

Multi-Task Fusion via Reinforcement Learning for Long-Term User Satisfaction in Recommender Systems

Qihua Zhang
Tencent Inc.
Shenzhen, China
kirahzhang@tencent.com

Junning Liu
Tencent Inc.
Beijing, China
korchinliu@tencent.com

Yuzhuo Dai
Tencent Inc.
Shenzhen, China
zoyadai@tencent.com

Yiyan Qi
Tencent Inc.
Shenzhen, China
yiyanqi@tencent.com

Yifan Yuan
Tencent Inc.
Shenzhen, China
yaphetyuan@tencent.com

Kunlun Zheng
Tencent Inc.
Beijing, China
kunlunzheng@tencent.com

Fan Huang*
Tencent Inc.
Shenzhen, China
sinohuang@tencent.com

Xianfeng Tan
Tencent Inc.
Beijing, China
victan@tencent.com

ABSTRACT

Recommender System (RS) is an important online application that affects billions of users every day. The mainstream RS ranking framework is composed of two parts: a Multi-Task Learning model (MTL) that predicts various user feedback, i.e., clicks, likes, sharings, and a Multi-Task Fusion model (MTF) that combines the multi-task outputs into one final ranking score with respect to user satisfaction. There has not been much research on the fusion model while it has great impact on the final recommendation as the last crucial process of the ranking. To optimize long-term user satisfaction rather than obtain instant returns greedily, we formulate MTF task as Markov Decision Process (MDP) within a recommendation session and propose a Batch Reinforcement Learning (RL) based Multi-Task Fusion framework (BatchRL-MTF) that includes a Batch RL framework and an online exploration. The former exploits Batch RL to learn an optimal recommendation policy from the fixed batch data offline for long-term user satisfaction, while the latter explores potential high-value actions online to break through the local optimal dilemma. With a comprehensive investigation on user behaviors, we model the user satisfaction reward with subtle heuristics from two aspects of user stickiness and user activeness. Finally, we conduct extensive experiments on a billion-sample level real-world dataset to show the effectiveness of our model. We propose a conservative offline policy estimator (Conservative-OPEstimator) to test our model offline. Furthermore, we take online experiments in a real recommendation environment to compare performance of different models. As one of few Batch RL researches applied in MTF task successfully, our model has also been deployed on a large-scale industrial short video platform, serving hundreds of millions of users.

KEYWORDS

Recommender system; Batch reinforcement learning; Multi-task fusion; Long-term user satisfaction

1 INTRODUCTION

With the information explosion on the Internet, Recommender Systems (RS) that aim to recommend potentially interesting items for users, are playing an increasing important role in various platforms including E-commerce sites [11, 19, 37, 38], videos sharing sites [4], social networks [13, 16], etc. There are usually two main stages in industrial RS: candidate generation and ranking [35]. The first stage selects hundreds or thousands of candidates from millions or even billions of items, while the second stage returns several top-ranked items for each user from the candidates.

Given a user query, the quality of the ranking results is a leading factor in affecting user satisfaction. Early works [23, 25, 36, 37] usually only consider a single instant metric, e.g., Click-Through Rate (CTR), and train a ranking model over this metric. In practice, it is usually hard to measure real user satisfaction with just one metric. For example, in news recommendation scenario, users may click the recommended news with click bait like the cover pictures or titles but quickly exit when they are not interested in the content. Clearly, only considering CTR will cause a lot of improper recommendations and this encourages us to explore multiple metrics through different user behaviors. When more than one metric is considered, an essential question arises, as to how should these metrics be combined to optimize ranking quality.

Model fusion is a popular approach to solve this problem [3, 6], which is usually composed of two parts: 1) a Multi-Task Learning model (MTL) [20, 29] that predicts multiple metrics associated with user satisfaction; 2) a Multi-Task Fusion model (MTF) [15, 24] that constructs a combination function based on those predicted scores and produces the final ranking. Compared with the MTL works on the first step, MTF algorithms are crucial for recommendation quality but there has been little good research on it.

In this paper, we focus on the MTF algorithms in RS. Given a fusion function $f(o_1, o_2, \dots, o_k)$ with respect to ranking scores o_1, o_2, \dots, o_k from different predicted tasks, a naive way to find out the optimal fusion weights is to use hyper-parameter searching techniques such as Grid Search and Bayesian Optimization [22].

*Corresponding author

These methods fail in large RS not only because of their inefficiency but also because they can't produce personalized weights for different users and different contexts to make accurate recommendations. One can solve these issues by bridging user states and fusion weights through neural networks and turn to optimize the network weight via Evolutionary Strategy [2]. However, all the above methods still focus on optimizing instant user satisfaction in a greedy way but ignore long-term user satisfaction. To reduce the expected regret of recommendations and improve the long-term rewards for RSs, we need to consider instant user satisfaction as well as delayed user satisfaction brought by the long-term utility of the recommendation. After all, the current recommendation may affect user preference later.

Intuitively, Reinforcement Learning (RL) is usually designed to maximize long-term rewards. However, applying RL in the large-scale online RS to optimize long-term user satisfaction is still a non-trivial problem. 1) The long-term user satisfaction is very complicated and can be measured in various user behaviors. How to build feasible reward according to different behaviors is challenging. 2) In order to learn an optimal RL policy effectively, the recommender agent requires a very large number of sequential interactions with real users to trial and error. However, on-policy RL would harm user experiences with the noise generated by online exploration during learning. 3) An alternative is to build a recommender agent offline through the logged data, which can mitigate the cost of the trial-and-error exploration. Unfortunately, traditional off-policy methods not only suffer from the Deadly Triad problem [30], i.e., the problem of instability and divergence arises whenever combining function approximation, bootstrapping, and offline training, but also suffer from serious extrapolation error[9] where state-action pairs not in the fixed batch data, also called out-of-distribution (OOD) training data, are erroneously estimated to have unrealistic values.

Considering the above problems, we propose a Batch RL based Multi-Task Fusion framework (BatchRL-MTF) to optimize long-term user satisfaction. Our model consists of two components: 1) Batch RL framework learns an optimal recommendation policy with high returns, strong robustness and less extrapolation error offline, which provides the fusion function with a set of personalized fusion weights trading off instant and long-term user satisfaction; 2) Online Exploration Policy interacts with real users online to discover high-reward fusion weights as offline batch samples. In particular, we formulate our MTF task as Markov Decision Process (MDP) within a recommendation session to model the sequential interaction between the user and the RS. We first comprehensively investigate different user behaviors, and deliberately design the reward based on these behaviors to optimize long-term user satisfaction from two aspects of user stickiness and user activeness. To reduce the learning costs of our model and the damage to user experience, we present Batch RL to enhance learning efficiency, mitigate extrapolation error and optimize accumulated rewards according to the history logs. In addition, we conduct online exploration that can find more high-value state-action pairs to ensure the sufficient training of our model, which not only reduces extrapolation error but also prevents our model from falling into local optimum. In experiments, we propose a conservative offline policy estimator (Conservative-OPEstimator) to verify the outstanding performance

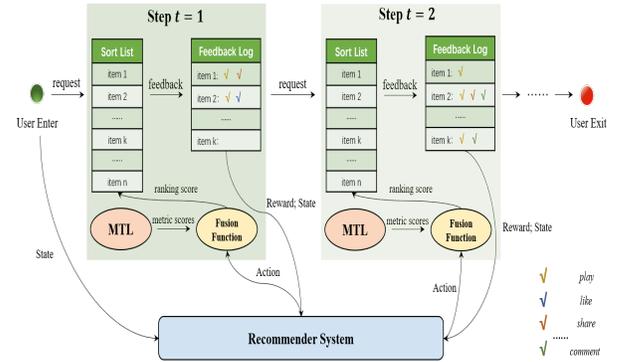


Figure 1: The User-RS interaction within a session-based recommendation.

of the model by evaluating the long-term rewards brought by it. Meanwhile, we conduct online evaluations for BatchRL-MTF with competitive baselines on a real-world short video platform. The significant improvements on user stickiness and user activeness exhibit the effectiveness of our BatchRL-MTF in practice. The major contributions of our work include:

- We formulate the session-based MTF task as an MDP and exploit Batch RL to optimize long-term user satisfaction, which is one of the few works of Batch RL applied in MTF task successfully.
- We design our reward function with respect to multiple user behaviors that relate to user stickiness and user activeness, and train our BatchRL-MTF based on history logs. Specially, our BatchRL-MTF contains two components: Batch RL framework and online exploration policy. The former learns an optimal recommendation policy offline and the latter explores potential high-value state-action pairs online. We show that our framework can mitigate the deadly triad problem and extrapolation error problem of traditional off-policy applied in practice RSs.
- We creatively propose a conservative offline policy estimator (Conservative-OPEstimator) to test our model offline, while conducting online experiments in real recommendation environment to demonstrate our model outperforms baselines greatly. In addition, BatchRL-MTF has been deployed in the short video recommendation platform and remarkably improved 2.550% app dwell time and 9.651% user positive-interaction rate.

2 PROBLEM FORMULATION

As shown in Figure 1, in the short video recommendation scenario, the RS agent interacts with a user $u \in U$ at discrete time steps within a recommendation session. At each time step t , the RS feeds a top-ranked item $i^{(t)} \in I$ (or item list $(i_1^{(t)}, \dots, i_l^{(t)})$) to user u and receives feedback vector $v^{(t)} = (v_1^{(t)}, \dots, v_m^{(t)})$, where I is the candidate set and $v_i^{(t)}, 1 \leq i \leq m$ is a specific user's behavior (e.g., video play time, like, sharing, etc.) on item $i^{(t)}$. The mainstream RS ranking pipeline is composed of two parts: an MTL model that predicts various user behaviors and an MTF model that combines the multi-task outputs, i.e., $o = (o_1, \dots, o_k)$, into one final ranking score. We employ Progressive Layered Extraction (PLE) model [29] with

excellent performance to conduct multi-task predictions. To output the personalized recommendation matching user preferences, we build a fusion function $f(o)$ that aggregates these predicted scores to model user satisfaction. Considering the effect of magnitude difference among o_1, \dots, o_k , we define the function as follow:

$$f(o|\alpha) = \sum_i^k \alpha_i \log(o_i + \beta_i), \quad (1)$$

where $\alpha = (\alpha_1, \dots, \alpha_k)$ is the fusion weights to be optimized and $\beta = (\beta_1, \dots, \beta_k)$ is a constant bias set by priori knowledge to smooth those o_1, \dots, o_k . In this paper, we aim to find out the optimal weights α that maximize long-term user satisfaction.

In particular, we study the above problem as Markov Decision Process (MDP) within the recommendation session, where an agent (the RS) interacts with the environment (users) by sequentially recommending items over time, to maximize the cumulative reward within the session. Formally, the MDP consists of a tuple of five elements $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$:

- **State space \mathcal{S} :** is a set of user state s , which is composed of two parts: user profile feature (e.g., age, male, location, etc.) and user's interaction history feature (e.g., like, sharing, skipping, etc.). In our RS, the latter is usually formulated by a user's various interactions with his/her watched last 500 videos.

- **Action space \mathcal{A} :** The action determines the ranking of the candidate list. In our problem, the action a is the fusion weight vector $\alpha = (\alpha_1, \dots, \alpha_k)$ to be optimized in Equation (1). With the fusion weights, the agent can further calculate the final ranking scores and return recommendation item(s).

- **Reward \mathcal{R} :** After the RS takes an action a_t at the state s_t , i.e., figure out the fusion scores and return a recommendation item to response user's request at time step t , a user will feedback with different behaviors. We define the instant reward $r(s_t, a_t)$ based on these behaviors.

- **Transition probability \mathcal{P} :** Transition probability $p(s_{t+1}|s_t, a_t)$ defines the state transition from s_t to s_{t+1} after taking action a_t . As the user state is defined as user profile and his/her interaction history, the next state s_{t+1} is determined after user feedback.

- **Discount factor γ :** Discount factor $\gamma \in [0, 1]$ is introduced to control the present value of future reward.

With the above notations and definitions, the task of Batch RL applied in MTF can be defined as: given the interaction history logs between RS agent and users in MDP form within the recommendation session, how to find a policy obtaining the optimal fusion weights α to maximize the cumulative reward of RS. In the rest of this paper, we omit the superscript (t) for simplicity.

3 PROPOSED FRAMEWORK

In this section, we propose a Batch RL framework for short video RSs to optimize long-term user satisfaction. Specifically, we design the reward based on multiple user behaviors and present our Batch RL architecture. Firstly, we discuss how to train our BatchRL-MTF offline via users' behavior log. Secondly, we also conduct online exploration to discover potential high-value state-action pairs in real recommendation environment. More importantly, we briefly describe the implementation of the proposed model on a real-world short-video recommender platform.

3.1 Reward Function

In a recommendation session, the RS interacts with users in turn, i.e., the agent takes an action a at the state s (i.e., recommending a video to a user with respect to MTL ranking scores) and then the user provides his/her feedback v . To comprehensively measure the instant satisfaction, we formulate the reward function as follow:

$$r(s, a) = \sum_{i=1}^m w_i v_i, \quad (2)$$

where w_i is the weight of feedback v_i . In our RS, feedback v_1, \dots, v_m includes video play time, play integrity and interaction behavior such as liking, sharing, commenting, liking, etc. And weights w_1, \dots, w_m are set via extensive statistic analysis on relationships between these feedback and future user app dwell time.

3.2 Batch Reinforcement Learning for MTF

Due to costly real-world interactions, it is difficult to apply on-policy and off-policy algorithms in large-scale industrial RSs. An alternative approach is learning an RL-based recommender agent solely from historical logs, which is also known as Batch Reinforcement Learning. However, a severe problem in Batch RL approaches is extrapolation error, which is caused by the distribution difference between training data and learned policy. To solve this problem, we exploit Batch-Constrained deep Q-learning model (BCQ) [9], based on the popular Actor-Critic architecture, to learn the optimal fusion weights in our framework.

3.2.1 Actor Network. The Actor network includes two sub-network: (a) the action generative network $G_\theta = \{E_{\theta_1}, D_{\theta_2}\}$ and (b) the action perturbation network $P_\omega(s, a, \rho)$.

Specially, network $G_\theta = \{E_{\theta_1}, D_{\theta_2}\}$ is a conditional variational auto-encoder (VAE) to generate candidate action set of which the distribution is similar to that of training samples \mathcal{B} , which reduces extrapolation error and enhances our framework robustness. G_θ can be further divided into an encoder block $E_{\theta_1}(z|s, a)$ and a decoder block $D_{\theta_2}(a|s, z)$. To be more specific, encoder $E_{\theta_1}(z|s, a)$ learns the latent distribution of $(s, a) \in \mathcal{B}$ and forces it to approximate $\mathcal{N}(0, 1)$ via KL divergence. $E_{\theta_1}(z|s, a)$ has two outputs, i.e., the mean value μ and standard deviation δ . Decoder $D_{\theta_2}(a|s, z)$ takes the latent vector $z \sim \mathcal{N}(\mu, \delta^2)$ and the user state s as inputs to output the action \hat{a} that is similar to $(s, a) \in \mathcal{B}$. Formally, we train the VAE based on the following objective on the log-likelihood of the dataset:

$$\theta \leftarrow \underset{\theta}{\operatorname{argmin}} \sum_{(s,a) \in \mathcal{B}} (a - G_\theta(s))^2 + D_{KL}(E_{\theta_1}(z|s, a)|\mathcal{N}(0, 1)). \quad (3)$$

To increase the diversity of the actions from VAE, we use network $P_\omega(s, a, \rho)$ to generate a perturbation $\xi \in [-\rho, \rho]$ with respect to state s and action a . Specifically, given user state s , we generate n actions $\{\hat{a}_i \sim G_\theta(s)\}_{i=1}^n$ as candidates. Meanwhile, $P_\omega(s, a, \rho)$ generates n perturbations $\{\xi_i\}_{i=1}^n$, $\xi_i = P_\omega(s, \hat{a}_i, \rho)$ that updates the actions as $\hat{a}_i + \xi_i$. As shown in Figure 2, the optimal action is selected from n perturbed actions as:

$$\pi(s) = \underset{\hat{a}_i + \xi_i}{\operatorname{argmax}} Q(s, \hat{a}_i + \xi_i), \quad (4)$$

Algorithm 1: Offline Training of BatchRL-MTF.

Input : the transition dataset \mathcal{B} , mini-batch size M , number of sampled actions n , perturbation bound ρ , target network update rate η_t , discount factor γ , delay update step L , number of training epochs EP

- 1 Initialize generative network $G_\theta = \{E_{\theta_1}, D_{\theta_2}\}$, perturbation network P_ω , and Critic network $Q_\phi = \{Q_{\phi_1}, Q_{\phi_2}\}$ with random parameters θ, ω, ϕ ;
- 2 Initialize target networks $P_{\omega'}, Q_{\phi'} = \{Q_{\phi'_1}, Q_{\phi'_2}\}$ with $\omega' \leftarrow \omega, \phi'_1 \leftarrow \phi_1, \phi'_2 \leftarrow \phi_2$;
- 3 **foreach** $1 \leq ep \leq EP$ **do**
- 4 Sample mini-batch of M transitions (s, a, r, s') from \mathcal{B} randomly;
- 5 Update G_θ according to Equation (3);
- 6 Sample n actions $\{\hat{a}'_i \sim G_\theta(s')\}_{i=1}^n$;
- 7 Generate n perturbed actions $\{\hat{a}'_i + P_\omega(s', \hat{a}'_i, \rho)\}_{i=1}^n$;
- 8 Update P_ω according to Equation (5);
- 9 Update Q_ϕ according to Equation (6);
- 10 **if** $ep \% L == 0$ **then**
- 11 $\omega' \leftarrow \eta_t \omega + (1 - \eta_t) \omega'$;
- 12 $\phi'_i \leftarrow \eta_t \phi_i + (1 - \eta_t) \phi'_i, i = 1, 2$;
- 13 **end**
- 14 **end**

and $P_\omega(s, a, \rho)$ is optimized through the deterministic policy gradient algorithm [28]:

$$\omega \leftarrow \underset{\omega}{\operatorname{argmax}} \sum_{\substack{s \in \mathcal{B} \\ \hat{a} \sim G_\theta(s)}} Q(s, \hat{a} + P_\omega(s, \hat{a}, \rho)). \quad (5)$$

3.2.2 Critic Network. The Critic network $Q_\phi(s, a)$ aims to estimate the cumulative reward of a state-action pair (s, a) . Following the common setting, we build four Critic networks during the learning process, i.e., two current networks Q_{ϕ_1}, Q_{ϕ_2} and two target networks $Q_{\phi'_1}, Q_{\phi'_2}$. Specifically, the goal of Critic network is to minimize TD-error in the bootstrapping way:

$$\phi_j \leftarrow \underset{\phi_j}{\operatorname{argmin}} \sum_{(s, a, s') \in \mathcal{B}} [y - Q_{\phi_j}(s, a)]^2, \quad j \in \{1, 2\}, \quad (6)$$

and the learning target y is set according to the Clipped Double Q-Learning [8] that reduces the overestimation bias:

$$y = r + \gamma \max_{a'} [\min_{j=1,2} Q_{\phi'_j}(s', a')], \quad (7)$$

$$a' \in \{\hat{a}'_i + P_{\omega'}(s', \hat{a}'_i, \rho), \hat{a}'_i \sim G_\theta(s')\}_{i=1}^n$$

where a' sampled from the generative model and outputted by the target action perturbation network.

3.2.3 Offline Model Training. We train our Batch RL model offline based on the pre-collected dataset and the learning algorithm is shown in Algorithm 1. We first construct a transition dataset \mathcal{B} based on history trajectories that record the online interaction between the user and the recommender agent. To improve the utilization efficiency of samples and speed up the convergence of our model, we used replay buffer to store historical transitions

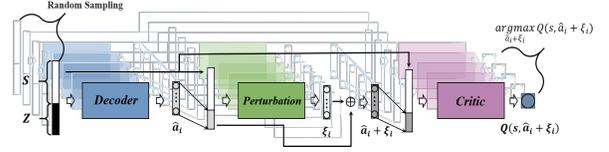


Figure 2: The BatchRL-MTF policy.

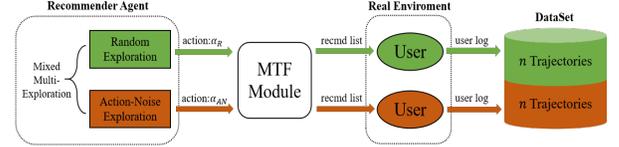


Figure 3: The online exploration policy of BatchRL-MTF framework.

during training process. Specially, replay buffer is a fixed-size queue, where new transitions will replace old ones. Based on the mini-batch sampled from \mathcal{B} , we update sub-networks G_θ, P_ω , and Q_ϕ in order. In addition, we perform soft update on target networks every L -iteration to reduce overestimation bias.

3.3 Online Exploration

Although the proposed Batch RL framework can reduce extrapolation error, we have considered the Actor network learns a fine-tuned policy based on training data generated by behavior policy. As a result, the agent's performance is limited by behavior policy. If behavior policy doesn't explore the real environment enough, it will cause extrapolation error and local optimum problems in target policy. To solve this problem, we propose an online exploration policy to discover potential high-reward state-action pairs in online serving. In particular, we first perform two types of exploration on two groups of users separately.

- **Random Exploration.** The agent randomly samples an action from Gaussian distribution to interact with real users while having no priori knowledge.

- **Action-Noise Exploration.** To improve exploration efficiency and exploit the priori knowledge of optimal actions, we conduct action-noise exploration policy, which adds Gaussian noise to the action outputted by the optimal target policy. Formally, we have:

$$\pi_{ep}(s) = \pi_t^*(s) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, 0.1). \quad (8)$$

where $\pi_{ep}(s)$ is exploration policy and $\pi_t^*(s)$ is the target policy of current model. The random exploration policy can enrich the diversity of actions, which prevents the local optimum problem; the action-noise exploration policy makes full use of the priori knowledge of optimal actions to explore nearby actions with potentially high value. To exploit the advantages of these two exploration policy, we propose an online exploration policy that combines them, called *Mixed Multi-Exploration policy*.

- **Mixed Multi-Exploration** is composed of random exploration and action-noise exploration, shown in Figure 3. It constructs the training dataset with equal amounts of trajectory samples collected by random exploration and action-noise exploration.

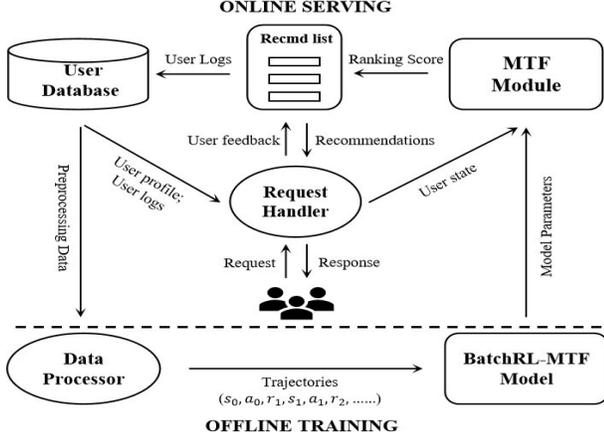


Figure 4: BatchRL-MTF framework of our short video recommender system.

Different from exploration with simulators, online exploration will receive feedback from real users, which results in more accurate and unbiased rewards. In experiments, we will show that training model with the dataset from our mixed multi-exploration policy will exhibit better performance.

3.4 Short Video Recommender System with BatchRL-MTF

We implement our model on a large-scale short video recommendation platform. As shown in Figure 4, the RS is composed of two components: offline training and online serving. The two components are connected through user database and online MTF module, which are used for collecting user-agent interaction logs and conducting MTF task based on the BatchRL-MTF model.

The offline training component is responsible for data preprocessing and BatchRL-MTF model training. To be more specific, data processor pulls recent user data including user profile and user history logs from user database and organises these data as interaction trajectories $\{\tau_i = (s_i^{(0)}, a_i^{(0)}, r_i^{(1)}, s_i^{(1)}, a_i^{(1)}, r_i^{(2)}, \dots)\}_{i=1}^N$, where each trajectory represents a user-agent interaction session. Then, we exploit Algorithm 1 to learn a new model based on these trajectories and update online model. In practice, we re-train our model daily based on the trajectory data collected in the past three days.

The online serving component provides personalized recommendations and collects user data. When receiving a user request, request handler retrieves the relevant information of this user from user database, constructs user state feature, and transfers user state to MTF module. According to the user state, the MTF module produces a final ranking score for each video candidate. Finally, request handler feeds top-ranked recommended video(s) to the user and collect his/her feedback.

4 EXPERIMENTS

4.1 Dataset

Our batch data is collected from a real-world short video recommendation platform, including about 3.142 million sessions and

11.155 million user-agent interactions. For offline experiments, to ensure that models fully learn the sequential information between recommendation sessions, we take the first 90% user sessions of the dataset in time order as the training dataset \mathcal{B} to train the model and the remaining 10% user sessions as testing dataset \mathcal{D} to evaluate the model. For online experiments, we exploit the batch dataset to learn an optimal model and deploy it in our online short video recommendation platform for one month to conduct A/B tests.

4.2 Implementation Details

In BatchRL-MTF, the input user state s is concatenated by user profiles feature and interaction feature from last 500 watched videos. Its output action a is 12-dimensional vector representing fusion weights α in Equation (1). All networks in BatchRL-MTF are MLP with ReLU activation function in hidden layers and are optimized based on Adam optimizer. For action perturbation network P_ω , we use Tanh activation function maps its output ψ to $[-1, 1]$ and let the perturbation bound $\rho = 0.15$. We set the reward discount factor as $\gamma = 0.95$. The initial learning rate η for the action generative network, the action perturbation network and critic networks is set to 0.1×10^{-2} , 0.1×10^{-3} and 0.2×10^{-3} separately; the soft update rate and the delay update step for target networks are $\eta_t = 0.05$ and $L = 10$ separately. The above parameters are determined by offline experiments for maximizing long-term returns. In addition, the replay buffer size, mini-batch size and training epochs in our training process are set to 100,000, $M = 256$, and $Ep = 300,000$ respectively.

4.3 Evaluation Setting

4.3.1 Offline policy estimator. Online A/B testing, a general method in the industry for evaluating new recommendation technologies, is also one of our preferred methods in comparison experiments. But it takes time and costs resources. Most importantly, the terrible policy could hurt user experiences. To overcome the above problems and accelerate iteration of new technologies, we propose an offline policy estimator to offline evaluate the performance of RL model. Inspired by Fitted Q Evaluation (FQE) algorithm [18, 31] and Conservative Q-Learning (CQL) algorithm [17], our Conservative Offline Policy Estimator, also called Conservative-OPEstimator, is designed as:

$$\hat{V}(\pi_e) = \frac{1}{n} \sum_{i=1}^n \sum_{\substack{s_i^0 \sim d^0(s) \\ a \sim \pi_e(a|s)}} \pi_e(a|s_i^0) \hat{Q}(s_i^0, a, \theta) \quad (9)$$

where $d^0(s)$ is the initial state distribution; $\pi_e(a|s)$ is the policy to be evaluated; $\hat{Q}(\cdot, \theta)$ is used to estimate how many present values of long-term revenues within an online real recommendation session will be produced by an initial state-action pair. To improve the accuracy of estimation, we resort to CQL algorithm to construct $\hat{Q}(\cdot, \theta)$ by function approximation, which punishes the estimated Q values of state-action pairs not in the dataset \mathcal{D} to prevent overestimation

of the policy value. $\hat{Q}(\cdot, \theta) = \lim_{k \rightarrow \infty} \hat{Q}_k$ where:

$$\begin{aligned} \hat{Q}_{k+1} &\leftarrow \underset{\theta}{\operatorname{argmax}} \alpha \cdot \mathcal{R}(\theta) + \frac{1}{2} \cdot \mathbb{E}_{s,a,r,s' \sim \mathcal{D}} \mathcal{T}^2(\theta) \\ \mathcal{R}(\theta) &= \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_e(a|s)} \hat{Q}_k(s, a, \theta) - \mathbb{E}_{s,a \sim \mathcal{D}} \hat{Q}_k(s, a, \theta) \\ \mathcal{T}(\theta) &= \hat{Q}_k(s, a, \theta) - [r + \gamma \mathbb{E}_{s' \sim \mathcal{D}, a' \sim \pi_e(a'|s')} \hat{Q}_k(s', a', \theta)] \end{aligned} \quad (10)$$

where $\mathcal{R}(\theta)$, called CQL regularizer, is used to penalize Q values of OOD data and $\mathcal{T}(\theta)$ is the standard Temporal-Difference (TD) error. \mathcal{D} is the testing dataset collected from our online recommendation platform, as mentioned in Section 4.1. α is the penalty factor that regulates the punishment degree to the CQL regularizer. To reduce the estimated error, we set α as 0.5×10^{-3} to learn \hat{Q}_k such that the expected value of a policy $\pi_e(a|s)$ under this Q-function lower-bounds its true value. We describe the implementation details of Conservative-OPEstimator in Appendix A.

In our experiments, we take the output value of Conservative-OPEstimator in Equation (9) as the offline evaluation metric. Specially, we use **Long-term user satisfaction per session** $\hat{V}(\pi_e)$ to measure the average present value of long-term user satisfaction generated by a recommendation policy during a session.

4.3.2 Online A/B testing. We deploy baselines to a large short video platform over a month of online A/B test and focus on long-term user satisfaction with the recommendation policies. We select two online metrics as follows to comprehensively test the model from two aspects of user stickiness and user activeness:

- **App dwell time (ADTime)** is the average APP usage time for all users within a day.
- **User positive-interaction rate (UPIRate)** is the percentage of video plays with positive user interactions during a day.

4.4 Compared Methods

4.4.1 Baselines. We compare our model with the common non-reinforcement learning algorithms and the advanced reinforcement learning algorithms.

- **Bayesian Optimization (BO)** fits the prior distribution of the objective function based on Gaussian Process Regression and samples the optimal weight with upper confidence bound.
- **Evolutionary Strategy (ES)** uses a simple network whose inputs are user state and outputs are personalized fusion weights. With the network, ES turns to search the optimal network parameters through natural gradient.
- **Twin Delayed Deep Deterministic Policy Gradient (TD3)** [5] is an advanced off-policy algorithm, which also learns two target networks to reduce the overestimation bias.
- **UWAC+TD3** connects TD3 with Uncertainty Weighted Actor-Critic(UWAC) [33] based on the well-established Bayesian uncertainty estimation methods, so that it can identify OOD training samples and reduce their weights to train a conservative Q function.
- **CQL+SAC** [17] combines CQL with Soft Actor-Critic (SAC) [14]. By regularizing the Q value of OOD action-state pairs, CQL algorithm learns a conservative, lower-bound Q function to reduce extrapolation error. SAC maximizes the cumulative rewards as well as the entropy of the policy to increase the stochastic of the policy and break away from the local optimum.

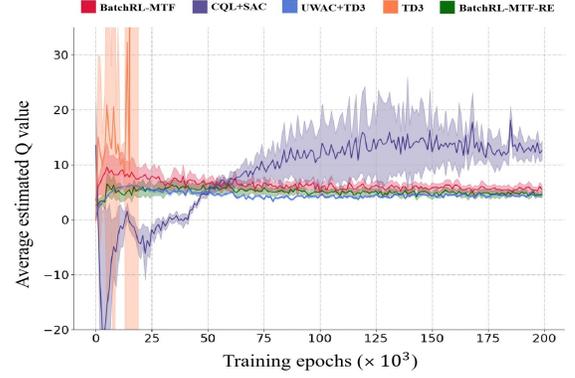


Figure 5: The estimated Q value curve of different policy.

Table 1: Offline and Online Experimental results. The online experiment results represents the improvement of other methods on ADTime and UPIRate compared with Bayesian Optimization (BO). ‘*’ means as the benchmark, BO results is 0. The best result is bold. ‘-’ means that the result is empty.

Methods	Metrics	online	
		ADTime	UPIRate
BO	-	*	*
ES	-	+0.376%	+0.402%
TD3	-648.162	-	-
UWAC+TD3	-297.053	-0.513%	+16.061%
CQL+SAC	5.194	+2.322%	+10.258%
BatchRL-MTF	4.126	+2.216%	+9.118%
BatchRL-MTF-RE	3.023	+0.862%	-1.282%
BatchRL-MTF-Rtime	-	+2.254%	+8.877%
BatchRL-MTF-Rintegrity	-	+1.996%	+9.464%
BatchRL-MTF-Rinteraction	-	+2.550%	+9.651%

4.4.2 Variations of our model. We also compare our model with two types of variants designed to illustrate the effects of online exploration and reward function on model performance, respectively.

- **BatchRL-MTF-RE** is trained only with the data from random exploration. We use this variant to illustrate the power of our mixed multi-exploration.
- **BatchRL-MTF-Rtime** aims to improve user stickiness by setting a larger affinity weight on video play time in Equation (2).
- **BatchRL-MTF-Rintegrity** aims to improve user stickiness by setting a larger affinity weight on video play integrity in Equation (2).
- **BatchRL-MTF-Rinteraction** aims to improve user activeness by setting larger affinity weights on interaction behaviors in Equation (2).

4.5 Offline Evaluation

In this section, we conduct extensive offline experiments to verify the excellent performance of our Batch RL model with strong robustness, high returns and less extrapolation error, while demonstrating

our online exploration that can discover high-value training samples to prevent policy from falling into local optimum. In addition, we test the sensitivity of parameters in our model and details of the analysis results are presented in Appendix B.

4.5.1 Effectiveness of Our Batch RL Framework. As mentioned, the traditional off-policy algorithms exhibit large error in the estimation of Q value when we use the fixed dataset for offline learning. In Figure 5, we show RL algorithms of the average estimated Q value and its variance during training, and further show the policy action distribution in Figure 6 to observe extrapolation error of RL models. All models are trained by samples from mixed multi-exploration policy and all action values are truncated within a pre-defined interval and then normalized to $[-1, 1]$.

Extrapolation error is the overestimation bias to Q value of OOD data, and it is constantly accumulated by iterating the Bellman backup operator. We notice that TD3 tends to output extreme action values that are out of the distribution of batch data in Figure 6 because its Q function overestimates Q values of OOD actions. The estimated Q value of TD3 in Figure 5 exponentially increases in the early stage and can not converge to a stable and reasonable value. In order to alleviate extrapolation error of TD3, we also try to optimize TD3 with UWAC algorithm. In Figure 6, different from TD3 which only generates OOD actions, UWAC+TD3 exploits Monte Carlo (MC) dropout to identify OOD actions and avoids learning from these actions. As a result, the model has lower probability to produce OOD actions. Unfortunately, UWAC+TD3 still have serious extrapolation error problem.

Both BatchRL-MTF and CQL+SAC perform well on reducing extrapolation error. However, these two methods apply different strategies during training process. CQL+SAC model chooses to penalize the Q values of unseen state and action pairs. In complex RS with noisy user feedback, this soft constraint strategy may not get rid of all OOD actions. Different from CQL+SAC, BCQ model introduces the action generative network that hard constraints the output actions around the seen ones and exploits online exploration to increase the diversity of actions. In Figure 5, we find that, although BCQ produces lower Q value, it is more stable than CQL+SAC and gives faster convergence during training, which is benefited from the direct constraints on actions. In next section, we also show that BCQ achieves stable improvements on both ADTime and UPIRate on our online RS.

In addition, we take our Conservative-OPEstimator to evaluate all models and the result are shown in Table 1. Conservative-OPEstimator gives a negative evaluation value $\hat{V}(\pi_e)$ to TD3 and UWAC + TD3 both with large extrapolation error, and thinks these models would hurt user experience. The above results also confirm our policy takes effect that Conservative-OPEstimator tries to punish the estimated Q values of OOD actions to avoid overestimation. Although CQL+SAC with a highest evaluation value, $\hat{V}(\pi_e)$, as an average metrics, can only represents the overall rewards. Compared to BCQ, which directly limits the distribution of actions, CQL+SAC has a more volatile output, which is fatal to large-scale personalized recommendation platforms that demand online stability. In a word, the overall performance of our model is optimal in considering of both stability and returns.

4.5.2 Effectiveness of Online Exploration. We also compare the output action distribution of our methods with different online exploration policies, i.e., BatchRL-MTF and BatchRL-MTF-RE, in Figure 7. Because Batch RL model based on behavior cloning constraints itself not to output OOD actions, the action distribution of both Batchr-MTF and Batchr-MTF-RE is more concentrated than that of their fixed batch data. However, the policy searches for the optimal action in the batch dataset, which limits the improvement of policy and causes the local optimal problem.

Therefore, our online exploration combines the extensive exploring of random exploration with the priori knowledge of action-noise exploration to construct high-quality dataset. It is obvious that the action distribution of BatchRL-MTF is more concentrated than that of behavior policy and BatchRL-MTF-RE. Compared with BatchRL-MTF-RE, BatchRL-MTF explores potential high-value state-action pairs based on the priori knowledge, which reduces those unnecessary exploration. The evaluation results of Conservative-OPEstimator in Table 1 also verify the effectiveness of our online exploration policy that focuses on searching for high-value samples while ensures the diversity of actions. Later, we will further show the efficiency of our online exploration via an online comparison experiment.

4.6 Online Evaluation

In this section, we deploy all compared models and our BatchRL-MTF on an online RS, except TD3 (it doesn't work and could seriously hurt user experience). We not only compare our model with other baselines, but also discuss the weights w_i in reward function on their effect of our model.

4.6.1 Comparison against Other Methods. We use BO as the benchmark for comparison and show the improvements of ADTime and UPIRate for other methods and the experimental results are shown in Table 1, where all improvements have statistical significance with p -value <0.05 . ES slightly improves the performance by considering user preferences, but is still a performance gap to RL models that aim to optimize long-term user satisfaction. For the evaluation of RL models, the online evaluation results are completely consistent with the offline evaluation results given by our Conservative-OPEstimator. UWAC+TD3 performs poorly, dragging down ADTime by 0.513%; the ADTime and UPIRate returns of CQL+SAC are both the highest in baselines, but not as stable as that of our model. The specific reasons for their performance have been analyzed in Section 4.5.1 and we will not repeat it.

Compared with BatchRL-MTF-RE, the mixed multi-exploration of our model significantly outperforms Gaussian Noise random exploration. To be more specific, random exploration increases 0.862% ADTime but decreases 1.282% UPIRate. The reason may lay in the large amount of noisy exploration on the whole action space that damaged the user experience. On the contrary, our BatchRL-MTF explore around optimal actions and effectively learns the optimal policy from the high-value samples.

4.6.2 Affinity Weight in Reward Function. To explore the impact of different affinity weights w_i in reward function, we deploy variants described in Section 4.4.2 on our RS and compare real user feedback. As shown in Table 1, compared with BatchRL-MTF, BatchRL-MTF-Rtime only gives a slight improvement on user stickiness,

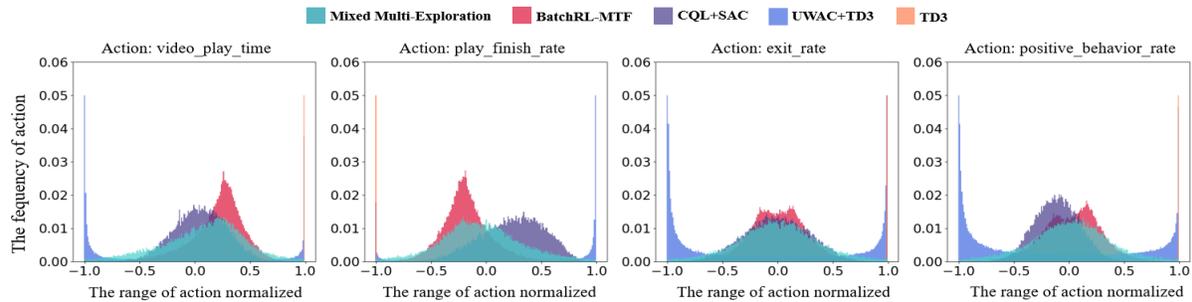


Figure 6: Action distribution of our model and other RL model of baselines. We select the most representative 4-dimensional actions video play time, play finish rate, exit rate and positive behavior rate for comparison.

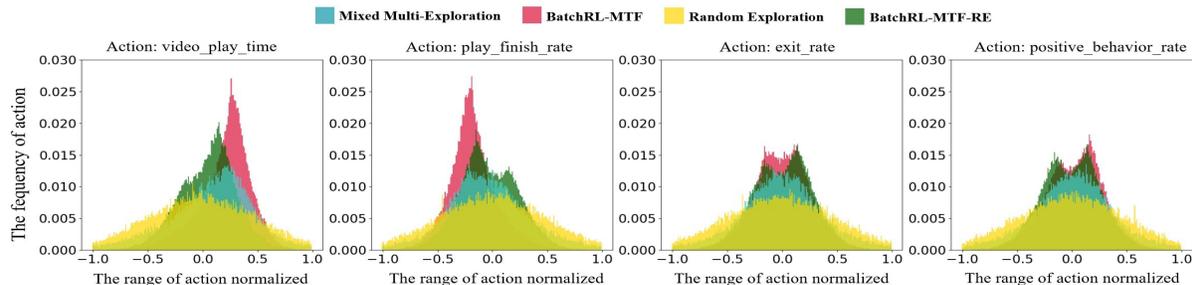


Figure 7: Action distribution of our model with different online exploration. We select the most representative 4-dimensional actions video play time, play finish rate, exit rate and positive behavior rate for comparison.

i.e., increasing ADTime by 0.037%; while harms user activeness significantly, i.e., decreasing UPIRate by 0.241%. Similarly, BatchRL-MTF-Rintegrity also exhibits a seesaw phenomenon between user stickiness and activeness. Clearly, the weights increasing of video play time or integrity may misleading the agent. For example, the agent will recommend more long videos to increase video play time or recommend more short videos to improve video play integrity. We notice that BatchRL-MTF-Rinteraction performs best among these variants, which improves ADTime and UPIRate by 0.333% and 0.534% respectively. The probable reason may be that user interaction behaviors are sparse but strong signals which can guide the agent to learn user preferences and improve user satisfaction.

5 RELATED WORK

Although there have been extensive studies on reinforcement learning based recommender systems [1, 7, 32, 34], we note that their problems are different from those studied in this paper. Particularly, the output (action) of the above methods is usually recommendation item(s), while in this paper we aim to figure out the optimal weights for model fusion of MTL predictions.

To figure out the optimal fusion weights, early works try to solve this problem via parameter searching algorithms such as Grid Search, Genetic Algorithm [21], and Bayesian Optimization [22]. For example, [12, 27] use Grid Search to iterate over all combinations of candidate parameter sets and select the optimal weights through A/B Test. Galuzzi et al. [10] propose to use Bayesian Optimization to optimize the number of latent factors, the regularization parameter,

and the learning rate. The main drawback of these methods is that they always produce unified fusion weights across different users and thus can not model user preferences. To search personalized fusion weights, Ribeiro et al. [26] propose using evolutionary algorithm to find Pareto-efficient hybrids. However, all above methods focus on instant returns but ignore long-term user satisfaction.

Recently, to maximize long-term user satisfaction, a few works try to search the optimal weights via reinforcement learning. Pei et al. [24] propose a reinforcement learning based model to maximize platform expected profit. To simplify the model, they use evolutionary strategy to solve the problem and thus the proposed method are still limited to optimize the profile of current recommendation. Han et al. [15] exploit off-policy reinforcement learning to find out the optimal weights between the predicted click-through rate and bid price of an advertiser. Because numerous interactions with the immature agents will harm user experiences, they build an environment simulator to generate users' feedback for training their model offline. However, the real recommendation environment is so complex that the simulator can't simulate it completely. In fact, the RL model based on the simulator will be difficult to adapt to the online environment and hurt user experience.

6 CONCLUSION

Multi-Task Fusion (MTF), which determines the final recommendation, is a crucial task in large RSs but has not received much attention by far in both academia and industry. In this paper, we propose BatchRL-MTF framework for MTF recommendation task,

which aims to optimize long-term user satisfaction. We first formulate our recommendation task as an MDP and deliberately design the reward function involving multiple user behaviors that represent user stickiness and user activeness. In order not to hurt online user experience, we exploit Batch RL model to optimize accumulated rewards and online exploration policy to discover potential valuable actions. Finally, we propose a Conservative-OPEstimator to test our model offline, while conduct online experiments in real recommendation environment for comparison of different models. Experiments show that our model has the advantages of high returns, strong robustness and small extrapolation error. In addition, we also explore the effect of affinity weight in the reward function on our model, and find that when the weight of user activeness feedback is increased, our model will obtain higher returns. Furthermore, we successfully implement our model in a large-scale short video platform, improving 2.550% app dwell time and 9.651% user positive-interaction rate.

REFERENCES

- [1] M Mehdi Afsar, Trafford Crump, and Behrouz Far. 2021. Reinforcement learning based recommender systems: A survey. *arXiv preprint arXiv:2101.06286* (2021).
- [2] Hans-Georg Beyer and Hans-Paul Schwefel. 2002. Evolution strategies—a comprehensive introduction. *Natural computing* 1, 1 (2002), 3–52.
- [3] David Carmel, Elad Haramaty, Arnon Lazerson, and Liane Lewin-Eytan. 2020. Multi-Objective Ranking Optimization for Product Search Using Stochastic Label Aggregation. In *WWW 2020*. 373–383.
- [4] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep neural networks for youtube recommendations. In *RecSys*. 191–198.
- [5] Stephen Dankwa and Wenfeng Zheng. 2019. Twin-delayed DDPG: A deep reinforcement learning technique to model a continuous movement of an intelligent robot agent. In *ICVISP*. 1–5.
- [6] Anlei Dong, Yi Chang, Zhaohui Zheng, Gilad Mishne, Jing Bai, Ruiqiang Zhang, Karolina Buchner, Ciya Liao, and Fernando Diaz. 2010. Towards recency ranking in web search. In *WSDM*. 11–20.
- [7] 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).
- [8] Scott Fujimoto, Herke Hoof, and David Meger. 2018. Addressing function approximation error in actor-critic methods. In *ICML*. PMLR, 1587–1596.
- [9] Scott Fujimoto, David Meger, and Doina Precup. 2019. Off-policy deep reinforcement learning without exploration. In *ICML*. PMLR, 2052–2062.
- [10] Bruno G Galuzzi, Ilaria Giordani, Antonio Candelieri, Riccardo Perego, and Francesco Archetti. 2020. Hyperparameter optimization for recommender systems through Bayesian optimization. *CMS* 17, 4 (2020), 495–515.
- [11] Yulong Gu, Zhuoye Ding, Shuaiqiang Wang, and Dawei Yin. 2020. Hierarchical User Profiling for E-commerce Recommender Systems. In *WSDM*. 223–231.
- [12] Yulong Gu, Zhuoye Ding, Shuaiqiang Wang, Lixin Zou, Yiding Liu, and Dawei Yin. 2020. Deep Multifaceted Transformers for Multi-objective Ranking in Large-Scale E-commerce Recommender Systems. In *CIKM*. 2493–2500.
- [13] Yulong Gu, Jiaying Song, Weidong Liu, and Lixin Zou. 2016. HLGPS: a home location global positioning system in location-based social networks. In *ICDM*. IEEE, 901–906.
- [14] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*. PMLR, 1861–1870.
- [15] Jianhua Han, Yong Yu, Feng Liu, Ruiming Tang, and Yuzhou Zhang. 2019. Optimizing Ranking Algorithm in Recommender System via Deep Reinforcement Learning. In *ALAM*. IEEE, 22–26.
- [16] Xinran He, Junfeng Pan, Ou Jin, Tianbing Xu, Bo Liu, Tao Xu, Yanxin Shi, Antoine Atallah, Ralf Herbrich, Stuart Bowers, et al. 2014. Practical lessons from predicting clicks on ads at facebook. In *ADKDD*. 1–9.
- [17] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. 2020. Conservative q-learning for offline reinforcement learning. *arXiv preprint arXiv:2006.04779* (2020).
- [18] Hoang Le, Cameron Voloshin, and Yisong Yue. 2019. Batch policy learning under constraints. In *International Conference on Machine Learning*. PMLR, 3703–3712.
- [19] Greg Linden, Brent Smith, and Jeremy York. 2003. Amazon.com recommendations: Item-to-item collaborative filtering. *IEEE Internet computing* 7, 1 (2003), 76–80.
- [20] Jiaqi Ma, Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H Chi. 2018. Modeling task relationships in multi-task learning with multi-gate mixture-of-experts. In *SIGKDD*. 1930–1939.
- [21] Melanie Mitchell. 1998. *An introduction to genetic algorithms*. MIT press.
- [22] Jonas Mockus. 1975. On Bayesian methods for seeking the extremum. In *Optimization techniques IFIP technical conference*. Springer, 400–404.
- [23] Wentao Ouyang, Xiuwu Zhang, Li Li, Heng Zou, Xin Xing, Zhaojie Liu, and Yanlong Du. 2019. Deep spatio-temporal neural networks for click-through rate prediction. In *SIGKDD*. 2078–2086.
- [24] Changhua Pei, Xinru Yang, Qing Cui, Xiao Lin, Fei Sun, Peng Jiang, Wenwu Ou, and Yongfeng Zhang. 2019. Value-aware recommendation based on reinforced profit maximization in e-commerce systems. *arXiv preprint arXiv:1902.00851* (2019).
- [25] Qi Pi, Weijie Bian, Guorui Zhou, Xiaoqiang Zhu, and Kun Gai. 2019. Practice on long sequential user behavior modeling for click-through rate prediction. In *SIGKDD*. 2671–2679.
- [26] Marco Tulio Ribeiro, Nivio Ziviani, Edleno Silva De Moura, Itamar Hata, Anisio Lacerda, and Adriano Veloso. 2014. Multiobjective pareto-efficient approaches for recommender systems. *TIST* 5, 4 (2014), 1–20.
- [27] Mario Rodriguez, Christian Posse, and Ethan Zhang. 2012. Multiple objective optimization in recommender systems. In *RecSys*. 11–18.
- [28] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. 2014. Deterministic policy gradient algorithms. In *ICML*. PMLR, 387–395.
- [29] Hongyan Tang, Junning Liu, Ming Zhao, and Xudong Gong. 2020. Progressive layered extraction (ple): A novel multi-task learning (mtl) model for personalized recommendations. In *RecSys*. 269–278.
- [30] Sebastian Thrun and Michael L Littman. 2000. Reinforcement learning: an introduction. *AI Magazine* 21, 1 (2000), 103–103.
- [31] Cameron Voloshin, Hoang M Le, Nan Jiang, and Yisong Yue. 2019. Empirical study of off-policy policy evaluation for reinforcement learning. *arXiv preprint arXiv:1911.06854* (2019).
- [32] Lu Wang, Wei Zhang, Xiaofeng He, and Hongyuan Zha. 2018. Supervised reinforcement learning with recurrent neural network for dynamic treatment recommendation. In *SIGKDD*. 2447–2456.
- [33] Yue Wu, Shuangfei Zhai, Nitish Srivastava, Joshua Susskind, Jian Zhang, Ruslan Salakhutdinov, and Hanlin Goh. 2021. Uncertainty Weighted Actor-Critic for Offline Reinforcement Learning. *arXiv preprint arXiv:2105.08140* (2021).
- [34] Xiangyu Zhao, Long Xia, Liang Zhang, Zhuoye Ding, Dawei Yin, and Jiliang Tang. 2018. Deep reinforcement learning for page-wise recommendations. In *RecSys*. 95–103.
- [35] Zhe Zhao, Lichan Hong, Li Wei, Jilin Chen, Aniruddh Nath, Shawn Andrews, Aditee Kumthekar, Maheswaran Sathiamoorthy, Xinyang Yi, and Ed Chi. 2019. Recommending what video to watch next: a multitask ranking system. In *RecSys*. 43–51.
- [36] Guorui Zhou, Na Mou, Ying Fan, Qi Pi, Weijie Bian, Chang Zhou, Xiaoqiang Zhu, and Kun Gai. 2019. Deep interest evolution network for click-through rate prediction. In *AAAI*, Vol. 33. 5941–5948.
- [37] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep interest network for click-through rate prediction. In *SIGKDD*. 1059–1068.
- [38] Lixin Zou, Long Xia, Zhuoye Ding, Jiaying Song, Weidong Liu, and Dawei Yin. 2019. Reinforcement learning to optimize long-term user engagement in recommender systems. In *SIGKDD*. 2810–2818.

Algorithm 2: Offline Evaluation Steps of Conservative Offline Policy Estimator $V(\pi_e)$.

Input : the transition dataset $\mathcal{D} = \{s_i, a_i, r_i, s'_i\}_{i=1}^N$, training batch size m , testing batch size n , discount factor γ , the penalty coefficient α , the policy π_e to be evaluated.

- 1 Initialize $Q_0(\cdot, \theta)$ randomly;
- 2 **foreach** $0 \leq k \leq K$ **do**
- 3 Sample training batch of m transitions $\{s_i, a_i, r_i, s'_i\}_{i=1}^m$ from \mathcal{D} randomly;
- 4 Compute current $\hat{Q}_i^e = Q_k(s_i, \pi_e(s_i))$ and $\hat{Q}_i = Q_k(s_i, a_i)$, target $\hat{y}_i = r_i + \gamma Q_k(s'_i, \pi_e(s'_i)) \quad \forall i$;
- 5 Construct training batch set:

$$\tilde{\mathcal{D}}_k = \{(s_i, a_i), \hat{Q}_i^e, \hat{Q}_i, \hat{y}_i\}_{i=1}^m;$$
- 6 Fit Q-function based on regression:

$$Q_{k+1} \leftarrow \underset{\theta}{\operatorname{argmax}} \alpha \cdot \left(\frac{1}{m} \sum_{i=1}^m \hat{Q}_i^e - \frac{1}{m} \sum_{i=1}^m \hat{Q}_i \right) + \frac{1}{2} \cdot \frac{1}{m} \sum_{i=1}^m (Q_k(s_i, a_i) - \hat{y}_i)^2;$$
- 7 **end**
- 8 Obtain fitted Q-function $\hat{Q}(\cdot, \theta) = \lim_{k \rightarrow K} \hat{Q}_k$;
- 9 Sample testing batch of n initial states $\{s_i^0\}_{i=1}^n$ from \mathcal{D} randomly;
- 10 Evaluate the offline policy π_e :

$$\hat{V}(\pi_e) = \frac{1}{n} \sum_{i=1}^n \sum_{a \sim \pi_e(a|s)} \pi_e(a|s_i^0) \hat{Q}(s_i^0, a).$$

Output: the offline evaluation result $\hat{V}(\pi_e)$.

A IMPLEMENTATION DETAILS OF CONSERVATIVE OFFLINE POLICY ESTIMATOR

In Algorithm 2, we elaborate the training and evaluation steps of our Conservative-OPEstimator. The Q network in Conservative-OPEstimator are MLP with ReLU activation function in hidden layers and are optimized based on Adam optimizer with the learning rate set to 0.1×10^{-3} . In addition, we set the training batch size, the testing batch size, the discount factor, the penalty coefficient and training epochs to $m = 512, n = 5000, \gamma = 0.95, \alpha = 0.5 \times 10^{-3}$ and $K = 5000$, respectively.

B SENSITIVITY ANALYSIS OF PARAMETERS

In this section, we will observe the influence on our model performance of two important model parameters, the perturbation bound

ρ and the learning rate η of the critic network, to figure out their optimal values. We evaluate Batch RL model class with the different value of ρ or η , and visualize the offline evaluation results in Figure 8.

B.1 Sensitivity of the Perturbation Bound ρ

As the behavior cloning-based method, Our Batch RL model reduces extrapolation error by constraining itself not to output OOD actions. It is essentially a fine-tuned policy based on training data generated by behavior policy. Once the batch data set constructed by the behavior policy is not the optimal set, our model will fall into the dilemma of local optimum. In order to explore potential high-value state-action pairs but avoid the action with noise to hurt user experience, we exploit the action perturbation network $P_\omega(s, a, \rho)$ to perturb actions cloned by VAE model, as discussed in Section 3.2.1. The perturbation bound ρ is a critical parameter that controls the range of exploration and influences the final output action.

As shown in Figure 8 (a), when $\rho = 0.15$, our model has the best performance and can achieve stable and high returns. With the decrease of ρ , the exploration for the action is conservative and narrow, which is conducive to obtain the model with stable returns but limits the model to seeking high returns; Instead, the exploration is more random and extensive, which is conducive to enrich the action candidate set but increases extrapolation error and the action with noise of the model. The above phenomenon also shows that ρ only adjusts the constraint range of the model on the output action, and the Q network of our model still overestimates the Q value of OOD actions. A large ρ means that there is a small constraint on the action. Therefore, our model may output OOD actions, resulting in extrapolation error.

B.2 Sensitivity of the Learning Rate η of the Critic Network.

The critic network of our Batch RL model, which can evaluate the value of state-action pairs, provides the actor network with the optimization direction and the basis for decision making. Therefore, it directly affects our model performance whether the critic network learns adequately and effectively.

In Figure 8 (b), we research how the learning rate η of the critic network influences the long-term rewards of our model. With the increase of learning rate η , the rewards of model fluctuates constantly. We note that, when the learning rate η is around 0.2×10^{-3} , the fluctuation flattens out and our model gains the highest returns. This shows that it is so appropriate to set $\eta = 0.2 \times 10^{-3}$ that the critic network does not fall into the local optimal because of a small η , nor miss the real optimal solution because of a large η . Meanwhile, the result ensures the critic network trained adequately that can accurately evaluate the value of state-action pairs.

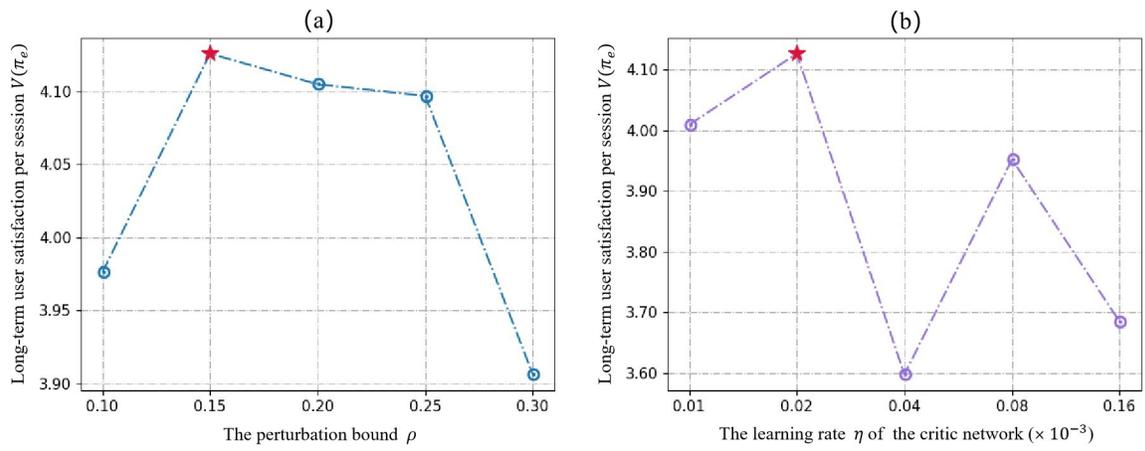


Figure 8: Parameter sensitivity of the perturbation bound ρ and the learning rate η of the critic network. The best result is marked by a red star.