

Value Penalized Q-Learning for Recommender Systems

Chengqian Gao,¹ Ke Xu² Peilin Zhao²,

¹ Shenzhen International Graduate School, Tsinghua University

² Tencent AI Lab

gcq19@mails.tsinghua.edu.cn, kaylakxu@tencent.com, masonzhao@tencent.com

Abstract

Scaling reinforcement learning (RL) to recommender systems (RS) is promising since maximizing the expected cumulative rewards for RL agents meets the objective of RS, i.e., improving customer’s long-term satisfaction. A key approach to this goal is offline RL, which aims to learn policies from logged data. However, the high-dimensional action space and the non-stationary dynamics in commercial RS intensify distributional shift issues, making it challenging to apply offline RL methods to RS. To alleviate the action distribution shift problem in extracting RL policy from static trajectories, we propose *Value Penalized Q-learning (VPQ)*, an uncertainty-based offline RL algorithm. It penalizes the unstable Q-values in the regression target by uncertainty-aware weights, without the need to estimate the behavior policy, suitable for RS with a large number of items. We derive the penalty weights from the variances across an ensemble of Q-functions. To alleviate distributional shift issues at test time, we further introduce the *critic* framework to integrate the proposed method with classic RS models. Extensive experiments conducted on two real-world datasets show that the proposed method could serve as a *gain plugin* for existing RS models.

1 Introduction

Practical recommender systems (RS) are usually trained for estimating a user’s immediate response to a slate of items, without considering the long-term utility. Reinforcement learning (RL) methods are designed to maximize the discounted cumulative rewards over time, which perfectly satisfies the original needs of RS. Driven by this motivation, there has been tremendous progress in developing reinforcement learning based recommender systems (RLRS) (Afsar, Crump, and Far 2021).

A promising approach for achieving RLRS is learning from the previously collected data because (1) there are lots of logged data available for training RL agents, and (2) it is expensive and risky to train commercial recommendation agents purely with online interactions. However, directly utilizing RL algorithms in the offline setting often results in poor performance, even for off-policy methods, which can leverage data collected by other policies in principle (Fujimoto, Meger, and Precup 2019; Levine et al. 2020). The main reason may be the overestimation for out-of-distribution (OOD) queries, induced by distributional shift, both at test and training time (Levine et al. 2020). Unlike the

general RL algorithms, offline RL (batch RL) concerns the problem of training agents from a static dataset, addressing the OOD predictions without access to interacting with the environment. The settings of offline RL meet the requirements of RLRS, making it achievable to build RLRS agents from logged data.

Although it is promising to utilize offline RL approach in RS, there are still challenges due to the different settings between RL and RS tasks. Firstly, the number of candidate items for RS agents is much larger than the action space in synthetic RL environments. With a larger action space, the number of in-distribution data points is much less than that of OOD queries, which represents as a sparse user-item interaction matrix. Therefore, the RLRS agents are more vulnerable to overestimation errors during training. Moreover, the large action space also intensifies the difficulty of estimating an accurate behavior policy for constraining the learned policy’s deviation from the logged data. As a result, solutions from offline RL domain are probably not suited for RS tasks. Secondly, the non-stationary dynamic of RS also differs from RL environments. In general offline RL problems, e.g. D4RL (Fu et al. 2020), trajectories are collected by policies with an environment whose dynamics are exactly as same as the test environment. While in recommendation scenarios, the environment dynamics constantly change as the user preferences changing over time.

This paper proposes Value Penalized Q-learning (VPQ), an uncertainty-based offline RL method, to fix the distributional shift issues during training. Then we combine it with different recommendation models, making a trade-off between the robust RS models and the alluring but risky RL methods.

Specifically, we use sequential or session-based recommendation models to *represent* the sequence of user-item interactions (states). We define the agent as a functional instance generating recommendations (actions) by the states received from users (environment) so as to maximize their long-term satisfaction, e.g., purchases and clicks. We then take two techniques to address the distributional shift issues for value-based RL methods in recommendation setting. Firstly, we penalize the unstable Q-values in regression target with uncertainty-aware weights, reducing the overestimation induced by the action distribution shift during the learning process. We train an ensemble of Q-functions with

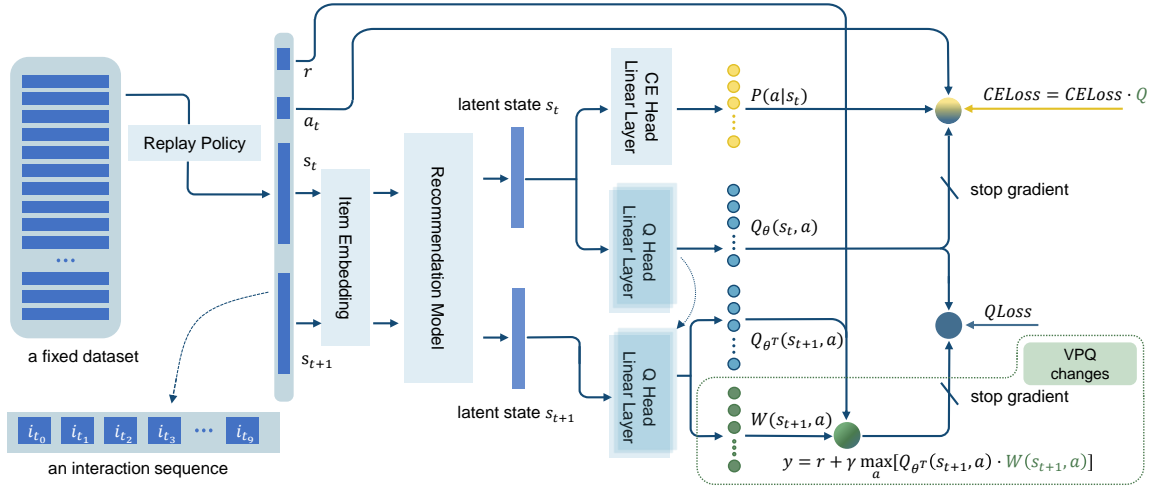


Figure 1: RLRS with VPQ. Recommendation models are employed to extract the hidden state for CE head (supervised head trained with cross-entropy loss to perform ranking) and Q head (trained for maximizing the accumulated rewards). We attain the conservative value function via penalizing the unstable targets by uncertainty-aware weights.

different random learning rates and set the variance of their predictions as the uncertainty quantifier. Compared with offline RL methods that subtract uncertainty from predictions, our method is empirically more effective for obtaining the conservative value function in the RS settings. Besides, we employ the critic framework (Xin et al. 2020) to make use of the learned value function while avoiding the performance degradation caused by the distributional shift issues at test time. The architecture of an RS instance integrated with VPQ is shown in Figure 1. At test time, we discard all other components except the CE head. Although it is a trade-off between the RL-based method and the supervised approach, experimental results show that we can still benefit from incorporating VPQ into RS models. We summarize the contributions of this work as follows.

- (1) We propose an uncertainty-based offline RL method, VPQ, to address the action distribution shift issue during training by penalizing the unstable predictions with uncertainty quantification.
- (2) We empirically show that, it is more effective to attain the conservative value function by multiplying the unstable predictions with uncertainty-aware weights in the recommendation scenarios.
- (3) Extensive experiments show that the benefit of integrating with VPQ contains better representation and reinforcing actions with long-term rewards.

2 Related work

There are several difficulties in building RLRS, such as the state representation (Lei et al. 2020; Liu et al. 2020; Wang et al. 2020b), the slated-based recommendation setting (Sunehag et al. 2015; Swaminathan et al. 2017; Chen et al. 2019; Gong et al. 2019; Ie et al. 2019) and the large-scale recommendation problem (Dulac-Arnold et al. 2015;

Zhao et al. 2019). However, there has been limited attention paid to the negative impact of offline setting, i.e., training RLRS agents from static datasets. Here we give a literature review on the general offline RL methods.

Constrain Policy Close to the Behavior Policy. Methods based on different distance metrics are proposed to compel the RL policy close to the behavior policy that produces the dataset (Levine et al. 2020), such as KL divergence (Jaques et al. 2019), Wasserstein distance (Wu, Tucker, and Nachum 2019), MMD (Kumar et al. 2019), Fisher divergence (Kostrikov et al. 2021) and Euclidean squared distance (Fujimoto and Gu 2021). In this approach, GCQ (Wang et al. 2020a) utilizes a neural generator to recover the distribution of dataset, which is similar to BCQ (Fujimoto, Meger, and Precup 2019). These methods may suffer from estimating the dataset trajectories distribution in RS.

Estimate Accurate/Conservative Q-values. Algorithms like Ensemble-DQN (Faußer and Schwenker 2015; Osband et al. 2016), QR-DQN (Dabney et al. 2018), and REM (Agarwal, Schuurmans, and Norouzi 2020) in this approach estimate Q-values by predicting the expectation of a distribution or an ensemble. Another instance is CQL (Kumar et al. 2020), which adds two regularization terms to the Q-function update so as to avoid overestimating OOD actions. Despite its good performance in general RL testbeds, CQL may suffer from high dimensional action space (minimization term in CQL) and non-user trajectories in the datasets (maximization term). The proposed VPQ can be categorized into this approach. Besides, a recently proposed method named UWAC (Wu et al. 2021) adds uncertainty-aware weights to the Q-function queries, combining the two main approaches.

3 Preliminaries

MDP on Recommendation Problem

We represent the recommendation problem as a Markov decision process (MDP) defined by a tuple $(\mathcal{S}, \mathcal{A}, T, d_0, r, \gamma)$:

- **State space** \mathcal{S} is a continuous space describing the user’s interaction history. The interaction sequence contains 10 items the user lastly interacted with. We then use the recommendation model to map it into a latent space as the state for Q heads and as the features for CE head.
- **Action space** \mathcal{A} is a finite set of actions of recommendation agents. For a state s_t , the candidate actions are all items in the dataset.
- **Initial state distribution** d_0 is the distribution that initial state follows.
- **Transition probability** T is a conditional probability with $T(s_{t+1}|s_t, a_t)$ describing the distribution over the next state s_{t+1} after taking action a_t at state s_t .
- **Reward function** r describes the immediate rewards. For example, agent will get a reward $r(s_t, a_t)$ when taking the action a_t at state s_t .
- **Discount factor** γ is a scalar that makes trade-off between immediate rewards and long-term rewards.

The goal of reinforcement learning is to find a policy $\pi(a_t|s_t)$, which defines a distribution over actions a_t conditioned on states s_t , so as to maximize the cumulative reward:

$$\mathbb{E}_{\tau \sim p_\pi(\tau)} \left[\sum_{t=0}^{|\tau|} \gamma^t r(s_t, a_t) \right], \quad (1)$$

where the trajectory distribution p_π for a transition probability T , a policy π and an initial state distribution d_0 is:

$$p_\pi(\tau) = d_0(s_0) \prod_{t=0}^{|\tau|-1} \pi(a_t|s_t) T(s_{t+1}|s_t, a_t) \quad (2)$$

Value-based RL

Value-based RL methods, such as DQN (Mnih et al. 2015), usually try to learn a state-action value function $Q_\theta(s, a)$ with parameter θ , defined as:

$$Q^\pi(s_t, a_t) = \mathbb{E}_{\tau \sim p_\pi} \left[\sum_{t'=t}^{|\tau|} \gamma^{t'-t} r(s_{t'}, a_{t'}) \right] \quad (3)$$

An improved policy $\pi'(a|s)$ can be easily obtained by $\pi'(a|s) = \arg \max_a Q_\theta(s, a)$. And the parameter θ for the value network is obtained by minimizing the following loss function with a given dataset D :

$$\mathcal{L}_\theta = \frac{1}{2} \mathbb{E}_{(s_t, a_t, r_t, s_{t+1}) \sim D} \left[(y_t - Q_\theta(s_t, a_t))^2 \right], \quad (4)$$

where the item $y_t = r_t + \gamma \max_{a'} Q_{\theta_T}(s_{t+1}, a')$ is the so called target, and Q_{θ_T} is the target network with frozen parameters from historical value function Q_θ .

Challenges: Distributional Shift in Value-based RL

A key source of instability of value-based RL methods applying on static datasets is the distributional shift during the training and test phase (Levine et al. 2020). Distributional shift happens when evaluating a trained policy on another different distribution. It can be categorized into two types: state distribution shift and action distribution shift. Here, we explain these issues within the value-based RL methods.

State distribution shift affects the performance of RL agents only at test time. The learned policy will always meet unseen states at test time since the dataset usually does not fully cover the state space. The output of the Q-function can be unstable as the loss function in Equation (4) makes no guarantees about the error when encountering OOD state inputs, resulting in performance degradation at test time (Lee et al. 2021). Unfortunately, the narrow distribution of the collected trajectories in RS datasets and the non-stationary environments exacerbate this problem.

Action distribution shift affects the agent both at test time and training time. The Q-function is trained on dataset state-action pairs but is evaluated on all valid actions for policy improvement and real actions generation. At test time, the unstable predictions on OOD actions may result in unreasonable actions as the side effect of the greedy policy. During training, the highest value among noisy predictions in target will be taken as the future reward and then spread to the learned Q-function, leading to a persistent overestimation bias (Jaques et al. 2019; Kumar et al. 2019). In the online setting, active learning around these overestimated state-action pairs would correct them. By contrast, such correction will never happen in the offline setting since the agents can no longer interact with the environment. Moreover, the large action space in RLRS exacerbates this problem. For example, the general RL testbed Arcade Learning Environment (Bellemare et al. 2013) consists of no more than 18 discrete actions, while the candidate recommending items in the Retailrocket dataset, a real-world commercial dataset, is 70,852 (please refer to Table 1 for the details).

4 Methods

To address the distributional shift issues for training agents on offline RS datasets, we propose VPQ. It learns a conservative value function by penalizing the unstable regression targets with uncertainty-aware weights. After that, we integrate the proposed method with classic RS models to alleviate the distributional shift issues at test time. The training paradigm of the integrated model is shown in Figure 1 and the proposed VPQ is illustrated as Algorithm 1.

Penalize the Unstable Q-values

The main concern of this work is how to *use* the uncertainty metric to penalize the unstable predictions in the target network, so as to alleviate the impact from the distributional shift in training agents for RLRS. In this part, we firstly give two different ways of using the uncertainty quantifier, and then we analyse the advantages of our method, finally we give details about the proposed method.

Two Different Ways of Using the Uncertainty Metric.

An intuitive approach for uncertainty-based offline RL is to regress to the following target (Levine et al. 2020; Buckman, Gelada, and Bellemare 2020; Jin, Yang, and Wang 2020):

$$y = r + \gamma \max_a (Q(s_{t+1}, a) - \lambda \text{Unc}), \quad (5)$$

where Unc quantifies the amount of uncertainty for each query. We denote this family of algorithms as *p-sub* for they remove the uncertainty by subtracting.

Our approach to penalizing the unstable prediction is:

$$y = r + \gamma \max_a (Q(s_{t+1}, a) \cdot W(s_{t+1}, a)), \quad (6)$$

where $W(s_t, a)$ is a uncertainty weight for each prediction. In our practice it is designed as:

$$W(s_t, a_t) = \frac{1}{1 + \lambda \text{Unc}}, \quad (7)$$

where $\lambda > 0$ is a scaling factor that controls the amount of uncertainty. Unc is defined as the standard deviation across an ensemble of target Q-functions for each state-action pair, in this work. We denote this formulation as *p-mul* for the stable predictions are calculated by multiplying the unstable values by uncertainty-aware weights.

The Advantage of the Proposed Penalty Formulation.

We here give an analysis on the two families of penalty formulations in a large action space.

With the assumption that unstable predictions on OOD actions for a given state, i.e., $Q(s, a)$, are i.i.d. random variables follow a normal distribution $x \sim \mathcal{N}(\mu, \sigma^2)$, then the target value can be described as:

$$y = r + \gamma \max_{1 \leq i \leq n; x_i \sim \mathcal{N}(\mu, \sigma^2)} x_i, \quad (8)$$

where n equals to the number of OOD actions, and we denote Q_T as the maximum value of the n samples. Since the exact calculation of Q_T can be very hard, we here use the approximate form from the literature on normal order statistics (Blom 1958; Harter 1961; Royston 1982), i.e.,

$$\mathbb{E} \left[\max_{1 \leq i \leq n; x_i \sim \mathcal{N}(\mu, \sigma^2)} x_i \right] \approx \mu + \sigma \Phi^{-1} \left(\frac{n - 0.375}{n - 2 \times 0.375 + 1} \right), \quad (9)$$

where Φ^{-1} is the inverse of the standard normal CDF.

By using the properties of the expectation and max operator, we have the expected value of the penalized Q_T with the *p-sub* formulation:

$$\mathbb{E} [Q_T - \lambda \sigma] = \mu + (C_0 - \lambda) \sigma \quad (10)$$

and the expectation of that with the *p-mul* form:

$$\mathbb{E} \left[\frac{Q_T}{1 + \lambda \sigma} \right] = \frac{1}{1 + \lambda \sigma} (\mu + \sigma C_0), \quad (11)$$

where C_0 is a constant number for a given n .

There are two benefits of the proposed penalty formulation. Firstly, the penalized Q-values from *p-mul* is more robust than those from *p-sub*. For *p-sub*, a large λ would incur target values less than zero, while a small scaling factor could not reduce the unstable predictions. Unfortunately,

Algorithm 1: VPQ: Value Penalized Q-learning

Input: K randomly initialized Q -functions with parameters $\{\theta_k\}_{k=1}^K$, a mixture distribution P_Δ for Random Mixture, a scale factor λ for VPQ and an offline dataset D

Output: $\{\theta_k\}_{k=1}^K$

- 1: Initialize all trainable parameters.
- 2: **while** not done **do**
- 3: Sample a mini-batch of transitions (s_t, a_t, r, s_{t+1}) randomly from D
- 4: Sample mixture weights $\{\alpha_k\}_{k=1}^K$ from P_Δ
- 5: Compute the sample standard deviation $\tilde{\sigma}_T(s_{t+1}, a)$
- 6: Compute random mixture of target networks:
 $\tilde{\mu}_T(s_{t+1}, a) = \sum_{k=1}^K \alpha_k Q_{\theta_k^T}(s_{t+1}, a)$
- 7: Compute the penalty factor:
 $W(s_{t+1}, a) = (1 + \lambda \tilde{\sigma}_T(s_{t+1}, a))^{-1}$
- 8: Compute the penalized target:
 $y_t = r + \gamma \max_a (\tilde{\mu}_T(s_{t+1}, a) \cdot W(s_{t+1}, a))$
- 9: Update each head with random weight α_k :
 $\mathcal{L}_\theta = \frac{1}{2} (y_t - \sum_{k=1}^K \alpha_k \cdot Q_{\theta_k}(s_t, a_t))^2$
- 10: **end while**

for training RLRS agents, C_0 increases with an increase in the dimension of action space. To reduce unstable predictions, such a method has to increase the λ in Equation (5), which increases the risk of producing negative Q-values. As we will show in the next part, Q-values that below than zero would harm the performance. By contrast, the penalty from *p-mul* would always keep the target value greater than zero, even with a large scale factor. Secondly, the proposed *p-mul* formulation could heavily penalize the unstable and high predictions, while the *p-sub* formulation mainly penalizes the unstable predictions with small values. To alleviate the distributional shift at training time, it is more crucial to reduce the uncertainty within the large values because of the max operator. Finally, to illustrate the difference between the two forms, we give a toy experiment in Figure 2.

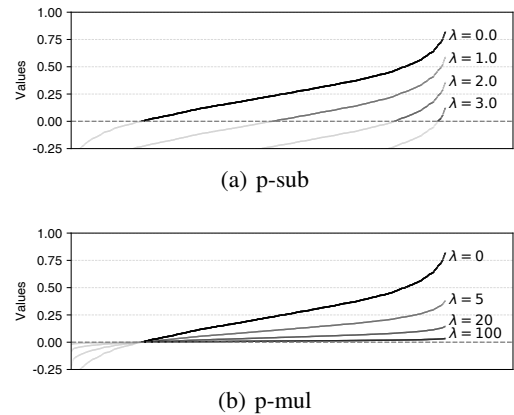


Figure 2: Penalize the simulated unstable Q-values in different ways. Heavy black lines are dense points sampled from a Gaussian distribution and sorted by their values. Scaling factor λ controls the intensity of penalization.

Details about the Proposed Algorithm. Inspired by the powerful penalty formulation, $p\text{-mul}$, we develop VPQ. For this paper, the uncertainty metric is defined as the standard deviation across an ensemble of Q-functions, i.e., a fully connected layer with K heads. We use the Random Ensemble Mixture (REM) technique (Agarwal, Schuurmans, and Norouzi 2020) to update each head with a random learning rate towards a shared regression target, improving the diversity of the learned Q-functions. The weights for mixing are sampled from a categorical distribution: first sample a set of K values i.i.d. from the uniform distribution $U(0, 1)$ and then normalize them.

Introducing the uncertainty eliminates the need to estimate the behavior policy for constraint. Besides, as an uncertainty-based offline RL method, VPQ may enjoy benefits from the progress in uncertainty measurement. Finally, we summarize our method in Algorithm 1.

The Critic Framework

To alleviate the action and state distribution shift at test time, we introduce the critic framework (Xin et al. 2020), which is inspired by the Actor-Critic framework (Konda 2002). Specifically, the Q-function learned from the proposed VPQ algorithm is employed to reweight the corresponding cross-entropy loss, i.e.,

$$\mathcal{L}_\phi = CE(s_t, a) \cdot \text{no_gradient}(Q(s_t, a)) \quad (12)$$

At test time, we only use the CE head to generate recommendations. The motivation is that the CE head is naturally free from the action distribution problem since it is optimized for all actions over each dataset state. Although it is a trade-off between the classic recommendation paradigms and the alluring RL methods, the critic framework may improve the performance of RS with stability.

5 Experiments

In this part we conduct experiments to verify the effectiveness of the proposed VPQ against online and offline RL methods in the recommendation setting.

Experimental Settings

Datasets. The datasets used in the experiments ¹ are Retailrocket ² and Yoochoose ³. We follow the data processing details in (Xin et al. 2020) and split each into two parts: 80% for training and 20% for testing. Statistics of the final datasets are summarized in Table 1.

Evaluation Metrics. We use two metrics: HR (hit ratio) for recall and NDCG (normalized discounted cumulative gain) for rank. We compute the ranking performance on clicked and ordered items separately. For example, we define HR for clicks as:

$$\text{HR}(\text{click}) = \frac{\# \text{ hits among clicks}}{\# \text{ clicks}} \quad (13)$$

¹We will release our code soon.

²<https://www.kaggle.com/retailrocket/e-commerce-dataset>

³<https://recsys.acm.org/recsys15/challenge/>

Table 1: Statistics of datasets used in the experiments.

| Statistics | Retailrocket | Yoochoose 200k |
|---------------------|--------------|----------------|
| # items | 70,852 | 26,527 |
| # training sessions | 156,419 | 160,000 |
| # test sessions | 39,105 | 40,000 |
| # purchase | 57,269 | 44,039 |
| # clicks | 1,176,680 | 1,109,869 |
| Averaged length | 6.31 | 5.77 |

Parameter Settings. Following (Xin et al. 2020), we use RS models to represent the states. The item embedding size and the dimension of latent representations (states) in the RS models are set as 64. For RL methods, the rewards r of clicked and purchased items are set as 0.2, 1.0, respectively. The discount factor γ is 0.5, for the short trajectories in RS task. For VPQ, the scale factor λ is set to 20, and the number of heads K is 15. The optimizer for minimizing Q loss and CE loss is Adam. The learning rate is set to 0.005. The mini-batch size is 256, and we evaluate the performance every 2000 batches.

Performance Gains from VPQ

To verify the effectiveness of VPQ in recommendation task, we integrate it with four classic RS methods: GRU (Hidasi et al. 2016), Caser (Tang and Wang 2018), NextItNet (Yuan et al. 2019) and SASRec (Kang and McAuley 2018). We also compare its performance with the following online RL and offline RL-based methods:

- **SAC** and **SQN** Self-Supervised Actor-Critic (SAC) and Self-Supervised Q-learning (SQN) (Xin et al. 2020) are RL plugin methods for RS. The main difference between the two methods is SAC learns a Q-function for reweighting the CE loss, while SQN only minimizes the TD error as an auxiliary loss for better representation.
- **REM** An RL method that trains an ensemble of Q-functions for stable predictions (Agarwal, Schuurmans, and Norouzi 2020). Our implementation uses a shared regression target for its robustness. Comparison with this approach tells us whether the improvement achieved by our method comes from the ensemble framework.
- **Minus** Obtaining a conservative value function with the $p\text{-sub}$ formulation in Equation (5). We choose the standard deviation across REM heads as the uncertainty. From the previous analysis, this method should face the dilemma of choosing an appropriate scale factor.
- **CQL** As a state-of-the-art algorithm in general offline RL, CQL (Kumar et al. 2020) performs well on both discrete and continuous control domains. Our implementation uses the $\text{CQL}(\mathcal{H})$ loss as the following:

$$\begin{aligned} \mathcal{L}_\theta = & \frac{1}{2} \mathbb{E}_{(s_t, a_t, r_t, s_{t+1}) \sim D} \left[(y_t - Q_\theta(s_t, a_t))^2 \right] \\ & + \alpha \mathbb{E}_{s \sim D} \left(\log \sum_{a \sim A} \exp(Q_\theta(s, a)) - \mathbb{E}_{a \sim \hat{\pi}_\beta} [Q_\theta(s, a)] \right) \end{aligned} \quad (14)$$

Table 2: Overall performance comparison on two recommendation datasets, averaged over 5 runs. NG is short for NDCG.

| Models | Retailrocket | | | | Yoochoose | | | | Total |
|--------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|-----------------------|
| | purchase | | click | | purchase | | click | | |
| | HR@20 | NG@20 | HR@20 | NG@20 | HR@20 | NG@20 | HR@20 | NG@20 | |
| GRU | 0.5414 | 0.3859 | 0.3067 | 0.196 | 0.62 | 0.3517 | 0.4771 | 0.2611 | 3.1401 (0.00%) |
| GRU-SAC | 0.5646 | 0.4219 | 0.3259 | 0.2127 | 0.6726 | 0.3957 | 0.4743 | 0.26 | 3.3276 (5.97%) |
| GRU-SQN | 0.5873 | 0.4291 | 0.3399 | 0.2186 | 0.645 | 0.3683 | 0.4909 | 0.2703 | 3.3495 (6.67%) |
| GRU-REM | 0.561 | 0.4168 | 0.3239 | 0.2119 | 0.6731 | 0.398 | 0.4752 | 0.2607 | 3.3206 (5.75%) |
| GRU-Minus | 0.5712 | 0.4241 | 0.3275 | 0.2143 | 0.6727 | 0.3978 | 0.4727 | 0.2591 | 3.3394 (6.35%) |
| GRU-CQL | 0.5251 | 0.3783 | 0.2957 | 0.1908 | 0.626 | 0.3565 | 0.4673 | 0.2558 | 3.0955 (-1.42%) |
| GRU-UWAC | 0.5618 | 0.409 | 0.3303 | 0.2126 | 0.6719 | 0.3973 | 0.4775 | 0.2636 | 3.3239 (5.85%) |
| GRU-VPQ | 0.5699 | 0.4263 | 0.3324 | 0.2188 | 0.6817 | 0.4062 | 0.4811 | 0.2654 | 3.3817 (7.69%) |
| Caser | 0.4226 | 0.3076 | 0.2813 | 0.1875 | 0.6607 | 0.3911 | 0.4613 | 0.2545 | 2.9665 (0.00%) |
| Caser-SAC | 0.4617 | 0.3336 | 0.3006 | 0.1985 | 0.6889 | 0.4173 | 0.4513 | 0.2472 | 3.0991 (4.47%) |
| Caser-SQN | 0.4263 | 0.3121 | 0.2789 | 0.1868 | 0.6667 | 0.3961 | 0.4626 | 0.2552 | 2.9847 (0.61%) |
| Caser-REM | 0.4711 | 0.339 | 0.3019 | 0.1999 | 0.6916 | 0.4216 | 0.4481 | 0.2452 | 3.1184 (5.12%) |
| Caser-Minus | 0.4618 | 0.3325 | 0.2991 | 0.1984 | 0.7029 | 0.433 | 0.4429 | 0.2419 | 3.1124 (4.92%) |
| Caser-CQL | 0.4699 | 0.3376 | 0.3048 | 0.2009 | 0.6705 | 0.4008 | 0.4598 | 0.2527 | 3.0970 (4.40%) |
| Caser-UWAC | 0.3884 | 0.2851 | 0.2588 | 0.1714 | 0.6758 | 0.4126 | 0.4509 | 0.2464 | 2.8893 (-2.60%) |
| Caser-VPQ | 0.4775 | 0.3454 | 0.3067 | 0.2036 | 0.7081 | 0.435 | 0.4459 | 0.2439 | 3.1661 (6.73%) |
| Next | 0.6564 | 0.4961 | 0.3329 | 0.213 | 0.5934 | 0.3352 | 0.4937 | 0.2715 | 3.3922 (0.00%) |
| Next-SAC | 0.6741 | 0.5354 | 0.3346 | 0.2149 | 0.5899 | 0.3383 | 0.473 | 0.2584 | 3.4186 (0.78%) |
| Next-SQN | 0.687 | 0.5276 | 0.3481 | 0.2215 | 0.5962 | 0.3386 | 0.4974 | 0.274 | 3.4905 (2.90%) |
| Next-REM | 0.6839 | 0.5387 | 0.3433 | 0.2216 | 0.6025 | 0.3509 | 0.4781 | 0.2623 | 3.4812 (2.62%) |
| Next-Minus | 0.6835 | 0.5432 | 0.3446 | 0.2227 | 0.6047 | 0.3511 | 0.4771 | 0.2622 | 3.4890 (2.85%) |
| Next-CQL | 0.6845 | 0.5399 | 0.3423 | 0.2215 | 0.6073 | 0.3509 | 0.4802 | 0.2635 | 3.4901 (2.89%) |
| Next-UWAC | 0.684 | 0.531 | 0.3452 | 0.2206 | 0.6155 | 0.3583 | 0.4861 | 0.2694 | 3.5102 (3.48%) |
| Next-VPQ | 0.699 | 0.5714 | 0.3556 | 0.233 | 0.6226 | 0.3672 | 0.4944 | 0.2744 | 3.6175 (6.64%) |
| SASRec | 0.6387 | 0.4599 | 0.3516 | 0.219 | 0.663 | 0.3733 | 0.5044 | 0.2761 | 3.4861 (0.00%) |
| SASRec-SAC | 0.6981 | 0.5629 | 0.3603 | 0.2347 | 0.6786 | 0.3955 | 0.5015 | 0.276 | 3.7077 (6.35%) |
| SASRec-SQN | 0.6749 | 0.501 | 0.3744 | 0.237 | 0.6707 | 0.3827 | 0.514 | 0.2839 | 3.6386 (4.37%) |
| SASRec-REM | 0.6655 | 0.5113 | 0.365 | 0.2361 | 0.6809 | 0.3986 | 0.5033 | 0.2779 | 3.6387 (4.37%) |
| SASRec-Minus | 0.6666 | 0.5125 | 0.3634 | 0.2353 | 0.6815 | 0.3965 | 0.5038 | 0.2783 | 3.6378 (4.35%) |
| SASRec-CQL | 0.7046 | 0.5599 | 0.3749 | 0.2404 | 0.6631 | 0.3805 | 0.5013 | 0.2751 | 3.6997 (6.12%) |
| SASRec-UWAC | 0.6657 | 0.5015 | 0.3715 | 0.2367 | 0.6786 | 0.3959 | 0.5126 | 0.285 | 3.6475 (4.63%) |
| SASRec-VPQ | 0.7171 | 0.5914 | 0.3785 | 0.2484 | 0.6841 | 0.4063 | 0.5081 | 0.2814 | 3.8153 (9.44%) |

- **UWAC** An actor-critic based offline RL algorithm that adds uncertainty-aware weights on Bellman loss and the actor loss (Wu et al. 2021). We incorporate this method to REM by adding its uncertainty-aware weights to the final Q-function estimates, i.e.,

$$\mathcal{L}_\phi = CE(s_t, a) \cdot \text{no_gradient}\left(\frac{\beta \cdot Q(s_t, a)}{\text{Unc}}\right) \quad (15)$$

All the above methods are embedded with the four classic RS models under the same critic framework, except SQN. We select the best hyper-parameter λ for Minus, α for CQL, and β for UWAC by sweeping on each dataset. λ for VPQ is set to 20 for both datasets. Table 2 shows the top-20 performance of prototypical models integrated with or without different online or offline RL methods.

Compared to the corresponding original RS methods, almost all combined methods perform better. The consistent improvement from SQN implies that minimizing the TD error could benefit the representation learning. Improvements from other RL methods is unstable, suggesting that accurate

Q-values may improve the performance while an unstable value function could harm the recommendation.

Compared to the RL methods, VPQ performs better. CQL performs not so well as the minimization item in its loss function would produce Q-values below zero, which is unreasonable in our setting and may impact its final performance. REM behaves similar to SAC, indicating that a simple ensemble method could not address the difficult distributional shift issues. The advantage compared to REM also illustrates that the performance gain from VPQ is not from the ensemble method but from successfully suppressing the unstable Q-values. Besides, the advantage compared to UWAC suggests it is more effective to avoid Q-function explosion by penalizing at bootstrapping time.

Compared to the Minus method, VPQ can produce more conservative and stable Q-values for reweighting the CE loss in our observation, while the mean and variance of Q-values from Minus and REM are higher than that of VPQ. This may suggest that the Minus methods could not reduce the unstable predictions in recommendation scenarios.

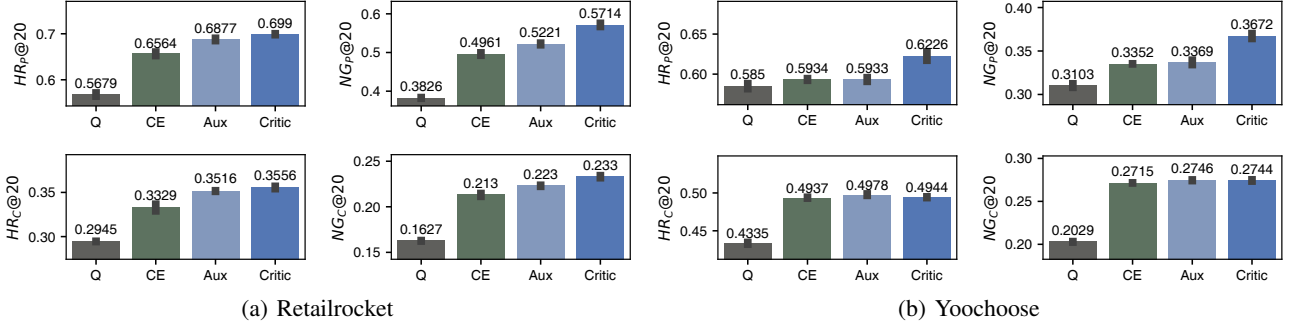


Figure 3: Comparison of performance when integrating Q head in different ways. Error bars show standard deviation.

Ablation Study: Where do the gains come from?

We compare four ways of using the learned Q-function: (1) *Q*, only train the Q head and generate recommendation according to Q-values, (2) *CE*, only optimize the CE head (the prototypical methods), (3) *Aux*, use the TD error as auxiliary loss for representation learning (SQN way), (4) *Critic*, use the learned Q-function as a critic for CE loss (SAC and VPQ way). The base model is NextItNet, and all hyper-parameters are the same as the previous experiments.

The results are shown in Figure 3. *Q* achieves the worst performance on both datasets, although we use the negative feedback (Zhao et al. 2018) technique to stabilize the training process. This verifies the necessity of considering the impact of distributional shift at test time. *Aux* surpasses *CE*, suggesting that minimizing the TD error would benefit representation learning. Moreover, by incorporating the learned Q-values with the CE loss, *Critic* outperforms *Aux*, achieves the best performance generally.

The results above suggest that improvement of incorporating with VPQ comes from better state representation (serving as an auxiliary loss) and exploiting more accurate Q-values without destroying performance boost from the former (successfully mitigating the distributional shift issues).

Hyper-parameter Study

We conduct experiments to investigate the effect of the scaling factor (λ in VPQ). We give the performance and the values of key indicators of VPQ with different λ in Figure 4. The base model is NextItNet, the dataset is Yoochoose and the key indicators are the mean and standard deviation of the Q-values queried for reweighting.

When λ increases from 1 to 20, the recommendation performance improves continually (Figure 4(a)) with the decreasing of the key indicators (Figure 4(b)), suggesting that the scaling factor λ controls the amount of penalization on unstable predictions and thereby affects the performance. When λ is greater than 20, increasing the penalization does not lead to a more conservative Q-function, i.e., the key indicators (almost) stop decrease with the λ increases. Instead, it hurts the performance of predicting items with high rewards (purchased items). Besides, the optimal λ for other three base models and the Retailrocket dataset is also 20.

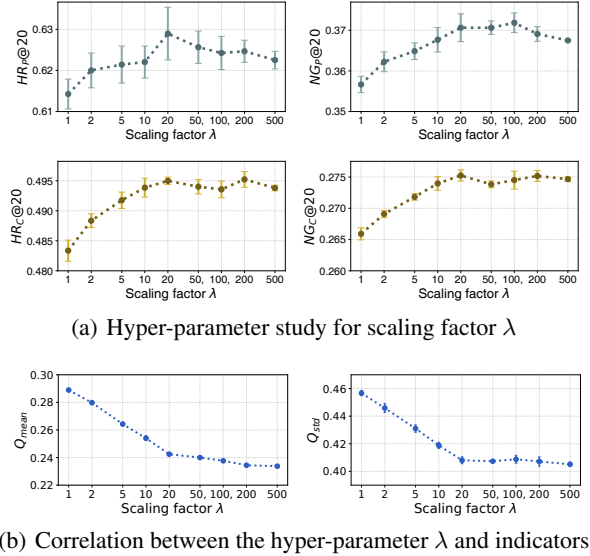


Figure 4: Sweeping the scaling factor on Yoochoose dataset. We plot the scaling factor λ using a base-e logarithmic scale. Error bars show standard deviation.

6 Conclusion

Scaling RL methods into recommendation tasks is challenging due to the severe distributional shift issues in RS settings. In this paper we proposed Value Penalized Q-learning (VPQ) to address the distributional shift during training. It is an uncertainty-based offline RL algorithm that avoids estimation of difficult behavior policy or constraints on the learned policy, and is suitable for RS with large item set. In addition, we illustrate its properties and empirically demonstrate its advantage over other offline RL methods in terms of robustness and effectiveness. The consistent improvement from incorporating with VPQ shows that it could be a gain plugin for recommendation models.

References

Afsar, M. M.; Crump, T.; and Far, B. 2021. Reinforcement learning based recommender systems: A survey. arXiv:2101.06286.

- Agarwal, R.; Schuurmans, D.; and Norouzi, M. 2020. An Optimistic Perspective on Offline Reinforcement Learning. In *ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, 104–114. PMLR.
- Bellemare, M. G.; Naddaf, Y.; Veness, J.; and Bowling, M. 2013. The Arcade Learning Environment: An Evaluation Platform for General Agents. *Journal of Artificial Intelligence Research*, 47: 253–279.
- Blom, G. 1958. *Statistical Estimates and Transformed Beta Variables*. Sweden: Almqvist & Wiksell, John Wiley & Sons, Inc.
- Buckman, J.; Gelada, C.; and Bellemare, M. G. 2020. The Importance of Pessimism in Fixed-Dataset Policy Optimization. arXiv:2009.06799.
- Chen, M.; Beutel, A.; Covington, P.; Jain, S.; Belletti, F.; and Chi, E. H. 2019. Top-K Off-Policy Correction for a REINFORCE Recommender System. In *WSDM 2019, Melbourne, VIC, Australia, February 11-15, 2019*, 456–464. ACM.
- Dabney, W.; Rowland, M.; Bellemare, M. G.; and Munos, R. 2018. Distributional Reinforcement Learning With Quantile Regression. In *AAAI 2018, New Orleans, Louisiana, USA, February 2-7, 2018*, 2892–2901. AAAI Press.
- Dulac-Arnold, G.; Evans, R.; van Hasselt, H.; Sunehag, P.; Lillicrap, T.; Hunt, J.; Mann, T.; Weber, T.; Degris, T.; and Coppin, B. 2015. Deep Reinforcement Learning in Large Discrete Action Spaces. arXiv:1512.07679.
- Faußer, S.; and Schwenker, F. 2015. Neural Network Ensembles in Reinforcement Learning. *Neural Processing Letters*, 41(1): 55–69.
- Fu, J.; Kumar, A.; Nachum, O.; Tucker, G.; and Levine, S. 2020. D4RL: Datasets for Deep Data-Driven Reinforcement Learning. arXiv:2004.07219.
- Fujimoto, S.; and Gu, S. S. 2021. A Minimalist Approach to Offline Reinforcement Learning. arXiv:2106.06860.
- Fujimoto, S.; Meger, D.; and Precup, D. 2019. Off-Policy Deep Reinforcement Learning without Exploration. In *ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, 2052–2062. PMLR.
- Gong, Y.; Zhu, Y.; Duan, L.; Liu, Q.; Guan, Z.; Sun, F.; Ou, W.; and Zhu, K. Q. 2019. Exact-K Recommendation via Maximal Clique Optimization. In *SIGKDD 2019, Anchorage, AK, USA, August 4-8, 2019*, 617–626. ACM.
- Harter, H. L. 1961. Expected values of normal order statistics. *Biometrika*, 48(1 and 2): 151–165.
- Hidasi, B.; Karatzoglou, A.; Baltrunas, L.; and Tikk, D. 2016. Session-based Recommendations with Recurrent Neural Networks. In *ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*.
- Ie, E.; Jain, V.; Wang, J.; Narvekar, S.; Agarwal, R.; Wu, R.; Cheng, H.; Chandra, T.; and Boutilier, C. 2019. SlateQ: A Tractable Decomposition for Reinforcement Learning with Recommendation Sets. In *IJCAI, 2019, Macao, China, August 10-16, 2019*, 2592–2599. ijcai.org.
- Jaques, N.; Ghandeharioun, A.; Shen, J. H.; Ferguson, C.; Lapedriza, À.; Jones, N.; Gu, S.; and Picard, R. W. 2019. Way Off-Policy Batch Deep Reinforcement Learning of Implicit Human Preferences in Dialog. arXiv:1907.00456.
- Jin, Y.; Yang, Z.; and Wang, Z. 2020. Is Pessimism Provably Efficient for Offline RL? arXiv:2012.15085.
- Kang, W.; and McAuley, J. J. 2018. Self-Attentive Sequential Recommendation. In *ICDM 2018, Singapore, November 17-20, 2018*, 197–206. IEEE Computer Society.
- Konda, V. 2002. *Actor-critic Algorithms*. Ph.D. thesis, Massachusetts Institute of Technology, Cambridge, MA, USA.
- Kostrikov, I.; Fergus, R.; Tompson, J.; and Nachum, O. 2021. Offline Reinforcement Learning with Fisher Divergence Critic Regularization. In *ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, 5774–5783. PMLR.
- Kumar, A.; Fu, J.; Soh, M.; Tucker, G.; and Levine, S. 2019. Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction. In *NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, 11761–11771.
- Kumar, A.; Zhou, A.; Tucker, G.; and Levine, S. 2020. Conservative Q-Learning for Offline Reinforcement Learning. In *NeurIPS 2020, December 6-12, 2020, virtual*.
- Lee, S.; Seo, Y.; Lee, K.; Abbeel, P.; and Shin, J. 2021. Offline-to-Online Reinforcement Learning via Balanced Replay and Pessimistic Q-Ensemble. arXiv:2107.00591.
- Lei, Y.; Pei, H.; Yan, H.; and Li, W. 2020. Reinforcement Learning based Recommendation with Graph Convolutional Q-network. In *SIGIR 2020, Virtual Event, China, July 25-30, 2020*, 1757–1760. ACM.
- Levine, S.; Kumar, A.; Tucker, G.; and Fu, J. 2020. Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems. arXiv:2005.01643.
- Liu, F.; Tang, R.; Li, X.; Zhang, W.; Ye, Y.; Chen, H.; Guo, H.; Zhang, Y.; and He, X. 2020. State Representation Modeling for Deep Reinforcement Learning based Recommendation. *Knowledge-Based Systems*, 205: 106–170.
- Mnih, V.; Kavukcuoglu, K.; Silver, D.; Rusu, A. A.; Veness, J.; Bellemare, M. G.; Graves, A.; Riedmiller, M. A.; Fidjeland, A.; Ostrovski, G.; Petersen, S.; Beattie, C.; Sadik, A.; Antonoglou, I.; King, H.; Kumaran, D.; Wierstra, D.; Legg, S.; and Hassabis, D. 2015. Human-level Control through Deep Reinforcement Learning. *nature*, 518(7540): 529–533.
- Osband, I.; Blundell, C.; Pritzel, A.; and Roy, B. V. 2016. Deep Exploration via Bootstrapped DQN. In *NeurIPS 2016, December 5-10, 2016, Barcelona, Spain*, 4026–4034.
- Royston, J. P. 1982. Expected normal order statistics (exact and approximate). *Journal of the Royal Statistical Society Series C (Applied Statistics)*, 31(2): 161–165.
- Sunehag, P.; Evans, R.; Dulac-Arnold, G.; Zwols, Y.; Visentin, D.; and Coppin, B. 2015. Deep Reinforcement Learning with Attention for Slate Markov Decision Processes with High-Dimensional States and Actions. arXiv:1512.01124.

Swaminathan, A.; Krishnamurthy, A.; Agarwal, A.; Dudík, M.; Langford, J.; Jose, D.; and Zitouni, I. 2017. Off-policy Evaluation for Slate Recommendation. In *NeurIPS 2017, December 4-9, 2017, Long Beach, CA, USA*, 3632–3642.

Tang, J.; and Wang, K. 2018. Personalized Top-N Sequential Recommendation via Convolutional Sequence Embedding. In *WSDM 2018, Marina Del Rey, CA, USA, February 5-9, 2018*, 565–573. ACM.

Wang, C.; Zhou, T.; Chen, C.; Hu, T.; and Chen, G. 2020a. Off-Policy Recommendation System Without Exploration. In *PAKDD 2020, Singapore, May 11-14, 2020, Proceedings, Part I*, volume 12084 of *Lecture Notes in Computer Science*, 16–27. Springer.

Wang, P.; Fan, Y.; Xia, L.; Zhao, W.; Niu, S.; and Huang, J. 2020b. KERL: A Knowledge-Guided Reinforcement Learning Model for Sequential Recommendation. In *SIGIR 2020, Virtual Event, China, July 25-30, 2020*, 209–218. ACM.

Wu, Y.; Tucker, G.; and Nachum, O. 2019. Behavior Regularized Offline Reinforcement Learning. arXiv:1911.11361.

Wu, Y.; Zhai, S.; Srivastava, N.; Susskind, J. M.; Zhang, J.; Salakhutdinov, R.; and Goh, H. 2021. Uncertainty Weighted Actor-Critic for Offline Reinforcement Learning. In *ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, 11319–11328. PMLR.

Xin, X.; Karatzoglou, A.; Arapakis, I.; and Jose, J. M. 2020. Self-Supervised Reinforcement Learning for Recommender Systems. In *SIGIR 2020, Virtual Event, China, July 25-30, 2020*, 931–940. ACM.

Yuan, F.; Karatzoglou, A.; Arapakis, I.; Jose, J. M.; and He, X. 2019. A Simple Convolutional Generative Network for Next Item Recommendation. In *WSDM 2019, Melbourne, VIC, Australia, February 11-15, 2019*, 582–590. ACM.

Zhao, X.; Gu, C.; Zhang, H.; Liu, X.; Yang, X.; and Tang, J. 2019. Deep Reinforcement Learning for Online Advertising in Recommender Systems. arXiv:1909.03602.

Zhao, X.; Zhang, L.; Ding, Z.; Xia, L.; Tang, J.; and Yin, D. 2018. Recommendations with Negative Feedback via Pair-wise Deep Reinforcement Learning. In *KDD 2018, London, UK, August 19-23, 2018*, 1040–1048. ACM.