**Output:** Valid recommendation page  $a_{\text{val}}^{\text{cur}}$ .

- 1: Generate proto-action  $a_{\text{pro}}^{\text{cur}}$  according Eq.(3) to Eq.(10)
- 2: **for**  $m = 1, M$  **do**
- 3:   Select the most similar item as  $e_m$  according to Eq.(14)
- 4:   Add item  $e_m$  into  $a_{\text{val}}^{\text{cur}}$  (at the same location as  $e_m$  in  $a_{\text{pro}}^{\text{cur}}$ )
- 5:   Remove item  $e_m$  from  $\mathcal{I}$
- 6: **end for**
- 7: **return**  $a_{\text{val}}^{\text{cur}}$

We present the mapping algorithm in Algorithm 1. The Actor first generates proto-action  $a_{\text{pro}}^{\text{cur}}$  (line 1). For each  $e_m$  in  $a_{\text{pro}}^{\text{cur}}$ , the RA selects the most similar item in terms of cosine similarity (line 3), and then adds this item into  $a_{\text{val}}^{\text{cur}}$  at the same position as  $e_m$  in  $a_{\text{pro}}^{\text{cur}}$  (line 4). Finally, the RA removes this item from the item-embedding space (line 5), which prevents recommending the same item repeatedly in one recommendation page.

Then the RA recommends the new recommendation page  $a_{\text{val}}^{\text{cur}}$  to user, and receives the immediate feedback (reward) from user. The reward  $r$  is the summation of rewards of all items in this page:

$$r = \sum_{m=1}^M \text{reward}(e_m). \quad (15)$$

**3.1.2 Offline Training Procedure.** When we use user's historical browsing data to train the proposed Actor-Critic framework, user's browsing history, the new recommendation page  $a_{\text{val}}^{\text{cur}}$  and user's corresponding feedback (reward)  $r$  are given in the data. Thus, there is a gap between  $a_{\text{pro}}^{\text{cur}}$  and  $a_{\text{val}}^{\text{cur}}$ , i.e., no matter what proto-action  $a_{\text{pro}}^{\text{cur}}$  outputted by the Actor, the valid-action  $a_{\text{val}}^{\text{cur}}$  is fixed. This will disconnect the Actor and the Critic.

From existing work [9, 20] and Section 3.1.1, we learn that  $a_{\text{pro}}^{\text{cur}}$  and  $a_{\text{val}}^{\text{cur}}$  should be similar, which is the prerequisite to connect the Actor and the Critic for training. Thus, we choose to minimize the difference between  $a_{\text{pro}}^{\text{cur}}$  and  $a_{\text{val}}^{\text{cur}}$ :

$$\min_{\theta} \sum_{b=1}^B (\|a_{\text{pro}}^{\text{cur}} - a_{\text{val}}^{\text{cur}}\|_F^2), \quad (16)$$

where  $B$  is the batch size of samples in each iteration of SGD. Equation(16) updates Actor's parameters in the direction of pushing  $a_{\text{pro}}^{\text{cur}}$  and  $a_{\text{val}}^{\text{cur}}$  to be similar. In each iteration, given user's browsing history, the new recommendation page  $a_{\text{val}}^{\text{cur}}$ , the RA generates

<sup>2</sup>In general, user's preference in current session should be related to user's last clicked/purchased items before the current session(say  $\mathcal{L}$ ). Thus for each item in  $\mathcal{L}$ , we collect a number of most similar items in terms of cosine similarity from the whole item space, and combine all collected items as the initial item-embedding space  $\mathcal{I}$  of current recommendation session. During the current session, when a user clicks or purchases an item, we will also add a number of its most similar items into the item-embedding space  $\mathcal{I}$ .

proto-action  $a_{pro}^{cur}$  and then minimizes the difference between  $a_{pro}^{cur}$  and  $a_{val}^{cur}$ , which can connect the Actor and the Critic. Next, we can follow conventional methods to update the parameters of Actor and Critic. The reward  $r$  is the summation of rewards of all items in page  $a_{val}^{cur}$ .

**3.1.3 Training Algorithm.** In this work, we utilize DDPG [20] algorithm to train the parameters of the proposed Actor-Critic framework. The Critic can be trained by minimizing a sequence of loss functions  $L(\theta^\mu)$  as:

$$L(\theta^\mu) = \mathbb{E}_{s, a, r, s'} [(r + \gamma Q_{\theta^\mu}(s', f_{\theta^{\pi'}}(s')) - Q_{\theta^\mu}(s, a))^2], \quad (17)$$

where  $\theta^\mu$  represents all parameters in Critic. The critic is trained from samples stored in a replay buffer [25]. Actions stored in the replay buffer are generated by valid-action  $a_{val}^{cur}$ , i.e.,  $a = \text{conv2d}(a_{val}^{cur})$ . This allows the learning algorithm to leverage the information of which action was actually executed to train the critic [9].

The first term  $y = r + \gamma Q_{\theta^\mu}(s', f_{\theta^{\pi'}}(s'))$  in Equation (17) is the target for the current iteration. The parameters from the previous iteration  $\theta^{\mu'}$  are fixed when optimizing the loss function  $L(\theta^\mu)$ . In practice, it is often computationally efficient to optimize the loss function by stochastic gradient descent, rather than computing the full expectations in the above gradient. The derivatives of loss function  $L(\theta^\mu)$  with respect to parameters  $\theta^\mu$  are presented as follows:

$$\nabla_{\theta^\mu} L(\theta^\mu) = \mathbb{E}_{s, a, r, s'} [(r + \gamma Q_{\theta^{\mu'}}(s', f_{\theta^{\pi'}}(s')) - Q_{\theta^\mu}(s, a)) \nabla_{\theta^\mu} Q_{\theta^\mu}(s, a)]. \quad (18)$$

We update the Actor using the policy gradient:

$$\nabla_{\theta^\pi} f_{\theta^\pi} \approx \mathbb{E}_s [\nabla_{\hat{a}} Q_{\theta^\mu}(s, \hat{a}) \nabla_{\theta^\pi} f_{\theta^\pi}(s)], \quad (19)$$

where  $\hat{a} = f_{\theta^\pi}(s)$ , i.e.,  $\hat{a}$  is generated by proto-action  $a_{pro}^{cur}$  ( $\hat{a} = \text{conv2d}(a_{pro}^{cur})$ ). Note that proto-action  $a_{pro}^{cur}$  is the actual action outputted by Actor. This guarantees that policy gradient is taken at the actual output of policy  $f_{\theta^\pi}$  [9].

The online training algorithm for the proposed framework DeepPage is presented in Algorithm 2. In each iteration, there are two stages, i.e., 1) transition generating stage (lines 7-10), and 2) parameter updating stage (lines 11-16). For transition generating stage (line 7): given the current state  $s_t$ , the RA first recommends a page of items  $a_t$  according to Algorithm 1 (line 8); then the RA observes the reward  $r_t$  and updates the state to  $s_{t+1}$  (lines 9); and finally the RA stores transitions  $(s_t, a_t, r_t, s_{t+1})$  into the replay buffer  $\mathcal{D}$  (line 10). For parameter updating stage (line 11): the RA samples minibatch of transitions  $(s, a, r, s')$  from  $\mathcal{D}$  (line 12), and then updates parameters of Actor and Critic (lines 13-16) following a standard DDPG procedure [20].

The offline training procedure is similar with Algorithm 2. The two differences are: 1) in line 8, offline training follows off-policy  $b(s_t)$ , and 2) before line 13, offline training first minimizes the difference between  $a_{pro}^{cur}$  and  $a_{val}^{cur}$  according to Equation (16).

## 3.2 The Test Procedure

After the training stage, the proposed framework learns parameters  $\Theta^\pi$  and  $\Theta^\mu$ . Now we formally present the test procedure of the proposed framework DeepPage. We design two test methods, i.e., 1) Online test: to test DeepPage in online environment where RA interacts with users and receive real-time feedback for the recommended items from users, and 2) Offline test: to test DeepPage based on user's historical browsing data.

---

### Algorithm 2 Parameters Online Training for DeepPage with DDPG.

---

- 1: Initialize actor network  $f_{\theta^\pi}$  and critic network  $Q_{\theta^\mu}$  with random weights
  - 2: Initialize target network  $f_{\theta^{\pi'}}$  and  $Q_{\theta^{\mu'}}$  with weights  $\theta^{\pi'} \leftarrow \theta^\pi, \theta^{\mu'} \leftarrow \theta^\mu$
  - 3: Initialize the capacity of replay buffer  $\mathcal{D}$
  - 4: **for**  $session = 1, G$  **do**
  - 5:   Receive initial observation state  $s_1$
  - 6:   **for**  $t = 1, T$  **do**
  - 7:     **Stage 1: Transition Generating Stage**
  - 8:     Select an action  $a_t$  according to **Alg.1** (policy  $f_{\theta^\pi}$ )
  - 9:     Execute action  $a_t$  and observe the reward  $r_t$  according to **Eq. (15)** and new state  $s_{t+1}$  according to **Section 2.2.2**
  - 10:     Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$
  - 11:     **Stage 2: Parameter Updating Stage**
  - 12:     Sample minibatch of  $N$  transitions  $(s, a, r, s')$  from  $\mathcal{D}$
  - 13:     Set  $y = r + \gamma Q_{\theta^{\mu'}}(s', f_{\theta^{\pi'}}(s'))$
  - 14:     Update Critic by minimizing  $\frac{1}{N} \sum_n (y - Q_{\theta^\mu}(s, a))^2$  according to **Eq. (18)**
  - 15:     Update Actor using the sampled policy gradient according to **Eq. (19)**
  - 16:     Update the target networks:
 
$$\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$$

$$\theta^{\pi'} \leftarrow \tau \theta^\pi + (1 - \tau) \theta^{\pi'}$$
  - 17:   **end for**
  - 18: **end for**
- 

**3.2.1 Online Test.** The online test algorithm in one recommendation session is presented in Algorithm 3. The online test procedure is similar with the transition generating stage in Algorithm 2. In each iteration of the recommendation session, given the current state  $s_t$ , the RA recommends a page of recommendations  $a_t$  to user following policy  $f_{\Theta^\pi}$  (line 4). Then the RA observes the reward  $r_t$  from user (line 5) and updates the state to  $s_{t+1}$  (line 6).

---

### Algorithm 3 Online Test for DeepPage.

---

- 1: Initialize Actor with the trained parameters  $\Theta^\pi$
  - 2: Receive initial observation state  $s_1$
  - 3: **for**  $t = 1, T$  **do**
  - 4:   Execute an action  $a_t$  according to **Alg.1** (policy  $f_{\Theta^\pi}$ )
  - 5:   Observe the reward  $r_t$  from user according to **Eq. (15)**
  - 6:   Observe new state  $s_{t+1}$  according to **Section 2.2.2**
  - 7: **end for**
- 

**3.2.2 Offline Test.** The intuition of our offline test method is that, for a given recommendation session (offline data), the RA reranks the items in this session. If the proposed framework works well, the clicked/purchased items in this session will be ranked at the top of the new list. The reason why RA only reranks items in this session rather than items in the whole item space is that for the offline dataset, we only have the ground truth rewards of the existing items in this session. The offline test algorithm in one recommendation session is presented in Algorithm 4. In each iteration of an offline test recommendation session, given the state  $s_t$  (line 2), the RA recommends an page of recommendations  $a_t$  following policy  $f_{\Theta^\pi}$  (lines 4). For each item  $e_i$  in  $a_t$ , we add it into



new recommendation list  $\mathcal{L}$  (line 6), and record  $e_i$ 's reward  $r_i$  from user's historical browsing data (line 7). Then we can compute the overall reward  $r_t$  of  $a_t$  (line 9) and update the state to  $s_{t+1}$  (line 10). Finally, we remove all items  $e_i$  in  $a_t$  from the item set  $\mathcal{I}$  of the current session (line 11).

---

**Algorithm 4** Offline Test of DeepPage Framework.

---

**Input:** Item set  $\mathcal{I} = \{e_1, \dots, e_N\}$  and corresponding reward set  $\mathcal{R} = \{r_1, \dots, r_N\}$  of a session.

**Output:** Recommendation list  $\mathcal{L}$  with new order

- 1: Initialize Actor with well-trained parameters  $\Theta^\pi$
  - 2: Receive initial observation state  $s_1$
  - 3: **while**  $|\mathcal{I}| > 0$  **do**
  - 4:   Select an action  $a_t$  according to **Alg.1** (policy  $f_{\Theta^\pi}$ )
  - 5:   **for**  $e_i \in a_t$  **do**
  - 6:     Add  $e_i$  into the end of  $\mathcal{L}$
  - 7:     Record reward  $r_i$  from user's historical browsing data
  - 8:   **end for**
  - 9:   Compute the overall reward  $r_t$  of  $a_t$  according to **Eq. (15)**
  - 10:   Execute action  $a_t$  and observe new state  $s_{t+1}$  according to **Section 2.2.2**
  - 11:   Remove all  $e_i \in a_t$  from  $\mathcal{I}$
  - 12: **end while**
- 

## 4 EXPERIMENTS

In this section, we conduct extensive experiments with a dataset from a real e-commerce company to evaluate the effectiveness of the proposed framework. We mainly focus on two questions: (1) how the proposed framework performs compared to representative baselines; and (2) how the components in Actor and Critic contribute to the performance.

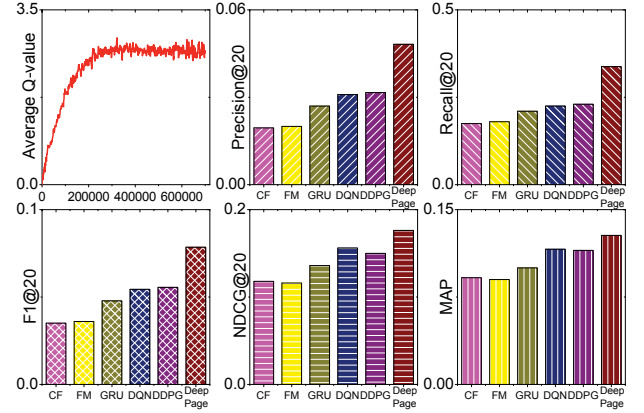
### 4.1 Experimental Settings

We evaluate our method on a dataset of September, 2017 from a real e-commerce company. We randomly collect 1,000,000 recommendation sessions (9,136,976 items) in temporal order, and use the first 70% sessions as the training/validation set and the later 30% sessions as the test set. For a new session, the initial state is collected from the previous sessions of the user. In this work, we leverage  $N = 10$  previously clicked/purchased items to generate the initial state. Each time the RA recommends a page of  $M = 10$  items (5 rows and 2 columns) to users<sup>3</sup>. The reward  $r$  of one skipped/clicked/purchased item is empirically set as 0, 1, and 5, respectively. The dimensions of item-embedding/category-embedding/ feedback-embedding are  $|E| = 50$ ,  $|C| = 35$ , and  $|F| = 15$ . We set the discounted factor  $\gamma = 0.95$ , and the rate for soft updates of target networks  $\tau = 0.01$ . For offline test, we select **Precision@20**, **Recall@20**, **F1-score@20** [12], **NDCG@20** [16] and **MAP** [38], as the metrics. For online test, we leverage the summation of all rewards in one recommendation session as the metric.

### 4.2 Performance Comparison for Offline Test

To answer the first question, we compare the proposed framework with the following representative baseline methods – **CF** [4]: collaborative filtering; **FM** [29]: Factorization Machines; **GRU** [15]: it utilizes the Gated Recurrent Units (GRU) to predict what user will click/purchase next based on the browsing histories (to make a fair

<sup>3</sup>This is based on offline historical data collected from mobile App, i.e., to fit the screen size of mobile phones, one page has only 5 rows and 2 columns.



**Figure 6: Overall performance comparison in offline test.**

comparison, it also keeps previous  $N = 10$  clicked/purchased items as initial states); **DQN** [24]: we use a Deep Q-network with five fully-connected layers in this baseline. The input is the concatenation of embeddings of users' historical clicked/purchased items (state) and a page of recommendations (action), and train this baseline by Eq.(1); and **DDPG** [9]: in this baseline, we use conventional Deep Deterministic Policy Gradient with five fully connected layers in both Actor and Critic.

We leverage offline training strategy to train DDPG and DeepPage as mentioned in Section 3.1.1. The results are shown in Figure 6. Figure 6 (a) illustrates the training process of DeepPage. We can observe that the framework approaches convergence when the model is trained by 500,000 offline sessions. We note that CF and FM perform worse than other baselines. These two baselines ignore the temporal sequence of the users' browsing history, while GRU can capture the temporal sequence, and DQN, DDPG and DeepPage are able to continuously update their strategies during the interactions. DQN and DDPG outperform GRU. We design GRU to maximize the immediate reward for recommendations, while DQN and DDPG are designed to achieve the trade-off between short-term and long-term rewards. This result suggests that introducing reinforcement learning can improve the performance of recommendations. Finally DeepPage performs better than conventional DDPG. Compared to DDPG, DeepPage jointly optimizes a page of items and uses GRU to learn user's real-time preference.

### 4.3 Performance Comparison for Online Test

Following [42], we build a simulated online environment (adapted to our case) for online test. We compare DeepPage with GRU, DQN and DDPG. Here we utilize online training strategy to train DDPG and DeepPage (both Actor-Critic framework) as mentioned in Section 3.1.2. Baselines are also applicable to be trained via the rewards generated by simulated online environment. Note that we use data different from the training set to build the simulated online environment to avoid over-fitting.

As the test stage is based on the simulator, we can artificially control the length of recommendation sessions to study the performance in short and long sessions. We define short sessions with 10 recommendation pages, while long sessions with 50 recommendation pages. The results are shown in Figure 7. We note DDPG performs similar to DQN, but the training speed of DDPG is much

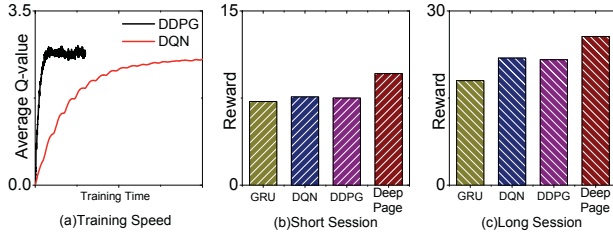


Figure 7: Overall performance comparison in online test.

faster than DQN, as shown in Figure 7 (a). This result indicates that Actor-Critic framework is suitable for practical recommender systems with enormous action space. In short recommendation sessions, GRU, DQN and DDPG achieve comparable performance. In other words, GRU model and reinforcement learning models like DQN and DDPG can recommend proper items matching users' short-term interests. In long recommendation sessions, DQN and DDPG outperform GRU significantly. GRU is designed to maximize the immediate reward for recommendations, while DQN and DDPG are designed to achieve the trade-off between short-term and long-term rewards. This result suggests that introducing reinforcement learning can improve the long-term performance of recommendations. Finally DeepPage performs better than conventional DQN and DDPG. DeepPage can learn user's real-time preference and optimize a page of items. We detail the effect of model components of DeepPage in the following subsection.

#### 4.4 Effectiveness of Components

To validate the effectiveness of each component, we systematically eliminate the corresponding model components by defining following variants of DeepPage – (1) **DeepPage-1**: this variant is to evaluate the performance of the embedding layers. We remove the embedding layers for items, categories and feedback; (2) **DeepPage-2**: in this variant, we evaluate the contribution of category and feedback information, hence, this variant does not contain one-hot indicator vectors of category and feedback; (3) **DeepPage-3**: This variant is to evaluate the effectiveness of GRU to generate initial state, so we eliminate the GRU in Figure 3; (4) **DeepPage-4**: In this variant, we evaluate the contribution of CNNs as shown in Figure 4, thus we remove CNNs and directly feed the outputs of embedding layers (concatenate embeddings of all items as one vector) into GRU units; (5) **DeepPage-5**: This variant is to evaluate the effectiveness of attention mechanism in Figure 4, therefore, we eliminate attention layer and use the hidden state of last GRU unit as the input of DeCNN; (6) **DeepPage-6**: In this variant, we evaluate the GRU to generate real-time state, thereby, we remove this GRU in Figure 4 and concatenate outputs of all CNNs as a vector, and feed it into DeCNN; and (7) **DeepPage-7**: This variant is to evaluate the performance of DeCNN to generate a new page of items, hence, we replace it by fully-connected layers, which output a concatenated vector of  $M$  item-embeddings.

The offline results are shown in Table 1. Note that we omit similar online observations because of the space limitation. DeepPage-1 and DeepPage-2 validate that incorporating category/feedback information and the embedding layers can boost the performance. DeepPage-3 and DeepPage-6 perform worse, which suggests that

Table 1: Performance comparison of different components.

Methods	Precision @20	Recall @20	F1score @20	NDCG @20	MAP
DeepPage-1	0.0479	0.3351	0.0779	0.1753	0.1276
DeepPage-2	0.0475	0.3308	0.0772	0.1737	0.1265
DeepPage-3	0.0351	0.2627	0.0578	0.1393	0.1071
DeepPage-4	0.0452	0.3136	0.0729	0.1679	0.1216
DeepPage-5	0.0476	0.3342	0.0775	0.1716	0.1243
DeepPage-6	0.0318	0.2433	0.0528	0.1316	0.1039
DeepPage-7	0.0459	0.3179	0.0736	0.1698	0.1233
<b>DeepPage</b>	<b>0.0491</b>	<b>0.3576</b>	<b>0.0805</b>	<b>0.1872</b>	<b>0.1378</b>

setting user's initial preference at the beginning of a new recommendation session, and capturing user's real-time preference in current session is helpful for accurate recommendations. DeepPage-5 proves that incorporating attention mechanism can better capture user's real-time preference than only GRU. DeepPage outperforms DeepPage-4 and DeepPage-7, which indicates that item display strategy can influence the decision making process of users. In a nutshell, DeepPage outperforms all its variants, which demonstrates the effectiveness of each component for recommendations.

## 5 RELATED WORK

In this section, we briefly review research related to our study, i.e., reinforcement learning for recommendations. The MDP-Based CF model can be viewed as approximating a partial observable MDP [32]. Mahmood et al. [23] adopted the RL technique to observe the responses of users in a conversational recommender. Taghipour et al. [36, 37] modeled web page recommendation as a Q-Learning problem and learned to make recommendations from web usage data. Sunehag et al. [34] addressed sequential decision problems with high-dimensional slate-action spaces. Zheng et al. [44] proposed a RL framework to do online news recommendation. Cai et al. [5, 6] employed a RL mechanism for impression allocation problem. Feng et al. [10] optimized ranking strategies collaboratively for multi-scenario recommendations. Choi et al. [8] used RL and biclustering for cold-start problem. Chen et al. [7] proposed strategies to improve the reward estimation in dynamic recommendation environments. Zhao et al. [42] captured both positive and negative feedback into RL based recommendations. Zhao et al. [43] optimized a set of recommendations with proper order.

## 6 CONCLUSION

In this paper, we propose a novel page-wise recommendation framework, DeepPage, which leverages Deep Reinforcement Learning to automatically learn the optimal recommendation strategies and optimizes a page of items simultaneously. We validate the effectiveness of our framework with extensive experiments based on data from a real e-commerce company. Further research directions include reducing the temporal complexity of mapping from proto-action to valid-action and handling multiple tasks such as search, bidding, advertisement and recommendation collaboratively in one reinforcement learning framework.

## ACKNOWLEDGEMENTS

This material is based upon work supported by the National Science Foundation (NSF) under grant number IIS-1714741 and IIS-1715940, and a grant from Criteo Faculty Research Award.



## REFERENCES

- [1] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).
- [2] Jie Bao, Yu Zheng, David Wilkie, and Mohamed Mokbel. 2015. Recommendations in location-based social networks: a survey. *Geoinformatica* 19, 3 (2015), 525–565.
- [3] Richard Bellman. 2013. *Dynamic programming*. Courier Corporation.
- [4] John S Breese, David Heckerman, and Carl Kadie. 1998. Empirical analysis of predictive algorithms for collaborative filtering. In *Proceedings of the Fourteenth conference on Uncertainty in artificial intelligence*. Morgan Kaufmann Publishers Inc., 43–52.
- [5] Qingpeng Cai, Aris Filos-Ratsikas, Pingzhong Tang, and Yiwei Zhang. 2018. Reinforcement Mechanism Design for e-commerce. In *Proceedings of the 2018 World Wide Web Conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 1339–1348.
- [6] Qingpeng Cai, Aris Filos-Ratsikas, Pingzhong Tang, and Yiwei Zhang. 2018. Reinforcement mechanism design for fraudulent behaviour in e-commerce. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence*.
- [7] Shi-Yong Chen, Yang Yu, Qing Da, Jun Tan, Hai-Kuan Huang, and Hai-Hong Tang. 2018. Stabilizing reinforcement learning in dynamic environment with application to online recommendation. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 1187–1196.
- [8] Sungwoon Choi, Heonseok Ha, Uiwon Hwang, Chanju Kim, Jung-Woo Ha, and Sungroh Yoon. 2018. Reinforcement Learning based Recommender System using Biclustering Technique. *arXiv preprint arXiv:1801.05532* (2018).
- [9] Gabriel Dulac-Arnold, Richard Evans, Hado van Hasselt, Peter Sunehag, Timothy Lillicrap, Jonathan Hunt, Timothy Mann, Theophane Weber, Thomas Degris, and Ben Coppin. 2015. Deep reinforcement learning in large discrete action spaces. *arXiv preprint arXiv:1512.07679* (2015).
- [10] Jun Feng, Heng Li, Minlie Huang, Shichen Liu, Wenwu Ou, Zhirong Wang, and Xiaoyan Zhu. 2018. Learning to Collaborate: Multi-Scenario Ranking via Multi-Agent Reinforcement Learning. (2018).
- [11] Huiji Gao, Jiliang Tang, Xia Hu, and Huan Liu. 2015. Content-Aware Point of Interest Recommendation on Location-Based Social Networks.. In *AAAI*. 1721–1727.
- [12] Asela Gunawardana and Guy Shani. 2009. A survey of accuracy evaluation metrics of recommendation tasks. *Journal of Machine Learning Research* 10, Dec (2009), 2935–2962.
- [13] Hao Guo, Xin Li, Ming He, Xiangyu Zhao, Guiquan Liu, and Guandong Xu. 2016. CoSoLoRec: Joint Factor Model with Content, Social, Location for Heterogeneous Point-of-Interest Recommendation. In *International Conference on Knowledge Science, Engineering and Management*. Springer, 613–627.
- [14] Milos Hauskrecht. 1997. Incremental methods for computing bounds in partially observable Markov decision processes. In *AAAI/IAAI*. 734–739.
- [15] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Dávid Szepesvári. 2015. Session-based recommendations with recurrent neural networks. *arXiv preprint arXiv:1511.06939* (2015).
- [16] Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated gain-based evaluation of IR techniques. *ACM Transactions on Information Systems (TOIS)* 20, 4 (2002), 422–446.
- [17] Michael Kearns, Yishay Mansour, and Andrew Y Ng. 2002. A sparse sampling algorithm for near-optimal planning in large Markov decision processes. *Machine learning* 49, 2 (2002), 193–208.
- [18] 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.
- [19] Omer Levy and Yoav Goldberg. 2014. Neural word embedding as implicit matrix factorization. In *Advances in neural information processing systems*. 2177–2185.
- [20] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. 2015. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971* (2015).
- [21] Greg Linden, Brent Smith, and Jeremy York. 2003. Amazon. com recommendations: Item-to-item collaborative filtering. *IEEE Internet computing* 7, 1 (2003), 76–80.
- [22] Thang Luong, Hieu Pham, and Christopher D Manning. 2015. Effective approaches to attention-based neural machine translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. 1412–1421.
- [23] Tariq Mahmood and Francesco Ricci. 2009. Improving recommender systems with adaptive conversational strategies. In *Proceedings of the 20th ACM conference on Hypertext and hypermedia*. ACM, 73–82.
- [24] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. 2013. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602* (2013).
- [25] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. *Nature* 518, 7540 (2015), 529.
- [26] Raymond J Mooney and Loriene Roy. 2000. Content-based book recommending using learning for text categorization. In *Proceedings of the fifth ACM conference on Digital libraries*. ACM, 195–204.
- [27] Andrew Y Ng and Michael Jordan. 2000. PEGASUS: A policy search method for large MDPs and POMDPs. In *Proceedings of the Sixteenth conference on Uncertainty in artificial intelligence*. Morgan Kaufmann Publishers Inc., 406–415.
- [28] Pascal Poupart and Craig Boutilier. 2005. VDCBPI: an approximate scalable algorithm for large POMDPs. In *Advances in Neural Information Processing Systems*. 1081–1088.
- [29] Steffen Rendle. 2010. Factorization machines. In *Data Mining (ICDM), 2010 IEEE 10th International Conference on*. IEEE, 995–1000.
- [30] Paul Resnick and Hal R Varian. 1997. Recommender systems. *Commun. ACM* 40, 3 (1997), 56–58.
- [31] Francesco Ricci, Lior Rokach, and Bracha Shapira. 2011. Introduction to recommender systems handbook. In *Recommender systems handbook*. Springer, 1–35.
- [32] Guy Shani, David Heckerman, and Ronen I Brafman. 2005. An MDP-based recommender system. *Journal of Machine Learning Research* 6, Sep (2005), 1265–1295.
- [33] Ramakrishnan Srikant, Sugato Basu, Ni Wang, and Daryl Pregibon. 2010. User browsing models: relevance versus examination. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 223–232.
- [34] Peter Sunehag, Richard Evans, Gabriel Dulac-Arnold, Yori Zwols, Daniel Visentin, and Ben Coppin. 2015. Deep Reinforcement Learning with Attention for Slate Markov Decision Processes with High-Dimensional States and Actions. *arXiv preprint arXiv:1512.01124* (2015).
- [35] Richard S Sutton and Andrew G Barto. 1998. *Reinforcement learning: An introduction*. Vol. 1. MIT press Cambridge.
- [36] Nima Taghipour and Ahmad Kardan. 2008. A hybrid web recommender system based on q-learning. In *Proceedings of the 2008 ACM symposium on Applied computing*. ACM, 1164–1168.
- [37] Nima Taghipour, Ahmad Kardan, and Saeed Shiry Ghidary. 2007. Usage-based web recommendations: a reinforcement learning approach. In *Proceedings of the 2007 ACM conference on Recommender systems*. ACM, 113–120.
- [38] Andrew Turpin and Falk Scholer. 2006. User performance versus precision measures for simple search tasks. In *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 11–18.
- [39] 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.
- [40] Yisong Yue and Carlos Guestrin. 2011. Linear submodular bandits and their application to diversified retrieval. In *Advances in Neural Information Processing Systems*. 2483–2491.
- [41] Xiangyu Zhao, Tong Xu, Qi Liu, and Hao Guo. 2016. Exploring the Choice Under Conflict for Social Event Participation. In *International Conference on Database Systems for Advanced Applications*. Springer, 396–411.
- [42] Xiangyu Zhao, Liang Zhang, Zhuoye Ding, Long Xia, Jiliang Tang, and Dawei Yin. 2018. Recommendations with Negative Feedback via Pairwise Deep Reinforcement Learning. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 1040–1048.
- [43] Xiangyu Zhao, Liang Zhang, Zhuoye Ding, Dawei Yin, Yihong Zhao, and Jiliang Tang. 2017. Deep Reinforcement Learning for List-wise Recommendations. *arXiv preprint arXiv:1801.00209* (2017).
- [44] Guanjie Zheng, Fuzheng Zhang, Zihan Zheng, Yang Xiang, Nicholas Jing Yuan, Xing Xie, and Zhenhui Li. 2018. DRN: A Deep Reinforcement Learning Framework for News Recommendation. In *Proceedings of the 2018 World Wide Web Conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 167–176.
- [45] Yu Zheng, Lizhu Zhang, Zhengxin Ma, Xing Xie, and Wei Ying Ma. 2011. Recommending friends and locations based on individual location history. *ACM Transactions on the Web* 5, 1 (2011), 5.