

AutoAssign+: Automatic Shared Embedding Assignment in Streaming Recommendation

Ziru Liu^{1†}, Kecheng Chen^{1†}, Fengyi Song¹, Bo Chen², Xiangyu Zhao^{1*}, Huifeng Guo² and Ruiming Tang^{2*}

¹City University of Hong Kong, Kowloon Tong, Hong Kong.

²Huawei Noah's Ark Lab, Shenzhen, China.

*Corresponding author(s). E-mail(s): xianzhao@cityu.edu.hk;
tangruiming@huawei.com;

Contributing authors: ziruliu2-c@my.cityu.edu.hk;
kechechen4-c@my.cityu.edu.hk; fengysong2-c@my.cityu.edu.hk;
chenbo116@huawei.com; huifeng.guo@huawei.com;

[†]These authors contributed equally to this work.

Abstract

In the domain of streaming recommender systems, conventional methods for addressing new user IDs or item IDs typically involve assigning initial ID embeddings randomly. However, this practice results in two practical challenges: (i) Items or users with limited interactive data may yield suboptimal prediction performance. (ii) Embedding new IDs or low-frequency IDs necessitates consistently expanding the embedding table, leading to unnecessary memory consumption. In light of these concerns, we introduce a reinforcement learning-driven framework, namely AutoAssign+, that facilitates Automatic Shared Embedding Assignment Plus. To be specific, AutoAssign+ utilizes an *Identity Agent* as an actor network, which plays a dual role: (i) Representing low-frequency IDs field-wise with a small set of shared embeddings to enhance the embedding initialization, and (ii) Dynamically determining which ID features should be retained or eliminated in the embedding table. The policy of the agent is optimized with the guidance of a *critic network*. To evaluate the effectiveness of our approach, we perform extensive experiments on three commonly used benchmark datasets. Our experiment results demonstrate that AutoAssign+ is capable of significantly enhancing recommendation performance by mitigating the cold-start problem. Furthermore, our framework yields

a reduction in memory usage of approximately 20-30%, verifying its practical effectiveness and efficiency for streaming recommender systems.

Keywords: Recommender Systems, Reinforcement Learning, Cold-Start, Streaming Recommendation

1 Introduction

With the rapid growth of personalized online applications, recommender systems have been widely implemented by various online businesses, including E-commerce websites, news platforms, online advertising, and so on [1, 2]. Among them, streaming recommendation [3, 4] is one of the common forms of recommender systems, where streaming data are constantly flowing into the recommendation models for training, thus better modeling the user’s current preferences. In addition, streaming recommendations are particularly important for time-sensitive items, such as news, as they allow for rapid identification and distribution of relevant content to interested users, which is critical for commercial information retrieval systems. Due to the ability to effectively capture the highly nonlinear relationship between user and item end-to-end, neural network-based models are rapidly becoming the mainstream of recommender systems. As shown in Fig. 1, existing deep recommendation models typically follow the “Embedding & Feature Interaction” paradigm [5]. The embedding layer serves as the encoder to represent sparse features in dense latent space, while the feature interaction layers serve to capture interactive signals among these features.

In a streaming recommender system, new items and users are continually added to the data corpus, creating a highly dynamic streaming environment that presents several challenges, which can be summarized as:

- **Cold-start:** The streaming recommender system is confronted with a constant influx of new users, many of whom can be classified as visitor-type users and possess extremely limited behavior information. Furthermore, the system is constantly updated with new items, yet there has not been enough interaction with these items to generate an adequate level of training data. The consequence of employing insufficiently trained new user/item embeddings is a significant decline in the performance of the recommendation model.
- **Interval interaction:** Users’ preferences and activity levels exhibit temporal fluctuations, manifesting as periods of “*active* → *inactive* → *active* → ...” behavior. Similarly, items display a comparable trend over time. In order to predict the reactivation of users/items accurately, it is crucial to maintain the inactive user/item ID parameters updated with the recommendation model in real-time. Additionally, effectively recycling ID parameters

of long-term inactive users/items can prevent the model from becoming overburdened with redundant data, leading to a significant reduction in model size and conserving space.

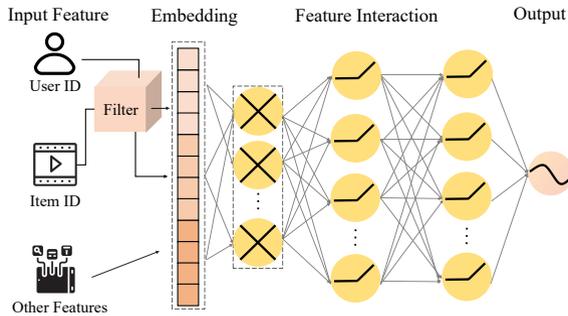


Fig. 1: A typical deep learning based recommendation model.

The challenges faced by streaming recommender systems can be attributed to the handling of (temporary) low-frequency user/item ID features. It is impractical and problematic to consider high and low-frequency IDs as equal entities. In industrial streaming recommender systems, a typical approach to address low-frequency IDs is to replace their embeddings temporarily with shared embeddings. As shown in Fig. 1, a pre-defined frequency threshold can be used to filter those low-frequency IDs before the embedding lookup. However, this straightforward approach has a few obstacles when applied to streaming recommendation systems.: (i) Determining optimal pre-defined thresholds for user/item IDs requires human expert knowledge and can be time-consuming; (ii) The distribution of IDs can change over time, resulting in an inadequate performance when a fixed threshold is used. One potential solution is to use an adaptive threshold that can be adjusted dynamically, based on the current distribution of IDs.

Recent research efforts have focused on addressing the cold-start problem in streaming recommender systems using auxiliary information [6–8]. Additionally, meta-learning has emerged as a popular approach for learning global and general information for relevant tasks, thereby facilitating the initialization of new, related tasks. For example, MetaEmb [9] generates an initial ID embedding for new items using item features, while MWUF [10] employs the average of all existing ID embeddings in conjunction with scaling and shifting functions to initialize cold-start IDs. However, the auxiliary information methods [6–8] rely on rich user profiles or item attributes, such as user social networks or item images, which may not be available in scenarios involving only user/item IDs. Moreover, meta-embedding-based approaches that use global parameters in the feature interaction layer or pre-trained parameters of the embedding table to enhance initialization for new IDs can introduce deviations in the embeddings of high-frequency IDs.

To tackle the aforementioned difficulties, we propose Automatic Shared Embedding Assignment Plus (**AutoAssign+**), a reinforcement learning-based approach, which utilizes an **Identity Agent** to dynamically and field-wisely assign shared embedding and unique embedding to the users/items ID features that appear anytime in the streaming recommendation scenario. Besides, one **critic network** is applied to optimize and fine-tune the policy of the agent based on TD error. In our previous version, the Identity Agent of **AutoAssign** generates the candidate embeddings based solely on the input features and doesn't take into account the performance of the recommendation model. In AutoAssign+, the critic network is introduced to optimize the Identity Agent by evaluating the quality of the candidate embeddings generated by the Identity Agent, allowing it to adjust and optimize its candidate embedding generation process to improve the recommendation model's performance. Notably, our Automatic Shared Embedding framework is designed to be "plug and play," allowing it to be applied to various deep recommendation models.

The main contributions of our work can be concluded as:

- We propose a new framework AutoAssign+ for streaming recommendation scenarios, which is compatible with various deep recommendation models to alleviate the low-frequency ID problem;
- AutoAssign+ utilizes an actor-critic structure from reinforcement learning along with hierarchical shared embeddings to dynamically assign optimal embeddings for each user/item ID based on their occurrence frequency.
- We performed extensive experiments on popular datasets to demonstrate the effectiveness and superiority of AutoAssign+. The results show that AutoAssign+ significantly outperforms the performance of AutoAssign while reducing storage space by 20%-30%.

2 Preliminary

2.1 Deep Recommendation Model

The deep recommendation model utilizes dense vectors (embeddings) of instances as input and makes predictions by feeding those embeddings into subsequent feature interaction components [11, 12]. The instance dense vectors can be constructed by concatenating the vector representations of each sparse (categorical) feature [13]. As illustrated in Fig. 1, for a given user-item pair, their IDs and corresponding features can be represented as:

$$\mathbf{x} = [\mathbf{x}^u, \mathbf{x}^i, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n], \quad (1)$$

where \mathbf{x}_u and \mathbf{x}^i are feature vectors for user and item, n is the number of other feature fields and \mathbf{x}_i is the high-dimensional sparse vector of i^{th} field. Then embedding layer transforms those sparse vectors into dense vectors:

$$\mathbf{e}_i = \mathbf{W}_i \mathbf{x}_i, \quad (2)$$

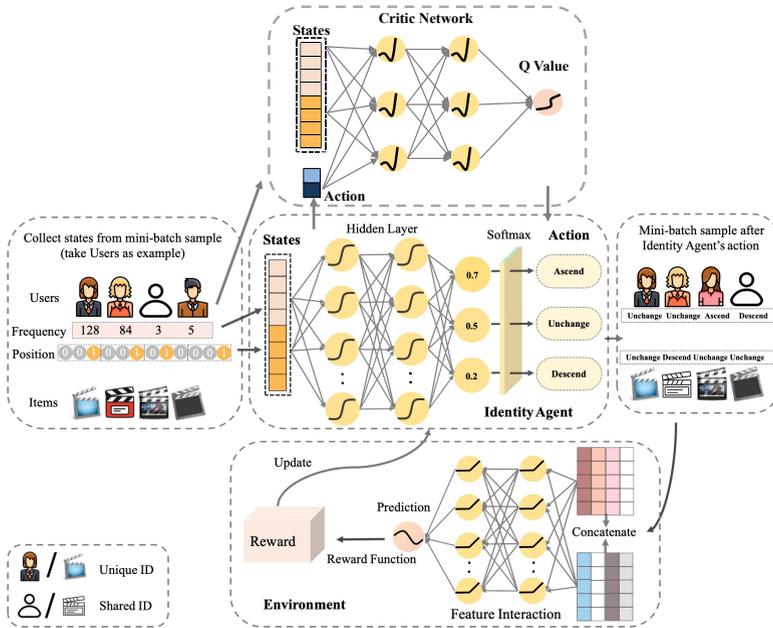


Fig. 2: The proposed overall framework contains an Identity Agent, which serves to assign a shared ID to the User and Item before the embedding layer.

where $\mathbf{W}_i \in \mathbb{R}^{d \times h}$ is a matrix (embedding table), d is the embedding dimension of dense vector and h is the number of feature value in i^{th} field. These embeddings are feed-in into the subsequent feature interaction layer to extract informative interactions. In the binary classification task, the predicted target is $y \in \{0, 1\}$ where 0 refers to dislike and 1 refers to like in explicit user behavior. We assume that the parameters of the deep recommendation model are θ . Depending on the feature interaction function $f(\cdot)$, we get the prediction \hat{y} :

$$\hat{y} = f(\theta, \mathbf{e}^u, \mathbf{e}^i, \mathbf{e}_n), \quad (3)$$

where \mathbf{e}^u and \mathbf{e}^i respectively denote the user ID and item ID embedding, and \mathbf{e}_n is the concatenate of n number other feature fields. The typical optimization target for binary classification tasks is Logloss or Mean Square Error.

In the context of deep recommendation models, user/item ID embeddings are crucial as they encode the latent features of specific users/items [14, 15]. However, in streaming recommendation scenarios, insufficient training of ID embeddings can arise due to two situations. Firstly, the emergence of new users and items is commonly referred to as the cold-start problem. Secondly, certain items or visitor-like users may have an inherently unpopular nature, resulting in a long-tailed distribution issue. The initialization of low-frequency and unpopular ID embeddings with random values cannot provide valid information, ultimately leading to poor recommendation performance.

2.2 Low-Frequency Filter

An intuitive approach to tackle the previously mentioned issue is to temporarily use a shared embedding to represent the low-frequency IDs. As shown in Fig. 1, a frequency threshold is adopted in a real-world industry recommender system to filter those IDs with a frequency lower than the threshold before performing the embedding lookup. Specifically, suppose the frequency threshold for User ID is τ and we set a shared ID ID_{shared} , the User ID is collected as $(ID_i, F(i))$, where i represents the i^{th} user and $F(\cdot)$ represent its current occurrence frequency. In low-frequency filter:

$$(ID_i, F(i)) = \begin{cases} ID_i & \text{if } F(i) > \tau \\ ID_{shared} & \text{if } F(i) \leq \tau \end{cases},$$

By training a shared embedding ID_{shared} on a large number of low-frequency IDs, it can obtain abundant generalized information. However, this method has some practical drawbacks. Firstly, determining a fixed threshold τ to identify low-frequency IDs requires expert knowledge or extensive search. Secondly, the frequency distribution of different ID fields (such as user ID and item ID) varies spatially and temporally, making it sub-optimal to use a fixed threshold to define low-frequency IDs.

3 Framework

The AutoAssign+ framework is designed to address the limitations of existing methods by automatically assigning user/item embeddings to low-frequency targets, which helps to alleviate the negative impact caused by the low-frequency ID features on the recommendation performance. In this section, we will provide a detailed description of our proposed approach.

3.1 Reinforcement Learning Setting

To address the embedding assignment problem, we formulate it as a Markov Decision Process (MDP) that can be analyzed within the reinforcement learning framework where we train the Identity Agent to make optimal embedding assignment decisions.

3.1.1 State

The entirety of the deep recommendation model's pipeline is considered to be the environment. With the available data sample obtained from the stream, the deep recommendation model generated an output using the Identity Agent's policy decision. Under this pipeline, we define a group of hierarchical candidate shared ID with size k , i.e., $\{S_1, S_2, \dots, S_k\}$, and its corresponding shared embedding table $\mathbf{E}_{shared} = [e_{s1}, e_{s2}, \dots, e_{sk}] \in \mathbb{R}^{d \times k}$, where d is the embedding dimension of shared ID and k is the size of shared ID group. In the case of this study, we choose $k = 2$. These shared embeddings are pre-trained from

the same data and parameters randomly. This hierarchical shared embedding setting acts as a buffer before the unique embedding is assigned to a specific user/item in the recommendation process.

The state is characterized as the present frequency of ID features along with their current position p_i , within the shared ID group of candidates (i.e., the i^{th} embedding vector of the shared embedding table), $State = (F(\cdot), p_i)$, where $i \in [1, k + 1]$ and $F(\cdot)$ is the current occurrence frequency of a certain ID. Note that: (i) there is an additional position $k + 1$, which refers to a unique ID already assigned to a certain ID, (ii) the state of a new coming ID is initialized as $State_{init} = (1, 1)$, which means that its current frequency is one and assigned with the first embedding vector in the shared embedding table.

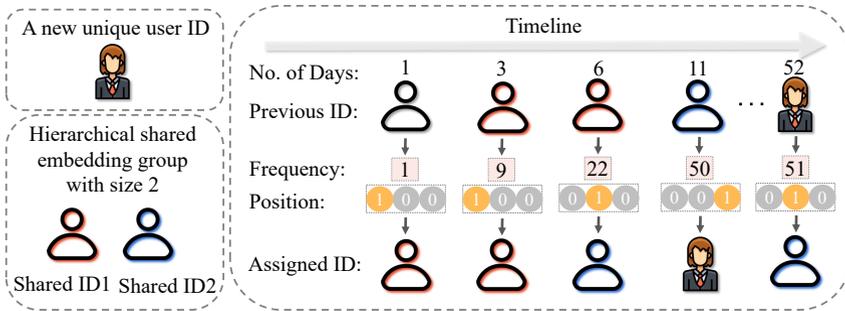


Fig. 3: Assign a shared ID / unique ID to a new user by AutoAssign+.

3.1.2 Action

The design and architecture of the Identity Agents are identical for both the user and the item. We use the user ID agent to illustrate the action as shown in Fig. 2. Given the state of a user ID, $(F(\cdot), p)$, we encode the frequency into a dense representation e^f and convert the discrete position of ID into one-hot vectors v with size $k + 1$. After that, these two transformed representations are concatenated and fed into multilayer perceptrons (MLP) with m hidden layers and activation function \tanh :

$$\begin{aligned}
 h_1 &= \tanh(\mathbf{W}_1[e^f : v]) + b_1 \\
 h_2 &= \tanh(\mathbf{W}_2 h_1) + b_2 \\
 &\dots \\
 h_m &= \tanh(\mathbf{W}_m h_{m-1}) + b_m \\
 \text{logits} &= \text{softmax}(h_m)
 \end{aligned} \tag{4}$$

Subsequently, the hidden state of the last layer denoted as h_m is further processed by the Softmax layer and outputs three probabilities that correspond to the actions of "Ascend", "Unchanged", and "Descend". Specifically, selecting "Ascend" as an action would indicate raising the position of the ID within the

shared embedding candidates set (e.g., from p_i to p_{i+1}). “Unchanged” means to remain in the same position. “Descend” means the Identity Agent decides to make a “rollback” on the position of the given ID (e.g., from p_{i+1} to p_i). Recall that the group of candidate shared IDs is hierarchical, then the motivation for designing the action of ascending or descending movements in this hierarchical group is shown in Fig. 3 The approach adopted for updating the position of an ID within the shared embedding candidates set is as follows: (i) when the cumulative frequency of a particular user ID reaches a specific threshold, the ID is elevated to a higher-frequency representation level; (ii) if an ID’s popularity declines, it will descend to a lower-frequency representation. It is noteworthy that (i) an ID’s position shifting from p_k to p_{k+1} indicates that it will no longer be a shared ID, and instead, it will be assigned its own ID and corresponding embedding, which will be initialized using the last shared embedding e_{sk} . This is based on the intuition that the embedding of shared IDs is generally well-trained and, therefore, possesses sufficient generality to provide suitable initialization for individual IDs. (ii) When certain IDs move from unique IDs to shared IDs, it helps in saving storage and avoiding parameter updating issues.

3.1.3 Reward

The primary objective of the Identity Agent is to enhance the recommendation performance by dynamically assigning either unique IDs or shared IDs to ID features that are present in the data stream. To this end, we propose employing a reward function that can effectively evaluate the model’s performance. Specifically, we define the reward as the prediction loss of the recommendation model utilized in our experiment. Given a pair of user u and item i and their current prediction loss L and their corresponding last T prediction losses $L^u = (L_1^u, L_2^u, \dots, L_T^u)$, $L^i = (L_1^i, L_2^i, \dots, L_T^i)$, where $L_t^{u/i}$ refers to the t^{th} prediction loss of user u / item i , the reward is then defined as the difference between the current loss L and the average of last T prediction losses:

$$R^u = \frac{1}{T} \sum_{t=1}^T L_t^u - L \quad (5)$$

$$R^i = \frac{1}{T} \sum_{t=1}^T L_t^i - L \quad (6)$$

The policy of the Identity Agent can be thought of as maximizing the reduction in the present prediction loss concerning the previous average loss. In this paper, we set $T = 30$, which corresponds to the prediction loss of the last 30 days. Using the past T prediction losses instead of just the last one has several benefits in the proposed approach. It allows the system to consider the historical performance of the model and incentivizes it to maintain a consistent level of accuracy over time. Also, using the past T prediction losses can help to reduce the impact of noise or fluctuations in the loss function that may occur

due to random variation in the input data or other factors. This approach allows the Identity Agent to make more informed decisions that result in lower prediction loss during the process of continuous learning and optimization.

3.2 AutoAssign+

The overall framework of AutoAssign+ includes two steps in the learning process, as shown in Fig. 2:

- **The Forward Step.** In the first step, the Identity Agent receives the corresponding states of each incoming ID and takes action to assign either a unique ID or one of the hierarchical shared IDs to each ID in a batch of user-item interactions. The action value and user-item combined features are further processed by the critic network for calculating the Q-value, which can fine-tune the performance of the Identity Agent. After the ID embeddings are assigned, the user and item IDs are processed using the embedding layer and concatenated, following which they are fed into the inference layer. The predicted output is then compared to the actual output, and the mean square error (MSE) loss function is used to calculate the overall loss.
- **The Backward Step.** We first update the critic network based on the gradients of Q-value and TD error until its convergent. Then the recommendation model parameters are updated based on the computed loss, and the Identity Agent parameters are updated by determining the reward based on the prediction loss. As a result, during continuous training and evaluation, the data-driven Identity Agent is fine-tuned and becomes capable of making more informed decisions.

3.2.1 Critic Network

The standard approach involves using a critic network to estimate the action value generated by the actor network, and then updating the actor network parameters based on this value. However, when using the Identity Agent as the actor network, the challenge lies in designing a suitable structure for the critic network to facilitate parameter updating. Our solution involves a shared-bottom layer in the critic network that simultaneously transforms the user-item features and action information. To achieve this, we first apply an embedding layer and an MLP structure to extract the features and then combine the resulting user-item feature and action information as input to a differentiable action value network that is parameterized by ϕ . This network outputs the estimated Q-value based on the state-action information. Specifically, given the current state s and action a , the Q-value is calculated as follows:

$$Q(s, a) = \mathbb{E}[R + \gamma V(s') \mid s, a] \quad (7)$$

$$= R + \gamma \sum_{s' \in \mathcal{S}} p_{s, a, s'} \cdot \max Q(s', a') \quad (8)$$

where $V(\cdot)$ is the state value function, and s', a' is the state and action value of next step. The hyperparameter γ is the discount rate set as 0.95 for our case. In our case, the action value a is estimated by the Identity Agent, and the next state s' is determined with a probability equal to 1. Therefore, the Q-value function in a multi-critic structure can be calculated as follows:

$$Q(s, a; \phi) = R + \gamma Q(s', a'; \phi) \quad (9)$$

Since the actor-critic network framework often faces the challenge of failing to converge, to address this issue, we incorporate the concept of Deterministic Policy Gradient Algorithms [16] by integrating target networks into the learning framework. The target networks share the exact same structure as the critic networks that have been proposed. We denote it as $Q(s, a; \tilde{\phi})$, which have lagging parameter $\tilde{\phi}$. The structure of critic network is shown in Fig. 4.

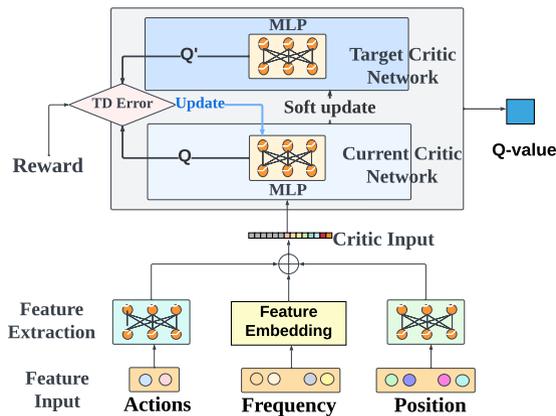


Fig. 4: Structure of Critic Network.

3.3 Overall Optimization

In this section, we will discuss the optimization process of the deep recommendation model, Identity Agent, and critic network, and present the overall optimization framework under the streaming recommendation setting.

3.3.1 Critic Network Updates.

ϕ is the crucial parameter in the critic network, which determines the action Q-value. Given transition. Given the action and state information, the TD target of the target critic network is derived from:

$$TD = R + \gamma Q(s', a'; \tilde{\phi}) \quad (10)$$

where $a' = \pi(s'; \tilde{\theta}_k)$ is the estimated next action from target Identity Agent. The Q-value generated from the current critic network, which estimates the current action value is defined as:

$$Q = Q(s, a; \phi) \quad (11)$$

We calculate the average TD error δ after the training among the batch:

$$\begin{aligned} \delta &= \frac{1}{b} \sum_{k=1}^b (TD_k - Q_k) \\ &= \frac{1}{b} \sum_{k=1}^b [R_k + \gamma Q_k(s', a'; \tilde{\phi}) - Q_k(s, a; \phi)] \end{aligned} \quad (12)$$

Then we update the current critic network for each task by the following gradient decent method with a learning rate α_ϕ

$$\phi \leftarrow \phi - \alpha_\phi \mathbf{I} \delta \nabla_\phi Q(s, a; \phi) \quad (13)$$

where $\nabla_\phi Q(s, a; \phi)$ is the gradient of the target Q-value. This completes the optimization of the current critic networks.

The target critic network is updated until the current critic network reaches the convergence condition towards the direction of parameters in current networks:

$$\tilde{\phi} = \beta \tilde{\phi} + (1 - \beta) \phi \quad (14)$$

where $\beta \in [0, 1]$ is the soft update rate.

3.3.2 Identity Agent Updates.

Let θ represent the parameters in the feature interaction model and ω represent the parameters in the embedding table. In a binary classification task, the Mean Square Error (MSE) can serve as the optimization target. Suppose we have a mini-batch of user-item pairs $\{u_j, i_j\}_{j=1}^N$ and their corresponding labels $\{y_j\}_{j=1}^N$, where $y_j \in \{0, 1\}$ and 0/1 denotes a negative/positive view about a certain item of a user. The MSE loss is defined as follows:

$$MSE(\theta, \omega) = \frac{1}{N} \sum_{j=1}^N (y_j - \hat{y}_j)^2, \quad (15)$$

where the predicted label of the deep recommendation model is denoted as \hat{y}_i .

With the Identity Agent parameterized by Ω and the reward function designed in Eq. (5), the Identity Agent maximizes the expectation of rewards as follows:

$$F(\Omega) = \mathbb{E}_{a \sim \Omega(a|s)}(a | s) \quad (16)$$

where a is the action and s is the state. Before the TD error δ converges to threshold ϵ , we update the Identity Agent parameterized by Ω through the gradients back-propagation of loss function for each layer after the forward process of each batch transitions:

$$\Omega \leftarrow \Omega + \alpha_{\Omega} \mathbf{I} \nabla_{\Omega} \mathcal{J}(\Omega) \quad (17)$$

where the loss for tower layers is defined by the negative of average Q-value $\mathcal{J}(\Omega) = -\frac{1}{b} \sum_{k=1}^b Q'(s, \Omega(a | s))$, which is generated from the critic network.

After the convergence of the critic network, we tend to optimize the Identity Agent based on the overall objective. Practically, it isn't easy to calculate the exact value of expectation reward, thus we use the Monte-Carlo sampling [17] to estimate its gradient along with the optimization algorithm of Eq. (16):

$$\nabla_{\Omega} F(\Omega) = \sum_a R(a | s) \nabla_{\Omega} \Omega(a | s) \quad (18)$$

$$= \sum_a R(a | s) \Omega(a | s) \nabla \log \Omega(a | s) \quad (19)$$

$$= \mathbb{E}_{a \sim \Omega(a|s)} [R(a | s) \nabla \log \Omega(a | s)] \quad (20)$$

$$\approx \frac{1}{N} \sum_{i=1}^N R(a | s) \nabla \log \Omega(a | s) \quad (21)$$

where N is the number of samples. After obtaining the gradient of rewards, the parameter of the Identity Agent network can then be updated with a learning rate α_I by:

$$\Omega \leftarrow \Omega + \alpha_I \nabla_{\Omega} F(\Omega). \quad (22)$$

Next, we present the optimization pipeline of the whole framework. The user-item intersection data are in-flowed boundlessly in the online recommender system in a data stream. We optimize the framework by iteratively collecting mini-batch data with a size of b and updating the parameters of the recommendation model, Identity Agent, and critic network alternately.

As shown in Algorithm 1, in the initialization stage, we initialize Θ , Ω and ϕ from the kaiming initialization [18] (line 1). The state of each ID feature that first time emerges in the data stream is initialized from (1,1). In the Ω updating stage, we sample the last batch in the historical transaction data as the validation batch (line 3) and sample action a_{val} $\Omega(D_{val})$ using the importance sampling [19] (line 4). The state and sample action information is further processed by the critic network to calculate Q-value and corresponding TD error σ , which updates the parameter of the current critic network ϕ (line 5). If the critic network is not converged to ϵ , we update the Identity Agent by the critic network (lines 6-8). Next we calculate the adjusted the position $(\hat{u}_j, \hat{i}_j)_{j=1}^b$ based on sampled action a_{val} (line 9). We then calculate

Algorithm 1 Algorithm for overall optimization

Data: Boundless data stream S in recommender system in the form of (userID, itemID, ground- truth label) with a mini-batch size b : $D = \{(u_j, i_j, y_j)\}_{j=1}^b$; Recommendation model parameter Θ ; Identity Agent parameter Ω ; Critic network parameter ϕ

Result: Fine-tuned parameter of recommendation model, Identity Agent, and critic network Θ^* , Ω^* and ϕ^*

- 1 Initialize Θ , Ω and ϕ from kaiming initialization [18];
- 2 **Repeat**
- 3 Sample a validation batch from the last transaction in the history data stream: $D_{val} = \{(u_j, i_j, y_j)\}_{j=1}^b$
- 4 Sample action a_{val} $\Omega(D_{val})$ by importance sampling;
- 5 Forward action a_{val} and state information into the critic network, calculate the TD error based on Eq. (12). Then updated parameter of the critic network by Eq. (13).
- 6 **if** $\delta \geq \epsilon$ **then**
- 7 | Update the Identity Agent based on Eq. (17)
- 8 **else**
- 9 | Adjust position of $(u_j, i_j)_{j=1}^b$ based on a_{val} to get $(\hat{u}_j, \hat{i}_j)_{j=1}^b$
- 10 | Calculate the reward by Eq. (5) using the adjusted \hat{D}_{val}
- 11 | Updated parameter of Identity Agent Ω by Eq. (22)
- 12 | Given the current batch of user-item data $D_{train} = \{(u_j, i_j, y_j)\}_{j=1}^b$
- 13 | Calculate the action: $a_{train} = \text{argmax}(\Omega(D_{train}))$
- 14 | Adjust position of $(u_j, i_j)_{j=1}^b$ according to a_{train}
- 15 | Updated the parameter of recommendation model Θ using the adjusted data $\hat{D}_{train} = \{(\hat{u}_j, \hat{i}_j, y_j)\}_{j=1}^b$.
- 16 **end**
- 17 Perform soft updates of critic network $\tilde{\phi} \leftarrow \beta\tilde{\phi} + (1 - \beta)\phi$
- 18 **Until** Θ \mathcal{E} Ω converge or encounter the end of the data stream;

the reward by Eq. (5) with the ground-truth label under the evaluation mode of the recommendation model (line 10). Then the parameter of Identity Agent Ω can be updated by Eq. (22) (line 11). Here we introduce the Θ updating stage: First, we collect the current batch of transaction data (line 12) and get the corresponding action (the highest probability among three actions) from the Identity Agent $a_{train} = \text{argmax}(\Omega(D_{train}))$ (line 13). Next we adjust and record the position of $(u_j, i_j)_{j=1}^b$ according to a_{train} (line 14). Then the parameter of recommendation model Θ can be updated using the adjusted data $\hat{D}_{train} = \{(\hat{u}_j, \hat{i}_j, y_j)\}_{j=1}^b$ (line 15). The whole optimization process can be terminated when meeting the satisfied criteria or at the end of the data stream.

4 Experiment

In this section, we will provide a detailed description of our experimental setup. Subsequently, we will conduct extensive experiments on three datasets to investigate four research questions:

- (1) What is the overall performance of our method among different datasets, and whether we have solved the significant cold-start problem?
- (2) How is the generalizability of our method? Is it in line with the actual streaming recommendation in the industry?
- (3) How does each component in our proposed method contributes to the final achievement?
- (4) Is our method able to allocate unnecessary memory usage? If the answer is affirmative, then what is the scale of specificity?

4.1 Experiment Settings

4.1.1 Dataset

We conduct the overall performance experiments on three popular datasets in recommendations.

MovieLens 25M¹: MovieLens 25M is a newly released stable benchmark dataset for personalized movie recommendations. It contains 25 million intersections and one million tag applications applied to 62,000 movies by 162,000 users. The ratings range from 1 to 5.

MovieLens Latest²: Similar to MovieLens 25M, it contains 27,000,000 ratings and 1,100,000 tag applications applied to 58,000 movies by 280,000 users. The major difference between this Dataset and MovieLens 25M is that this dataset is highly sparse as it only has a density of 0.16%, which is suitable for exploring the cold-start problem.

Netflix Price³: Netflix price open competition dataset for predicting user ratings on movies. The movie rating file contains more than 100 million ratings from 480,000 anonymous Netflix users and more than 17,000 movies. Data are collected from October 1998 to December 2005.

In order to simulate the real-world streaming recommendation scenario, we sort the samples uniformly based on their timestamps and use one epoch from the start to the end without multiple iterations. It is important to construct the dataset through timestamps to align with the real-world recommendation scenario, where user interests and preferences are sequential and change over time, avoiding the issue of data traversal. After sorting, we use the first 80% of the samples for parameter training and evaluate the accuracy and loss of the prediction results over the last 20% of the data. In the testing phase, we predict and record the performance alternately and continuously train our model [20], as mentioned in the previous optimization section.

¹<https://grouplens.org/datasets/MovieLens/>

²<https://grouplens.org/datasets/MovieLens/>

³<https://www.kaggle.com/datasets/netflix-inc/netflix-price-data>

Table 1: Overall information of Datasets

Dataset	Total User	Total Item	Ratings	Density
MovieLens Latest	280,000	58,000	27,000,000	0.16%
MovieLens 25M	162,000	62,000	25,000,000	0.25%
Netflix Price	480,000	17,000	100,000,000	1.22%

4.1.2 Evaluation Metrics

Since the focus of our study is on the binary classification task, we have converted the ratings of the three datasets to binary labels where 1 (>3) represents the user likes the item and 0 (≤ 3) represents the user dislikes the item. To evaluate the performance of our model, we have used Mean Square Error (MSE) loss, Accuracy, and Area Under the Curve (AUC) as the evaluation metrics.

4.1.3 Implementation details

In our experiments, we adopt one hidden layer in our agent network with a size of 512, and the embedding dimension of frequency is 32. For the recommendation model, we adopt an embedding size of 128 for both the user and item and two hidden layers with a size of 512 and use *LeakyRelu* as the activation function. Both the agent networks and recommendation networks use Adam optimizer and 0.0001 and 0.001 as the learning rates, respectively. The structure of the critic network is an input embedding layer with dimension 128, a $128 \times 512 \times 256$ Multi-Layer Perceptron (MLP) as the bottom layer, and a $256 \times 128 \times 64 \times 1$ MLP as the tower layer applying Adam optimizer with learning rate 0.0001, the default soft update rate $\beta = 0.2$ and batch size is 500. According to the experimental results, we set the size of the user/item shared embedding candidates group to 1 and 2, respectively. To clearly demonstrate the effectiveness of our method and solely investigate the difference between well-trained and under-trained embeddings, we only utilize the user and item embeddings as the input of our recommendation model in the horizontal comparison of these datasets. The implementation code is available online to ease reproducibility.⁴

4.1.4 Baselines

- **Origin** This model does not employ a threshold to filter out low-frequency ID features. It is used as a reference to demonstrate the negative impact and harm that low-frequency IDs and the cold-start problem can have on model performance.
- **LFF (Grid/Random)** LFF method is used to filter low-frequency IDs by setting a fixed frequency threshold. Those IDs that occur less than this threshold are assigned one shared embedding. The best two thresholds for User and Item are selected through grid search and random search.
- **AutoAssign** Automatic Shared Embedding Assignment method with only one agent network, which is the original version.

⁴<https://github.com/Applied-Machine-Learning-Lab/AutoAssign-Plus>

4.2 Overall Performance

Table 2 presents the overall performance of the evaluated methods. From the results, we can observe that

Table 2: Overall Performance

Dataset	Model	MSE Loss	Accuracy	AUC	Time (min)	Improve (%)
MovieLens 25M	Origin	0.2155	0.6578	0.7105	10	—
	LFF-Grid	0.2079	0.6747	0.7314	110	2.56
	LFF-Random	0.2074	0.6754	0.7322	231	2.67
	AutoAssign	<u>0.2052</u>	<u>0.6792</u>	<u>0.7364</u>	18	3.25
	AutoAssign+	0.2046	0.6817*	0.7386*	39	3.63
MovieLens Latest	Origin	0.2157	0.6567	0.7108	12	—
	LFF-Grid	0.2077	0.6749	0.7315	140	2.77
	LFF-Random	0.2073	0.6756	0.7328	294	2.87
	AutoAssign	<u>0.2050</u>	<u>0.6802</u>	<u>0.7377</u>	20	3.56
	AutoAssign+	0.2040*	0.6821*	0.7422*	41	3.87
Netflix Price	Origin	0.1914	0.7062	0.7640	26	—
	LFF-Grid	0.1819	0.7234	0.7883	280	2.43
	LFF-Random	0.1818	0.7235	0.7883	588	2.45
	AutoAssign	<u>0.1792</u>	<u>0.7293</u>	<u>0.7946</u>	45	3.27
	AutoAssign+	0.1756*	0.7349*	0.8093*	87	4.06

Bold denotes the highest score, and the underline indicates the best result of the baselines.

* represents the significance level p -value < 0.05 comparing with the best baselines.

- Our method AutoAssign+ outperforms the baseline Origin in all evaluation metrics with a significant average relative improvement of 3.85% in Accuracy. This finding supports the effectiveness of the proposed method in improving the predictive ability of the recommendation model by dynamically and automatically assigning shared embedding to low-frequency features using the critic network.
- Baseline LFF-Grid and LFF-Random achieved a remarkable result compared to the baseline Origin (roughly 2.6% relative improvement in average Accuracy). It demonstrates that the intuition of using shared ID embeddings to represent those low-frequency IDs is simple yet effective.
- Our proposed method outperforms all the baselines. The LFF-Grid and LFF-Random methods require an extensive searching process and yield sub-optimal results due to several limitations. (1) They use stationary thresholds that cannot capture the complex and dynamic distribution of user and item IDs in the streaming environment. (2) These methods use only frequency as the standard of criterion and have only one shared embedding, ignoring the alteration of the distribution of ID features in space and time. In contrast, the AutoAssign+ method has a group of hierarchical candidate shared embeddings that capture the fine-grained information of features with different frequencies. Temporally, AutoAssign+ has a fallback mechanism that

captures the pattern of users who have had no interactive data for a long time but already has used a unique ID and reassigns the shared embedding to them.

- AutoAssign+ outperforms AutoAssign on both three datasets. The Identity Agent of AutoAssign generates the candidate embeddings based solely on the input features and doesn't take into account the performance of the recommendation model. By incorporating the critic network into the optimization process, AutoAssign+ can generate more accurate and effective embeddings for the low-frequency IDs, leading to better performance of the recommendation model compared to AutoAssign.
- The train time of grid search and random search is time-consuming compared to the Origin method, while the AutoAssign+ requires only 3.4 times of the Origin method, which is acceptable in a real industrial scenario.

4.3 Cold-Start Stage

This experiment compares the performance of our proposed AutoAssign+ method with other baselines in terms of the cold-start stage in recommendation systems. As shown in Fig. 5, we choose the MovieLens Latest dataset, and the average accuracy of different user and movie frequencies was plotted, varying from 0 to 1,000. The results showed that our proposed AutoAssign+ method consistently outperformed other baselines in terms of accuracy. Specifically, in the user part, the gap between our method and other baselines was more significant, with an average accuracy improvement of 1.1% compared to LFF. In the movie part, the baseline LFF suffered a dramatic downward trend in the beginning stage due to the low number of movies, while AutoAssign+ reached a stable and much higher performance. The reason for this is that a fixed shared embedding threshold makes all movies use one ID embedding, causing a severe deviation. However, AutoAssign+ can dynamically make different judgments for each movie by minimizing the different losses that result from using the shared ID and smoothly avoiding the deviation caused by the single frequency information, even when using frequency information similarly.

4.4 Ablation Study

An ablation experiment was conducted on MovieLens Latest dataset to demonstrate how each component of AutoAssign+ contributed to the overall performance. The method "No Descend" involved reducing the output of the agent network to only two actions: "Ascend" and "Unchanged," while "Single SE" meant setting the number of shared embeddings of both userID and itemID to only 1. The results in Table 3 showed that both AutoAssign (No Descend) and AutoAssign (Single SE) produced a minor gap from AutoAssign but still achieved a superior performance to the baseline LFF. However, AutoAssign (Single SE + No Descend) performed worse than LFF. The reward of the RL-based Identity agent was based on the last five prediction losses of each ID, and the behavior of different IDs dynamically changed in both spatial

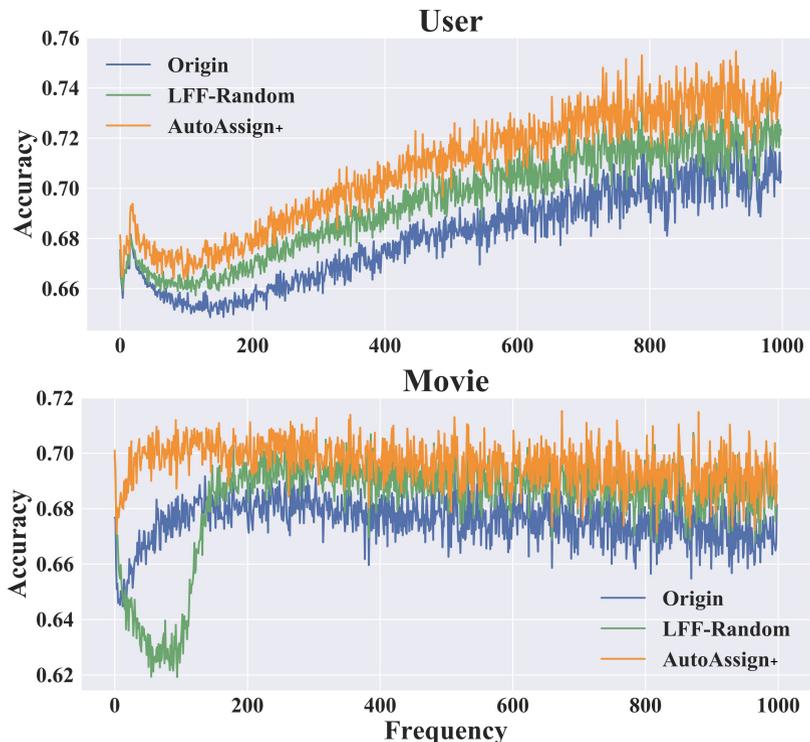


Fig. 5: Accuracy over different frequencies.

and temporal manners. In such a case, the use of both uni-directional action and a single shared embedding would make the agent’s decision excessively judgmental, generating the same deviation as LFF. Therefore, the hierarchical group of shared embedding and descend action had to be used simultaneously to achieve a better result. The critic network in AutoAssign+ further improved the model performance by optimizing the Identity Agent’s policy, enhancing the stability and efficiency of the framework.

Table 3: Effectiveness of Different Components

Method	MSE loss	Accuracy	Improve
Origin	0.2155	0.6578	—
AutoAssign(No Descend)	0.2059	0.6785	3.15%
AutoAssign(Single SE)	0.2056	0.6793	3.26%
AutoAssign(Single SE + No Descend)	0.2082	0.6732	2.34%
AutoAssign	0.2052	0.6801	3.40%
AutoAssign+	0.2047	0.6819	3.66%

4.5 Practicality Comparison

To ensure fairness in the evaluation of the low-frequency filter, we apply grid search and random search methods to optimize its performance. As shown in Fig. 6, the frequency thresholds for the user and movie were selected from a range of 5 to 200 with an interval of 10 in the grid search, and the same threshold was applied to all three datasets. We observed that the performance of the low-frequency filter decreased monotonically when the frequency threshold was set to 20 or higher. To further improve the effectiveness of the low-frequency filter, we conducted a random search for different threshold combinations on the user and movie sides in the interval of 5-20. The threshold group with the highest Accuracy was selected as the experimental group for overall performance, as presented in Fig. 7. Also, note that in Table 2, the train time of AutoAssign is much faster than grid search and random search, proving our model’s efficiency in real-world recommendations.

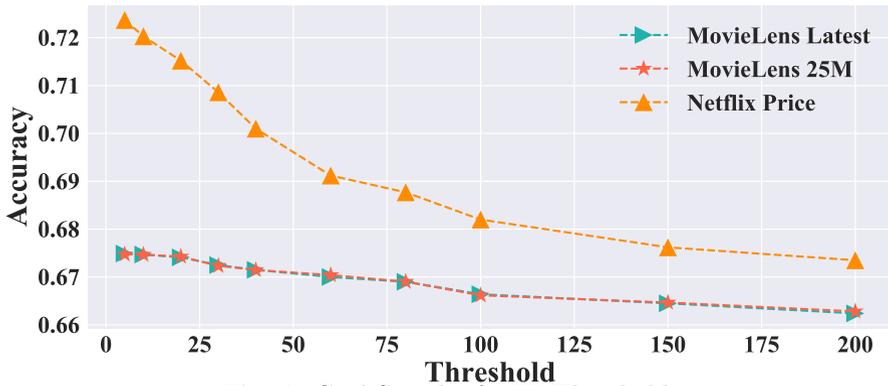


Fig. 6: Grid Search of LFF Threshold.

4.6 Parameter Reduction

Reducing memory consumption is a crucial factor in designing efficient recommender systems, especially for online platforms. The embedding layer is one of the most significant contributors to the total number of parameters in deep recommender models. The proposed Automated Shared Embedding approach in AutoAssign+ aims to reduce memory consumption by assigning shared embeddings to low-frequency IDs, thus avoiding unnecessary allocation of unique embeddings to these IDs. Table 4 presents the parameter usage of the embedding layer in AutoAssign+ and the deduction ratio achieved by using the shared embedding approach. The *Deduction* column shows the percentage of parameters saved in different datasets, and the *Total* column shows the overall deduction ratio. The results demonstrate that the shared embedding approach reduces the total number of parameters in the embedding layer by 20% to 30%, depending on the dataset. At the same time, the model’s accuracy is improved

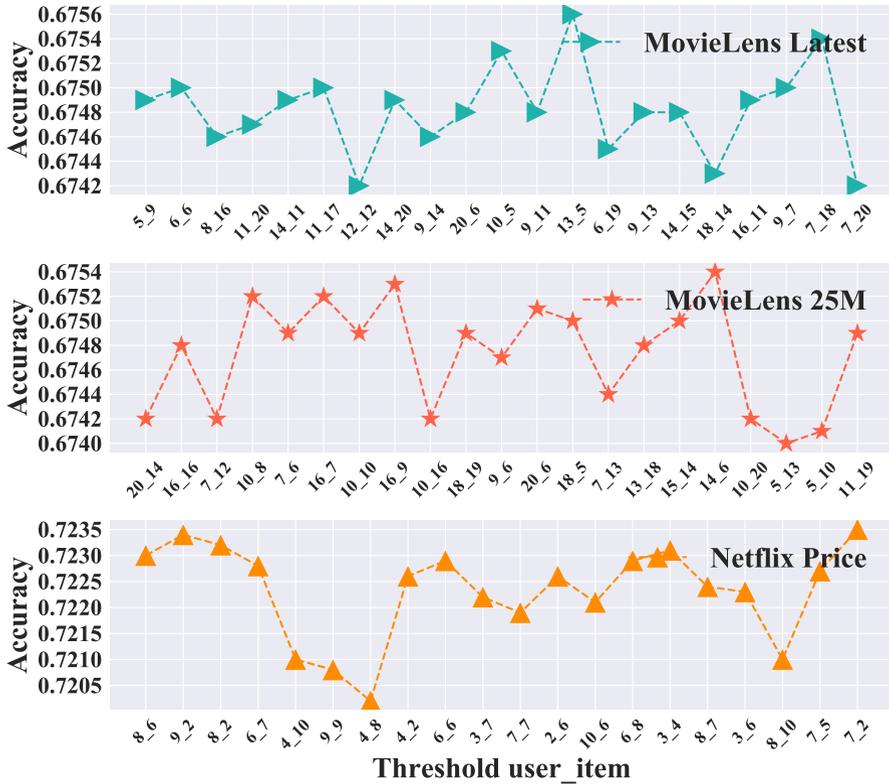


Fig. 7: Random Search of LFF Threshold.

due to the enhanced representative ability of the shared embeddings. Therefore, the Automated Shared Embedding approach in *AutoAssign+* can reduce memory consumption without compromising the model’s accuracy.

Table 4: Parameter Deduction on *AutoAssign+*

Dataset(field)	Origin Para	<i>AutoAssign+</i> Para	Deduction	Total
ML-Latest(user)	35,840,000	25,553,937	28.7%	31.3%
ML-Latest(movie)	7,424,000	4,164,872	43.9%	
ML-25M(user)	20,736,000	16,153,349	22.1%	23.1%
ML-25M(movie)	7,936,000	5,833,016	26.5%	

4.7 Case Study

In this subsection, we present a comparison between our proposed *AutoAssign+* method and the LFF baseline to illustrate the effectiveness of our approach. On the left-hand side instance, we consider a new user from the MovieLens-25M dataset, whose frequency increases to 12 at timestamp $t + 2$, which is below the threshold of 20. In this case, LFF assigns a shared ID to the user. However, at timestamp $t + 4$, the user frequency increases to 33,

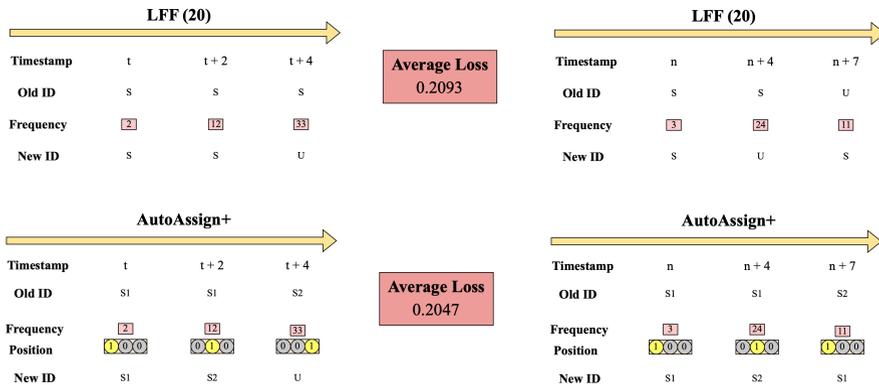


Fig. 8: Case Study on MovieLens-25M.

which is above the threshold, leading LFF to assign a specific user ID to this user. In contrast, AutoAssign+ adopts a different ID assignment strategy. At timestamp $t + 2$, when the user frequency increases to 12 with a position shift from p_1 to p_2 , AutoAssign+ assigns a shared ID with a higher hierarchy to the user. This hierarchical candidate shared ID setting with reinforcement learning powered embedding assignment results in a reduction in the average loss from 0.2093 to 0.2047, which highlights the efficiency of AutoAssign+.

5 Related Work

5.1 Cold-Start problem

The purpose of the recommender system is to recommend a set of items that the user may be interested in. However, if the user interaction data is limited, the performance of the recommendation model will be significantly reduced, which is called the cold-start problem. Cold-start problem is a ubiquitous and challenging problem in personalized recommendation, where extensive studies have been done concerning this issue [21–25]. Some content-based methods [6] make relevance between zero shot-learning and cold-start recommendation and propose a low-rank linear auto-encoder to solve the cold-start problem using the user’s auxiliary information. Internal and contextual attention networks (ICAN) [8] strengthen the interaction of the feature domain and use auxiliary information among multiple queues to get a better cold-start performance in the matching stage. Meta-learning is a common method to learn the global and general information for pertinent tasks and serves to speed up the initialization of new relative tasks. MetaEmb [9] and MeLU [26] apply Model-Agnostic Meta-Learning, where MetaEmb trains a generator to initialize embedding, and MeLU learns the initialization parameters of the whole model to solve the cold-start problem. MWUF [10] uses the average pooling of all items to initialize the new item embedding and uses two meta-network to enhance their representation. These methods mentioned above have their own limitations and are discussed in the introduction section.

5.2 Reinforcement Learning Based Recommendation

Numerous studies have explored the integration of Reinforcement Learning (RL) with Recommender Systems (RS) [27–34]. Instead of optimizing immediate user feedback similar to traditional learning-to-rank approaches [35], RL-based Recommender Systems seeks to optimize the cumulative reward function, which estimates multiple rounds of interactions between the recommendation policy and the user response environment. Specifically, the problem of sequential user-system interactions can be formulated as a Markov Decision Process (MDP) [36]. Various RL solutions have been investigated under this formulation, including tabular-based methods that store and update a table that represents the estimated value or quality of each state-action pair [37–39], value-based methods utilized to assess the effectiveness of a specific action or state [40–43], policy gradient methods that optimize the recommendation policy based on long-term reward [44, 45], and actor-critic methods [46] that simultaneously learn an action evaluator network and action generator, which is based on policy gradient [47–50]. Some of the major difficulties in employing RL for recommender systems include the vast state and action space [51, 52], uncertainty in the user environment [53, 54], exploration effectiveness and efficiency [55], and creating suitable reward function that caters to diverse behaviors [56]. Our work addresses reward function design while simultaneously enhancing the agent’s performance by utilizing the actor-critic framework.

6 Conclusion

In this study, we analyze the cold-start problem in streaming deep recommender systems, which results in poorly trained ID embeddings and reduces prediction performance, as well as leads to unnecessary memory usage in the model. To address these issues, we propose a framework called Automatic Shared Embedding Assignment Plus (AutoAssign+), which includes a critic network-enhanced Identity Agent that can automatically and field-wisely assign shared IDs to low-frequency IDs. AutoAssign+ reduces the human effort required in the time-consuming search process and expert knowledge and takes action to improve performance and reduce the model’s parameter usage by 20%-30%. The proposed framework is independent of the inference layer, making it easily applicable to various existing recommendation models that use an embedding lookup layer. We demonstrate the effectiveness of AutoAssign+ in addressing the cold-start problem and the practicality of their approach through extensive experiments on three popular datasets. However, the selection of actions (ascend/unchanged/descend) made by the actor network is somewhat opaque and not easily interpretable, which makes it difficult to gain a clear understanding of how the system is making decisions. This is one of the directions for further improvement.

ACKNOWLEDGEMENT

This research was partially supported by Huawei (Huawei Innovation Research Program), APRC - CityU New Research Initiatives (No.9610565, Start-up Grant for New Faculty of City University of Hong Kong), CityU - HKIDS Early Career Research Grant (No.9360163), Hong Kong ITC Innovation and Technology Fund Midstream Research Programme for Universities Project (No.ITS/034/22MS), SIRG - CityU Strategic Interdisciplinary Research Grant (No.7020046, No.7020074). We thank MindSpore [57] for the partial support of this work, which is a new deep learning computing framework.

References

- [1] An, M., Wu, F., Wu, C., Zhang, K., Liu, Z., Xie, X.: Neural news recommendation with long-and short-term user representations. In: ACL, pp. 336–345 (2019)
- [2] Ren, K., Zhang, W., Rong, Y., Zhang, H., Yu, Y., Wang, J.: User response learning for directly optimizing campaign performance in display advertising. In: CIKM, pp. 679–688 (2016)
- [3] Guo, L., Yin, H., Wang, Q., Chen, T., Zhou, A., Quoc Viet Hung, N.: Streaming session-based recommendation. In: KDD, pp. 1569–1577 (2019)
- [4] He, B., He, X., Zhang, Y., Tang, R., Ma, C.: Dynamically expandable graph convolution for streaming recommendation. In: Proceedings of the ACM Web Conference 2023. WWW '23, pp. 1457–1467. Association for Computing Machinery, New York, NY, USA (2023). <https://doi.org/10.1145/3543507.3583237>. <https://doi.org/10.1145/3543507.3583237>
- [5] Guo, H., Chen, B., Tang, R., Zhang, W., Li, Z., He, X.: An embedding learning framework for numerical features in ctr prediction. In: KDD (2021)
- [6] Li, J., Jing, M., Lu, K., Zhu, L., Yang, Y., Huang, Z.: From zero-shot learning to cold-start recommendation. In: AAAI, vol. 33, pp. 4189–4196 (2019)
- [7] Mo, K., Liu, B., Xiao, L., Li, Y., Jiang, J.: Image feature learning for cold start problem in display advertising. In: IJCAI (2015)
- [8] Xie, R., Qiu, Z., Rao, J., Liu, Y., Zhang, B., Lin, L.: Internal and contextual attention network for cold-start multi-channel matching in recommendation. In: IJCAI, pp. 2732–2738 (2020)

- [9] Pan, F., Li, S., Ao, X., Tang, P., He, Q.: Warm up cold-start advertisements: Improving ctr predictions via learning to learn id embeddings. In: SIGIR, pp. 695–704 (2019)
- [10] Zhu, Y., Xie, R., Zhuang, F., Ge, K., Sun, Y., Zhang, X., Lin, L., Cao, J.: Learning to warm up cold item embeddings for cold-start recommendation with meta scaling and shifting networks. In: SIGIR, pp. 1167–1176 (2021)
- [11] Liu, Q., Tian, F., Zheng, Q., Wang, Q.: Disentangling interest and conformity for eliminating popularity bias in session-based recommendation. *Knowl. Inf. Syst.* **65**(6), 2645–2664 (2023). <https://doi.org/10.1007/s10115-023-01839-0>
- [12] He, M., Wang, J., Ding, T., Shen, T.: Conversation and recommendation: knowledge-enhanced personalized dialog system. *Knowl. Inf. Syst.* **65**(1), 261–279 (2023). <https://doi.org/10.1007/s10115-022-01766-6>
- [13] Zhang, W., Qin, J., Guo, W., Tang, R., He, X.: Deep learning for click-through rate estimation. arXiv preprint arXiv:2104.10584 (2021)
- [14] Du, H., Tang, Y., Cheng, Z.: An efficient joint framework for interacting knowledge graph and item recommendation. *Knowl. Inf. Syst.* **65**(4), 1685–1712 (2023). <https://doi.org/10.1007/s10115-022-01808-z>
- [15] Chen, J., Zheng, L., Chen, S.: User view dynamic graph-driven sequential recommendation. *Knowl. Inf. Syst.* **65**(6), 2541–2569 (2023). <https://doi.org/10.1007/s10115-023-01840-7>
- [16] Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., Riedmiller, M.: Deterministic policy gradient algorithms. In: International Conference on Machine Learning, pp. 387–395 (2014). Pmlr
- [17] Hastings, W.K.: Monte carlo sampling methods using markov chains and their applications (1970)
- [18] He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: ICCV, pp. 1026–1034 (2015)
- [19] Goertzel, G.: Quota sampling and importance functions in stochastic solution of particle problems. Technical report (1949)
- [20] Liu, H., Zhao, X., Wang, C., Liu, X., Tang, J.: Automated embedding size search in deep recommender systems. In: SIGIR, pp. 2307–2316 (2020)
- [21] Feng, J., Xia, Z., Feng, X., Peng, J.: Rbpr: A hybrid model for the new user cold start problem in recommender systems. *Knowledge-Based Systems*

- 214**, 106732 (2021). <https://doi.org/10.1016/j.knosys.2020.106732>
- [22] Wahab, O.A., Rjoub, G., Bentahar, J., Cohen, R.: Federated against the cold: A trust-based federated learning approach to counter the cold start problem in recommendation systems. *Information Sciences* **601**, 189–206 (2022)
- [23] Jeevamol, J., Renumol, V.: An ontology-based hybrid e-learning content recommender system for alleviating the cold-start problem. *Education and Information Technologies* **26**, 4993–5022 (2021)
- [24] Liu, H., Wang, L., Li, P., Qian, C., Zhao, P., Wu, X.: Relation-propagation meta-learning on an explicit preference graph for cold-start recommendation. *Knowledge-Based Systems* **272**, 110579 (2023). <https://doi.org/10.1016/j.knosys.2023.110579>
- [25] Wang, X., Peng, Z., Wang, S., Yu, P.S., Fu, W., Xu, X., Hong, X.: Cdlfm: cross-domain recommendation for cold-start users via latent feature mapping. *Knowledge and Information Systems* **62**, 1723–1750 (2020)
- [26] Lee, H., Im, J., Jang, S., Cho, H., Chung, S.: Melu: Meta-learned user preference estimator for cold-start recommendation. In: *KDD*, pp. 1073–1082 (2019)
- [27] Liu, Z., Tian, J., Cai, Q., Zhao, X., Gao, J., Liu, S., Chen, D., He, T., Zheng, D., Jiang, P., *et al.*: Multi-task recommendations with reinforcement learning. In: *Proceedings of the ACM Web Conference 2023*, pp. 1273–1282 (2023)
- [28] Afsar, M.M., Crump, T., Far, B.: Reinforcement learning based recommender systems: A survey. *ACM Computing Surveys (CSUR)* (2021)
- [29] Wang, Y., Sharma, M., Xu, C., Badam, S., Sun, Q., Richardson, L., Chung, L., Chi, E.H., Chen, M.: Surrogate for long-term user experience in recommender systems. In: *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 4100–4109 (2022)
- [30] Zhang, Q., Liu, J., Dai, Y., Qi, Y., Yuan, Y., Zheng, K., Huang, F., Tan, X.: Multi-task fusion via reinforcement learning for long-term user satisfaction in recommender systems. In: *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 4510–4520 (2022)
- [31] Zhao, X., Zhang, L., Ding, Z., Yin, D., Zhao, Y., Tang, J.: Deep reinforcement learning for list-wise recommendations. *arXiv preprint arXiv:1801.00209* (2017)

- [32] Zhao, X., Xia, L., Zhang, L., Ding, Z., Yin, D., Tang, J.: Deep reinforcement learning for page-wise recommendations. In: RecSys, pp. 95–103 (2018). ACM
- [33] Zhao, X., Gu, C., Zhang, H., Yang, X., Liu, X., Liu, H., Tang, J.: Dear: Deep reinforcement learning for online advertising impression in recommender systems. In: AAAI, vol. 35, pp. 750–758 (2021)
- [34] Liu, H., Cai, K., Li, P., Qian, C., Zhao, P., Wu, X.: Redrl: A review-enhanced deep reinforcement learning model for interactive recommendation. *Expert Systems with Applications*, 118926 (2022)
- [35] Liu, T.-Y., *et al.*: Learning to rank for information retrieval. *Foundations and Trends® in Information Retrieval* **3**(3), 225–331 (2009)
- [36] Shani, G., Heckerman, D., Brafman, R.I., Boutilier, C.: An mdp-based recommender system. *Journal of Machine Learning Research* **6**(9) (2005)
- [37] Joachims, T., Freitag, D., Mitchell, T., *et al.*: Webwatcher: A tour guide for the world wide web. In: IJCAI (1), pp. 770–777 (1997). Citeseer
- [38] Mahmood, T., Ricci, F.: Learning and adaptivity in interactive recommender systems. In: Proceedings of the Ninth International Conference on Electronic Commerce, pp. 75–84 (2007)
- [39] Moling, O., Baltrunas, L., Ricci, F.: Optimal radio channel recommendations with explicit and implicit feedback. In: Proceedings of the Sixth ACM Conference on Recommender Systems, pp. 75–82 (2012)
- [40] Taghipour, N., Kardan, A., Ghidary, S.S.: Usage-based web recommendations: a reinforcement learning approach. In: Proceedings of the 2007 ACM Conference on Recommender Systems, pp. 113–120 (2007)
- [41] Zheng, G., Zhang, F., Zheng, Z., Xiang, Y., Yuan, N.J., Xie, X., Li, Z.: Drn: A deep reinforcement learning framework for news recommendation. In: Proceedings of the 2018 World Wide Web Conference, pp. 167–176 (2018)
- [42] Zhao, X., Zhang, L., Ding, Z., Xia, L., Tang, J., Yin, D.: Recommendations with negative feedback via pairwise deep reinforcement learning. In: KDD, pp. 1040–1048 (2018). ACM
- [43] Ie, E., Jain, V., Wang, J., Narvekar, S., Agarwal, R., Wu, R., Cheng, H.-T., Chandra, T., Boutilier, C.: Slateq: A tractable decomposition for reinforcement learning with recommendation sets (2019)
- [44] Chen, M., Beutel, A., Covington, P., Jain, S., Belletti, F., Chi, E.H.: Top-k

- off-policy correction for a reinforce recommender system. In: Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining, pp. 456–464 (2019)
- [45] Chen, H., Dai, X., Cai, H., Zhang, W., Wang, X., Tang, R., Zhang, Y., Yu, Y.: Large-scale interactive recommendation with tree-structured policy gradient. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 33, pp. 3312–3320 (2019)
- [46] Liu, F., Tang, R., Li, X., Zhang, W., Ye, Y., Chen, H., Guo, H., Zhang, Y.: Deep reinforcement learning based recommendation with explicit user-item interactions modeling. arXiv preprint arXiv:1810.12027 (2018)
- [47] Peters, J., Schaal, S.: Natural actor-critic. *Neurocomputing* **71**(7-9), 1180–1190 (2008)
- [48] Bhatnagar, S., Ghavamzadeh, M., Lee, M., Sutton, R.S.: Incremental natural actor-critic algorithms. *Advances in neural information processing systems* **20** (2007)
- [49] Degris, T., Pilarski, P.M., Sutton, R.S.: Model-free reinforcement learning with continuous action in practice. In: 2012 American Control Conference (ACC), pp. 2177–2182 (2012). IEEE
- [50] Liu, L., Cai, L., Zhang, C., Zhao, X., Gao, J., Wang, W., Lv, Y., Fan, W., Wang, Y., He, M., Liu, Z., Li, Q.: Linrec: Linear attention mechanism for long-term sequential recommender systems. In: Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval. SIGIR '23, pp. 289–299. Association for Computing Machinery, New York, NY, USA (2023). <https://doi.org/10.1145/3539618.3591717>. <https://doi.org/10.1145/3539618.3591717>
- [51] Dulac-Arnold, G., Evans, R., van Hasselt, H., Sunehag, P., Lillicrap, T., Hunt, J., Mann, T., Weber, T., Degris, T., Coppin, B.: Deep reinforcement learning in large discrete action spaces. arXiv preprint arXiv:1512.07679 (2015)
- [52] Liu, F., Tang, R., Li, X., Zhang, W., Ye, Y., Chen, H., Guo, H., Zhang, Y., He, X.: State representation modeling for deep reinforcement learning based recommendation. *Knowledge-Based Systems* **205**, 106170 (2020)
- [53] Ie, E., Hsu, C.-w., Mladenov, M., Jain, V., Narvekar, S., Wang, J., Wu, R., Boutilier, C.: Recsim: A configurable simulation platform for recommender systems (2019) [arXiv:1909.04847](https://arxiv.org/abs/1909.04847) [cs.LG]
- [54] Zhao, X., Xia, L., Tang, J., Yin, D.: Deep reinforcement learning for search, recommendation, and online advertising: a survey. *ACM SIGWEB*

Newsletter (Spring), 1–15 (2019)

- [55] Chen, M., Chang, B., Xu, C., Chi, E.H.: User response models to improve a reinforce recommender system. In: Proceedings of the 14th ACM International Conference on Web Search and Data Mining, pp. 121–129 (2021)
- [56] Zou, L., Xia, L., Ding, Z., Song, J., Liu, W., Yin, D.: Reinforcement learning to optimize long-term user engagement in recommender systems. In: Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pp. 2810–2818 (2019)
- [57] MindSpore (2020). <https://www.mindspore.cn/>