

# Dynamic Knapsack Optimization Towards Efficient Multi-Channel Sequential Advertising

Xiaotian Hao<sup>\*1</sup> Zhaoqing Peng<sup>\*2</sup> Yi Ma<sup>\*1</sup> Guan Wang<sup>3</sup> Junqi Jin<sup>2</sup> Jianye Hao<sup>1</sup> Shan Chen<sup>2</sup>  
Rongquan Bai<sup>2</sup> Mingzhou Xie<sup>2</sup> Miao Xu<sup>2</sup> Zhenzhe Zheng<sup>4</sup> Chuan Yu<sup>2</sup> Han Li<sup>2</sup> Jian Xu<sup>2</sup> Kun Gai<sup>2</sup>

## Abstract

In E-commerce, advertising is essential for merchants to reach their target users. **The typical objective is to maximize the advertiser's cumulative revenue over a period of time under a budget constraint.** In real applications, an advertisement (ad) usually needs to be exposed to the same user multiple times until the user finally contributes revenue (e.g., places an order). However, existing advertising systems mainly focus on the immediate revenue with single ad exposures, ignoring the contribution of each exposure to the final conversion, thus usually falls into suboptimal solutions. In this paper, we **formulate the sequential advertising strategy optimization as a dynamic knapsack problem.** We **propose a theoretically guaranteed bilevel optimization framework**, which significantly reduces the solution space of the original optimization space while ensuring the solution quality. To improve the exploration efficiency of reinforcement learning, we also **devise an effective action space reduction approach.** Extensive offline and online experiments show the superior performance of our approaches over state-of-the-art baselines in terms of cumulative revenue.

## 1. Introduction

In E-commerce, online advertising plays an essential role for merchants to reach their target users, in which Real-time Bidding (RTB) (Zhang et al., 2014; 2016; Zhu et al., 2017) is an important mechanism. In RTB, each advertiser is al-

lowed to bid for every individual ad impression opportunity. Within a period of time, there are a number of impression opportunities (user requests) arriving sequentially. For each impression, each advertiser offers a bid based on the impression **value** (e.g., revenue) and competes with other bidders in real-time. The advertiser with the highest bid wins the auction and thus display ad and enjoys the impression value. Displaying an ad also associates with a **cost**: in Generalized Second-Price (GSP) Auction (Edelman et al., 2007), the winner is charged for fees according to the second highest bid. The typical advertising objective for an advertiser is to maximize its cumulative revenue of winning impressions over a time period under a fixed budget constraint.

In a digital age, to drive conversion, advertisers can reach and influence users across various channels such as display ad, social ad, paid search ad (Ren et al., 2018). As illustrated in Figure 9, the user's decision to convert (purchase a product) is usually driven by multiple interactions with ads. Each ad exposure would influence the user's preferences and interests, and therefore contributes to the final conversion. However, existing advertising systems (Yuan et al., 2013; Zhang et al., 2014; Ren et al., 2017; Zhu et al., 2017; Jin et al., 2018; Ren et al., 2019) mainly focus on maximizing the single-step revenue, while ignoring the contribution of previous exposure to the final conversion, and thus usually falls into suboptimal solutions. The reason is that simply optimizing the total immediate revenue cannot guarantee the maximization of long-term cumulative revenue. Besides, there exist some works (Boutilier & Lu, 2016; Du et al., 2017; Cai et al., 2017; Wu et al., 2018) which optimize the overall revenue under an extra-long (billions) request sequence using a single Constrained Markov Decision Process (CMDP) (Altman, 1999). However, the optimization of these methods above is myopic as they ignore the mental evolution of each user and long-term advertising effects. The learning is particularly inefficient as well.

Apart from the myopic approaches, there exists some literatures considering the long-term effect of each ad exposure. Multi-touch attribution (MTA) (Ji & Wang, 2017; Ren et al., 2018; Du et al., 2019) study the credits assignment to the previous ad displays before conversion. However, these

<sup>\*</sup>Equal contribution <sup>1</sup>College of Intelligence and Computing, Tianjin University, Tianjin, China <sup>2</sup>Alimama, Alibaba Group, Beijing, China <sup>3</sup>Department of Automation, Tsinghua University, Beijing, China <sup>4</sup>Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai, China. Correspondence to: Junqi Jin <junqi.jjq@alibaba-inc.com>, Jianye Hao <jianye.hao@tju.edu.cn>.

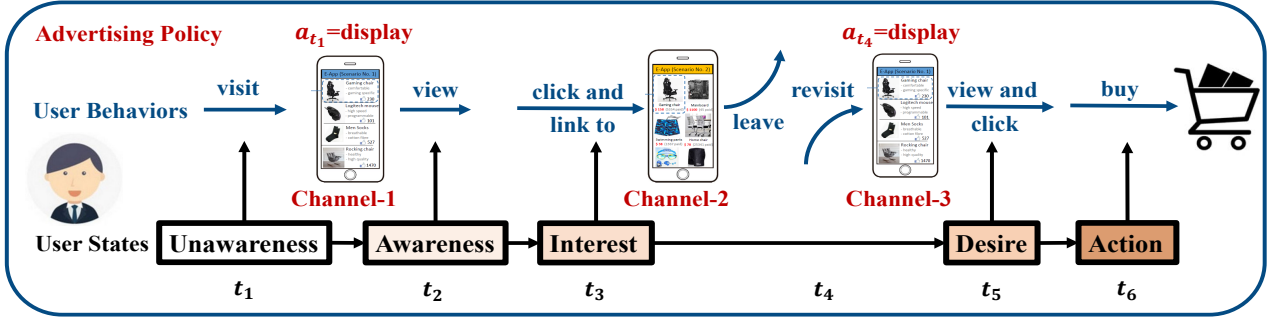


Figure 1. An illustration of the sequential multiple interactions (across different channels) between a user and an ad. Each ad exposure has long-term influence on the user’s final purchase decision.

methods only attend to figure out the contribution of each ad exposure, while not providing methods to optimize the strategies. Besides, since all media channels could affect users’ conversions, Li et al. (2018); Nuara et al. (2019) propose multi-channel budget allocation algorithms to help advertisers understand how particular channels contribute to user conversions. They optimize the budget allocation among all channels accordingly to maximize the overall revenue. However, the granularity of their optimizations is too coarse. They only optimize the budget allocation in the channel level and do not specifically optimize the advertising sequence for each user, which could lead to suboptimal overall performance.

Considering the shortcomings of existing works, we aim at optimizing the budget allocation of an advertiser among all users such that the cumulative revenue of the advertiser could be maximized, by explicitly taking into consideration the long-term influence of ad exposures to individual users. This problem consists of two levels of coupled optimization: **bidding strategy learning for each user** and **budget allocation among users, which we termed as Dynamic Knapsack Problem**. Different from traditional Knapsack problem, a number of challenges arise: 1) Given the estimated long-term value and cost for each user, the optimization space of the budget allocation grows exponentially in the number of users. Besides, since different advertising policies for each user will lead to different long-term values and costs, the overall optimization space is extremely large. 2) The long-term cumulative value and cost for each user are unknown, which are difficult to make accurate estimations.

To address the above challenges, we propose a novel bilevel optimization framework: **Multi-channel Sequential Budget Constrained Bidding (MSBCB)**, which transforms the original bilevel optimization problem into an equivalent two-level optimization with significantly reduced searching space. The higher-level only needs to optimize over one dimensional variable and the lower-level learns the optimal bidding policy for each user and computes the correspond-

ing optimal budget allocation solution. For the lower-level, we derive an optimal reward function with theoretical guarantee. Besides, we also propose an **action space reduction approach to significantly increase the learning efficiency of the lower-level**. Finally, extensive offline analyses and online A/B testing conducted on one of the world’s largest E-commerce platforms, Taobao, show the superior performance of our algorithm over state-of-the-art baselines.

## 2. Formulation: Dynamic Knapsack Problem

Within a time period of  $k$  days, we assume that there are  $N$  users  $\{i = 1, \dots, N\}$  visiting the E-commerce platform. Each user may interact with the app multiple times and trigger multiple advertising requests. During the sequential interactions between an ad and a user, each ad exposure could influence the user’s mind and therefore contributes to the final conversion. **Given a fixed ad, for each user  $i$** , we build a separate Markov Decision Process (MDP) (Sutton & Barto, 2018) to model its sequential interactions with the same ad. We use  $\pi_i$  to denote the advertising policy of the ad towards user  $i$ , which takes **user  $i$ ’s state as input and outputs the auction bid**. Details of the MDP will be discussed in Section 3.2. For the fixed ad, we define  $V_G(i|\pi_i)$  and  $V_C(i|\pi_i)$  as the expected long-term cumulative value and cost for each user  $i$  under policy  $\pi_i$ . Formally,

$$\begin{aligned} V_G(i|\pi_i) &= \mathbb{E}[G_i|\pi_i] = \mathbb{E}\left[\sum_{t=0}^{T_i} v_t|\pi_i\right] \\ V_C(i|\pi_i) &= \mathbb{E}[C_i|\pi_i] = \mathbb{E}\left[\sum_{t=0}^{T_i} c_t|\pi_i\right] \end{aligned} \quad (1)$$

where  $v_t$  and  $c_t$  represent the value (i.e., the revenue) and cost obtained from each request  $t$  according to policy  $\pi_i$ ,  $G_i = \sum_{t=0}^{T_i} v_t$  and  $C_i = \sum_{t=0}^{T_i} c_t$  represent the long-term cumulative value and cumulative cost,  $T_i$  is the length of the interaction sequence between user  $i$  and the current ad.

Given the above definitions, **for an advertiser**, our target is

to maximize its long-term cumulative revenue over  $k$  days under a budget constraint  $B$ , which is formulated as:

$$\begin{aligned} \max_{\Pi} \max_{\mathcal{X}} \sum_{i=1}^N x_i V_G(i|\pi_i) \\ \text{s.t.} \sum_{i=1}^N x_i V_C(i|\pi_i) \leq B \end{aligned} \quad (2)$$

where  $\Pi = \{\pi_1, \dots, \pi_N\}$ ,  $\mathcal{X} = \{x_1, \dots, x_N\}$ , and  $x_i \in \{0, 1\}$  indicates whether the user  $i$  is selected. Since whether displaying an ad to user  $i$  does not have any impact on user  $j$ 's behaviors,  $V_G(i|\pi_i)$ ,  $V_C(i|\pi_i)$  and  $\pi_i$  among different users are independent. Thus, given any fixed advertising policy  $\Pi = \{\pi_1, \dots, \pi_N\}$ ,  $V_G(i|\pi_i)$  and  $V_C(i|\pi_i)$  for each user  $i$  are fixed and the inner optimization of Equation (2) can be viewed as a classic knapsack problem. The items to be put into the knapsack is the users. However, different advertising policies would lead to different  $V_G(i|\pi_i)$ s and  $V_C(i|\pi_i)$ s for each user, thus here we define Equation (2) as a **Dynamic Knapsack Problem** where the value and cost of each item in the knapsack are dynamic. From the perspective of optimization, Formulation (2) is a typical bilevel optimization, where the optimization of  $\Pi$  is embedded (nested) within the optimization of  $\mathcal{X}$ . This bilevel optimization is challenging due to the following reasons:

- (1) The optimization space of the joint  $\Pi$  is continuous (for the bid space is continuous). The optimization space of  $\mathcal{X}$  is discrete, which grows exponentially in the number of users (hundreds of millions). Therefore, **the solution space of the combination of  $\Pi$  and  $\mathcal{X}$  is enormous** and thus is difficult or even impossible to optimize directly.
- (2) **The value of  $V_G(i|\pi_i)$  and  $V_C(i|\pi_i)$  are unknown and variable**, efficient approaches are required to estimate these values online under limited samples.

### 3. Methodology: MSBCB Framework

#### 3.1. Bilevel Decomposition and Proof of Correctness

Based on the above analysis, the bilevel optimization (2) is computationally prohibitive and cannot be solved directly. In this paper, we first decompose it into an equivalent two-level sequential optimization process. When taking a fixed policy  $\Pi$  as input, we denote the optimal solution of the degraded and static Knapsack Problem as  $K = \text{KP}(\Pi)$ . Further, the global optimal solution of Problem (2) could be defined as:

$$K^* = \max_{\pi_1, \pi_2, \dots, \pi_N} \text{KP}(\Pi) \quad (3)$$

where  $\pi_1, \dots, \pi_N$  are independent variables and  $K^*$  is the global optimal solution. To obtain  $K^*$ , we must firstly specify the form of the function  $\text{KP}(\Pi)$ .

When taking a fixed policy  $\Pi$  as input, computing  $\text{KP}(\Pi)$  is a classic static knapsack problem. However, another challenge in online advertising is that the user requests are arriving sequentially in real time and thus real-time decision makings are required. **Complicated algorithms (e.g. dynamic programming) are not applicable due to the incompleteness of all users values and costs.**

On the contrary, **the Greedy algorithm could compute a greedy solution without completely knowing the whole set of candidate users beforehand.** We will discuss this latter. Besides, the Greedy algorithm can achieve nearly optimal solution in the online advertising (Zhang et al., 2014; Wu et al., 2018). As proved by Dantzig (1957), if  $\forall i \in 1, \dots, N$ ,  $V_C(i|\pi_i) \leq (1 - \lambda)B$ ,  $0 \leq \lambda \leq 1$ , i.e., the cumulative cost for each user is much less than the budget, **the Greedy algorithm achieves an approximation ratio of  $\lambda$** , which means the **greedy solution is at least  $\lambda$  times of the optimal solution  $K$ .** The closer the  $\lambda$  gets to 1, the higher the quality of the greedy solution will be. **In online advertising,  $\lambda$  is usually greater than 99.9%.** Thus, the greedy solution is approximately optimal. We provide the detailed data and proof in Section B.1 of the Appendix. Therefore, in this paper, we refer to the Greedy algorithm, i.e.,  $\text{KP}(\Pi) \leftarrow \text{Greedy}(\Pi)$ .

We define  $\text{CPR}_i = \frac{V_G(i|\pi_i)}{V_C(i|\pi_i)}$  as the Cost-Performance Ratio of each user  $i$ . The greedy solution is computed by:

- (1) Sorting all users according to the Cost-Performance Ratio  $\text{CPR}_i$  in a descending order;
- (2) Pick users from top to bottom until the cumulative cost violates the budget constraint.

$V_G(i \pi_i)$	$V_C(i \pi_i)$	$\text{CPR}_i = V_G(i \pi_i)/V_C(i \pi_i)$
20	2	10
18	2	9
16	2	8
14	2	7
12	2	6
10	2	5
8	2	4

**Budget Constraint:  $B = 8$**       **Sorting in descending order**  
**CPR<sub>i</sub> threshold:  $\text{CPR}_{thr} = 7$**

Figure 2. The solution computing process of the Greedy algorithm.

An illustration is shown in Figure 2. In this example, the budget constraint  $B = 8$ . We denote the  $\text{CPR}_i$  of the last picked user as  $\text{CPR}_{thr}$ , the threshold of the cost-performance ratio. In this example, the  $\text{CPR}_{thr} = 7$ . The advantage is that the **Greedy algorithm only selects users whose  $\text{CPR}_i \geq \text{CPR}_{thr}$ .** **If we could estimate the  $\text{CPR}_{thr}$  beforehand, the Greedy algorithm could compute the solution online, without completely knowing the values and costs of all users.**

Now that  $\text{KP}(\Pi) \leftarrow \text{Greedy}(\Pi)$  and the Greedy algorithm

prefers users with larger  $\text{CPR}_i$  (only pick users whose  $\text{CPR}_i \geq \text{CPR}_{\text{thr}}$ ), according to Equation 3, to further improve the solution quality, an intuitive way is to optimize  $\pi_i$  for each user  $i$  such that each  $\text{CPR}_i$  could be maximized, i.e.,  $\pi_i' = \arg\max_{\pi_i} \text{CPR}_i$ . However, this intuition is incorrect. Maximizing the  $\text{CPR}_i$  of each user cannot guarantee that the greedy solution  $K = \text{Greedy}(\Pi)$  could be maximized. Next, we show that given all users' CPRs are maximized, we can still further improve the solution quality by increasing certain users' allocated budgets and decreasing their CPRs in exchange for greater overall cumulative value. Before we go into the details, we firstly give Lemma 1.

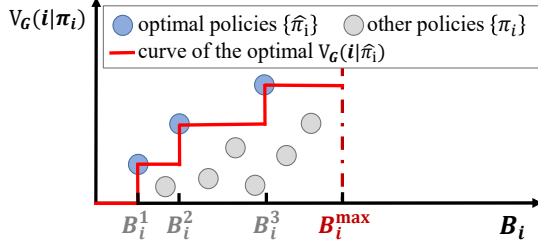


Figure 3.  $V_G(i|\hat{\pi}_i)$  is monotonic with  $V_C(i|\hat{\pi}_i)$ .

**Lemma 1.** For each user  $i$ , the cumulative value  $V_G(i|\hat{\pi}_i)$  increases monotonically with the increase of cost  $V_C(i|\hat{\pi}_i)$  within the range of all possible optimal policies  $\{\hat{\pi}_i\}$ .

*Proof.* We assume that the maximum budget allocated to each user  $i$  as  $B_i \in [0, B_i^{\max}]$ , where  $B_i^{\max}$  is the maximum cost user  $i$  can consume. Then, for each user  $i$ , within the current budget constraint  $B_i$ , the optimal advertising policy  $\hat{\pi}_i$  must be the one which could maximize the cumulative value, i.e.,  $\hat{\pi}_i = \arg\max_{\pi_i} V_G(i|\pi_i)$ , s.t.  $V_C(i|\pi_i) \leq B_i$ . Obviously, as  $B_i$  moves from 0 to  $B_i^{\max}$ , we will get a set of optimal policies  $\{\hat{\pi}_i\}$ , whose cost  $V_C(i|\hat{\pi}_i)$  and value  $V_G(i|\hat{\pi}_i)$  are both increasing. An illustration is shown in Figure 3. Thus we complete the proof.

As illustrated in Figure 4, each user's  $\text{CPR}_i$  (the width of each rectangular slice) is maximized initially. According to Lemma 1, for a user  $i$ , if we increase  $V_C(i|\pi_i)$  by  $\Delta V_C(i)$ , i.e., increase the height of user  $i$  by  $\Delta V_C(i)$ , the corresponding  $V_G(i|\pi_i)$  will also increase. We denote this increase in value as  $\Delta V_G(i)$ . Since there is a budget limit, a small increased height  $\Delta V_C(i)$  will squeeze out a small area nearby the  $\text{CPR}_{\text{thr}}$ , whose height is also  $\Delta V_C(i)$  and width is  $\text{CPR}_{\text{thr}}$ <sup>1</sup>. We denote the increased area by reshaping user  $i$  as  $\Delta V_G^+ = \Delta V_G(i)$  and the decreased area due to extrusion as  $\Delta V_G^- = \text{CPR}_{\text{thr}} * \Delta V_C(i)$ . Overall, if  $\Delta V_G^+ > \Delta V_G^-$ , the total area will be further increased. For

<sup>1</sup>Since  $\Delta V_C(j) \ll B$ , the area squeezed out could be considered as a tiny and smooth change and the width of the last user is approximately equal to  $\text{CPR}_{\text{thr}}$

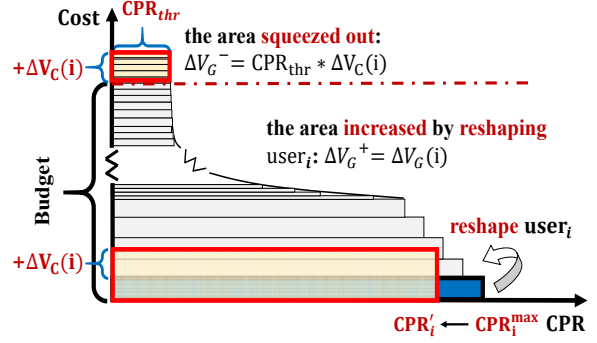


Figure 4. The x-axis denotes each user's  $\text{CPR}_i$  and y-axis denotes the cumulative cost of the Greedy algorithm. All users are sorted in descending order and arranged from bottom to top. Each rectangular slice's area (in gray) represents  $V_G(i|\pi_i) = \text{CPR}_i * V_C(i|\pi_i)$ , where  $\text{CPR}_i$  and  $V_C(i|\pi_i)$  are the width and height. Note that, the height of each rectangular slice is much less than the budget constraint, i.e.,  $V_C(i|\pi_i) \ll B$ . The red dashed line marks the position of the budget constraint. The total area of all rectangular slices under the red dashed line constitutes the greedy solution.

any user  $i$ ,  $\Delta V_G^+ > \Delta V_G^-$  yields:

$$\Delta V_G(i) > \text{CPR}_{\text{thr}} * \Delta V_C(i) \quad (4)$$

where  $\Delta V_G(i)$  and  $\Delta V_C(i)$  are caused by the change of  $\pi_i$ , e.g., from  $\pi_i'$  to  $\pi_i''$ . We denote  $\Delta V_G(i)$  as  $V_G(i|\pi_i'') - V_G(i|\pi_i')$  and  $\Delta V_C(i)$  as  $V_C(i|\pi_i'') - V_C(i|\pi_i')$ . We conclude that the greedy solution  $K = \text{Greedy}(\Pi')$  can be further improved if there exists any user  $i$  whose current policy  $\pi_i'$  can be further improved to  $\pi_i''$  such that  $\Delta V_G(i) > \text{CPR}_{\text{thr}} * \Delta V_C(i)$ . Otherwise, the current solution is optimal. Finally, we provide the definition of the optimal  $\pi_i^*$  in Theorem 1.

**Theorem 1.** Under the Greedy paradigm ( $K = \text{Greedy}(\Pi)$ ), for any given  $\text{CPR}_{\text{thr}}$ , the optimal advertising policy  $\pi_i^*$  for each user  $i$  is the one which could maximize  $V_G(i|\pi_i) - \text{CPR}_{\text{thr}} * V_C(i|\pi_i)$ . In other words,  $\pi_i^*$  is defined as:

$$\pi_i^* = \arg\max_{\pi_i} [V_G(i|\pi_i) - \text{CPR}_{\text{thr}} * V_C(i|\pi_i)] \quad (5)$$

We denote  $\Pi^* = \{\pi_1^*, \dots, \pi_N^*\}$ . The corresponding solution  $K_{\text{greedy}}^* = \text{Greedy}(\Pi^*)$  is the optimal Greedy solution of the Dynamic Knapsack Problem defined in Equation (2).

**Proof of Theorem 1.** We define  $\Pi^* = \{\pi_1^*, \dots, \pi_N^*\}$ , where  $\pi_i^*$  is defined according to Equation (5),  $\forall i \in \{1, \dots, N\}$ . We prove Theorem 1 by contradiction. Given the threshold  $\text{CPR}_{\text{thr}}$ , we firstly assume that  $\text{Greedy}(\Pi^*)$  is not the optimal greedy solution of the Dynamic Knapsack Problem, which means we could at least find a user  $i$ , whose policy  $\pi_i^*$  could be further improved to policy  $\pi_i''$  such that the overall area is increased. This means we could find a better policy  $\pi_i''$  for user  $i$  such that  $\Delta V_G(i) > \text{CPR}_{\text{thr}} * \Delta V_C(i)$  according to Equation (4), where  $\Delta V_G(i) = V_G(i|\pi_i'') - V_G(i|\pi_i^*)$



and  $\Delta V_C(i) = V_C(i|\pi_i'') - V_C(i|\pi_i^*)$  ( $V_G(i|\hat{\pi}_i)$  increases monotonically with the increase of  $V_C(i|\hat{\pi}_i)$  according to Lemma 1). Further,  $\Delta V_G(i) > \text{CPR}_{\text{thr}} * \Delta V_C(i)$  yields:

$$\begin{aligned} [V_G(i|\pi_i'') - \text{CPR}_{\text{thr}} * V_C(i|\pi_i'')] > \\ [V_G(i|\pi_i^*) - \text{CPR}_{\text{thr}} * V_C(i|\pi_i^*)] \end{aligned} \quad (6)$$

Equation (6) indicates that

$$\pi_i^* \neq \underset{\pi_i}{\operatorname{argmax}} [V_G(i|\pi_i) - \text{CPR}_{\text{thr}} * V_C(i|\pi_i)]$$

which contradicts the definition of  $\pi_i^*$  in Equation (5). Thus, the theorem statement is obtained.

---

**Algorithm 1** MSBCB Framework.
 

---

- 1: **Input:** an initial  $\text{CPR}_{\text{thr}}$ ;
  - 2: **Output:** optimal greedy solution of the Dynamic Knapsack Problem;
  - 3: **for** each period until convergence **do**
  - 4: Taking the current estimated  $\text{CPR}_{\text{thr}}$  as input, the agent optimizes the advertising policy  $\pi_i$  for each user  $i$  according to Section 3.2 and acquires the optimal  $\Pi^* = \{\pi_1^*, \dots, \pi_N^*\}$ .
  - 5: Based on the current estimated  $\text{CPR}_{\text{thr}}$  and the obtained  $\Pi^*$ , the agent calculates the greedy solution according to Section 3.3 and collects the actual feedback cost and the predefined budget.
  - 6: Update the estimated  $\text{CPR}_{\text{thr}}$  towards  $\text{CPR}_{\text{thr}}^*$  by minimizing the gap between the actual feedback cost and the budget according to Section 3.4.
  - 7: **end for**
- 

We present the overall MSBCB framework in Algorithm 1, which involves a two-level sequential optimization process. **(1) Lower-level:** Given any  $\text{CPR}_{\text{thr}}$ , we could obtain the optimal advertising policy  $\Pi^*$  following Equation 5 of Theorem 1, which will be discussed in Section 3.2. Then, based on  $\text{CPR}_{\text{thr}}$  and the optimized  $\Pi^*$ , we could acquire the Greedy solution by selecting users whose  $\text{CPR}_i \geq \text{CPR}_{\text{thr}}$ , which will be detailed in Section 3.3. **(2) Higher-level:** However, the current  $\text{CPR}_{\text{thr}}$  might  $\neq \text{CPR}_{\text{thr}}^*$ , which means selecting all users whose  $\text{CPR}_i \geq \text{CPR}_{\text{thr}}$  might violate the budget constraint or lead to a substantial budget surplus. Thus, we optimize the current  $\text{CPR}_{\text{thr}}$  towards  $\text{CPR}_{\text{thr}}^*$  in Section 3.4. Overall, the optimization space of  $\mathcal{X}$  is reduced from  $2^N$  to a one-dimensional continuous variable  $\text{CPR}_{\text{thr}}$ . We conclude that Algorithm 1 could iteratively converge to a unique and approximate optimal solution. We present the proof of convergence in Section B.3 of the Appendix.

### 3.2. Lower-level Advertising Policy Optimization with Reinforcement Learning

Given a threshold  $\text{CPR}_{\text{thr}}$  as input, we aim to acquire the optimal advertising policy  $\pi_i^*$  defined in Equation (5) of

Theorem 1. Combining the definitions of  $V_G(i|\pi_i)$  and  $V_C(i|\pi_i)$  with Equation (5), we have

$$\begin{aligned} \pi_i^* &= \underset{\pi_i}{\operatorname{argmax}} [V_G(i|\pi_i) - \text{CPR}_{\text{thr}} * V_C(i|\pi_i)] \\ &= \underset{\pi_i}{\operatorname{argmax}} (\mathbb{E}[G_i|\pi_i] - \text{CPR}_{\text{thr}} * \mathbb{E}[C_i|\pi_i]) \\ &= \underset{\pi_i}{\operatorname{argmax}} \mathbb{E}[(G_i - \text{CPR}_{\text{thr}} * C_i) | \pi_i] \\ &= \underset{\pi_i}{\operatorname{argmax}} \mathbb{E}\left[\sum_{t=0}^{T_i} (v_t - \text{CPR}_{\text{thr}} * c_t) | \pi_i\right] \end{aligned} \quad (7)$$

Accordingly, we define  $r_t = v_t - \text{CPR}_{\text{thr}} * c_t$ , i.e., value –  $\text{CPR}_{\text{thr}} * \text{cost}$ , as the immediate profit acquired at each step  $t$ . The objective of Equation (7) is to obtain the optimal advertising policy  $\pi_i^*$  which could maximize the expected long-term cumulative profit. To solve this sequential decision making problem, we formulate it as an MDP and use Reinforcement Learning (RL) (Sutton & Barto, 2018) techniques to acquire the optimal policy  $\pi_i^*$ .

We consider an episodic MDP, where an episode starts with the first interaction between a user and an ad, and ends up with a purchase or exceeding the maximum step  $T_i$  as:

- **State  $\mathcal{S}$ :** The state  $s_t$  should in principle reflect the user request status, ad info, user-ad interaction history info and the RTB environment.
- **Action  $\mathcal{A}$ :** The action each agent can take in the RTB platform is the bid, which is a real number between 0 and the upper bound  $\text{bid}_{\text{max}}$ , i.e.,  $a_t \in [0, \text{bid}_{\text{max}}]$ .
- **Reward  $\mathcal{R}(\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R})$ :** The immediate reward at step  $t$  is defined as  $r_t = v_t - \text{CPR}_{\text{thr}} * c_t$ .
- **Transition probability  $\mathcal{P}(\mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1])$ :** Transition probability is defined as the probability of state transitioning from  $s_t$  to  $s_{t+1}$  when taking action  $a_t$ .
- **Discount factor  $\gamma$ :** The bidding agent aims to maximize the total discounted reward  $R_j = \sum_{k=t}^{T_i} \gamma r_k$  from step  $t$  onwards, where  $\gamma \in [0, 1]$ .

For each user  $i$ , we define the state-action value function  $Q(s, a) = \mathbb{E}[R_i | s, a, \pi_i]$  as the expected cumulative reward achieved by following the advertising policy  $\pi_i$ . The MDP can be solved using existing Deep Reinforcement Learning (DRL) algorithms such as DQN (Mnih et al., 2013), DDPG (Lillicrap et al., 2015) and PPO (Schulman et al., 2017). After sufficient training, we would acquire the optimized advertising policies  $\Pi^* = \{\pi_1^*, \dots, \pi_N^*\}$  for all users.

### 3.3. Lower-level User Selection by Greedy Algorithm

Taking the current  $\text{CPR}_{\text{thr}}$  and the optimized advertising policies  $\Pi^* = \{\pi_1^*, \dots, \pi_N^*\}$  as inputs, we aim to obtain the

greedy solution of the Dynamic Knapsack Problem. In reality, we cannot know all users' request sequences and their values and costs beforehand because the user requests are arriving sequentially in real time. Thus, many complicated methods depending on the completeness of all users' data, e.g., the dynamic programming approach (Martello et al., 1999), are not applicable. Even the traditional Greedy algorithm cannot be applied either. Fortunately, **the greedy solution could be computed online** in an easy way: given the threshold  $\text{CPR}_{\text{thr}}$ , the agent only has to select users online whose CPRs are greater than the threshold (an illustration is shown in Figure 2). Therefore, we only have to estimate the  $\text{CPR}_i = \frac{V_G(i|\pi_i)}{V_C(i|\pi_i)}$  for each user  $i$ . To acquire  $V_G(i|\pi_i)$  and  $V_C(i|\pi_i)$ , besides  $Q(s, a)$ , **we also maintain two other state value functions  $V_G(s)$  and  $V_C(s)$  according to the Bellman Equation** (Sutton & Barto, 2018), where  $V_G(s) = \mathbb{E}[G_j|s, \pi_j]$  and  $V_C(s) = \mathbb{E}[C_j|s, \pi_j]$ .

### 3.4. Higher-level Optimization by Feedback Control

However, the current estimated threshold  $\text{CPR}_{\text{thr}}$  might have some bias from the optimal  $\text{CPR}_{\text{thr}}^*$ . Thus, selecting all users whose  $\text{CPR}_i \geq \text{CPR}_{\text{thr}}$  might violate the budget constraint or lead to a substantial budget surplus. Only when the estimated  $\text{CPR}_{\text{thr}}$  is exactly the same with the optimal  $\text{CPR}_{\text{thr}}^*$ , the actual total advertising cost will be equal to the budget. To achieve this, we design a feedback control mechanism, i.e., a **PID controller** (Åström & Hägglund, 1995), to dynamically adjust the  $\text{CPR}_{\text{thr}}$  towards  $\text{CPR}_{\text{thr}}^*$  according to actual feedback of the overall cost. The core formula is:

$$\text{CPR}_{\text{thr}}^* = \left[ 1 + \alpha_1 \left( \frac{\text{cost}_t}{B} - 1 \right) + \alpha_2 \left( \frac{\text{cost}_{t-n:t}}{n*B} - 1 \right) \right] \quad (8)$$

where  $\text{cost}_t$  is the actual feedback cost of the current period,  $B$  is the budget,  $\text{cost}_{t-n:t}$  and  $n*B$  are the overall cost and the overall budget of the most recent  $n$  periods.  $\alpha_1$  and  $\alpha_2$  are two learning rates. The main idea is when the actual cost exceeds (is less than) the budget, the threshold  $\text{CPR}_{\text{thr}}$  will be increased (decreased) accordingly such that less (more) users will be selected, which will reduce (increase) the cost in turn. The first term  $\alpha_1 \left( \frac{\text{cost}_t}{B} - 1 \right)$  is designed to keep up with the latest changes. The second term  $\alpha_2 \left( \frac{\text{cost}_{t-n:t}}{n*B} - 1 \right)$  is designed to stabilize learning.

### 3.5. Action Space Reduction for RL in Advertising

However, when applying the RL approaches mentioned in Section 3.2 to online advertising, one typical issue is that **the sample utilization is inefficient**. The main reason is that the action space of the agent is continuous, thus the range of  $[0, \text{bid}_{\text{max}}]$  needs to be fully explored in all states. To resolve this problem, we reduce the magnitude of the continuous action space (i.e.,  $a_t \in [0, \text{bid}_{\text{max}}]$ ) to a binary one (i.e.,  $\hat{a}_t \in \{0, 1\}$ ) by making full use of the prior knowledge in advertising, which greatly improves the sample utilization

of the RL approaches. Specifically, since different bids  $a_t$  can only result in two different outcomes  $\hat{a}_t \in \{0, 1\}$ , where  $\hat{a}_t = 1$  or 0 indicates whether the ad is displayed to the user, we only have to evaluate the different expected returns resulted by  $\hat{a}_t = 1$  or  $\hat{a}_t = 0$  for  $Q(s, a)$ . We denote the greedy action  $\hat{a}_t^*$  based on the current value estimations as:

$$\hat{a}_t^* = \begin{cases} 1 & \text{if } Q(s, \hat{a}_t = 1) > Q(s, \hat{a}_t = 0) \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

Then, to obtain an executable bid, for  $\hat{a}_t^* = 0$ , we could offer a low enough bid, e.g.,  $a_t = 0$ , to make sure that it is impossible to win the auction. **For  $\hat{a}_t^* = 1$ , we propose an optimal bid function which could output a bid greater than the second highest bid while not overbidding.**

In detail, we maintain two state-action value functions  $Q_G(s, \hat{a}_t) = \mathbb{E}[G_i|s, \hat{a}_t, \pi_i]$  and  $Q_C(s, \hat{a}_t) = \mathbb{E}[C_i|s, \hat{a}_t, \pi_i]$ . Since the reward function is defined as  $r_t = v_t - \text{CPR}_{\text{thr}} * c_t$ , we have  $Q(s, \hat{a}_t) = Q_G(s, \hat{a}_t) - \text{CPR}_{\text{thr}} * Q_C(s, \hat{a}_t)$ . Then  $Q(s, \hat{a}_t = 1) > Q(s, \hat{a}_t = 0)$  yields:

$$\begin{aligned} [Q_G(s, \hat{a}_t = 1) - \text{CPR}_{\text{thr}} * Q_C(s, \hat{a}_t = 1)] &> \\ [Q_G(s, \hat{a}_t = 0) - \text{CPR}_{\text{thr}} * Q_C(s, \hat{a}_t = 0)] \end{aligned} \quad (10)$$

If  $\hat{a}_t = 0$ , the expected immediate cost is 0 (since the ad is not exposed). If  $\hat{a}_t = 1$ , we denote the expected immediate cost as  $\mathbb{E}[c_t|\hat{a}_t = 1]$ , whose value depends on the pricing model. In online advertising, typical pricing models includes **CPM** (Cost Per Mille, the advertiser bid for impressions and is charged based on impressions), **CPC** (Cost Per Click, the advertiser bid for clicks and is charged based on clicks) and **CPS** (Cost Per Sales, the advertiser bid for conversions and is charged based on conversions). If CPM is used,  $\mathbb{E}[c_t|\hat{a}_t = 1] = \text{bid}_t^{\text{2nd}}$ , where  $\text{bid}_t^{\text{2nd}}$  denotes the second highest bid in the auction. If CPC is used,  $\mathbb{E}[c_t|\hat{a}_t = 1] = \text{bid}_t^{\text{2nd}} * \text{pCTR}$ , where pCTR represents the predicted Click-Through Rate. If CPS is used,  $\mathbb{E}[c_t|\hat{a}_t = 1] = \text{bid}_t^{\text{2nd}} * \text{pCTR} * \text{pCVR}$ , where pCVR represents the predicted Conversion Rate. For ease of presentation, we take CPM for an example. Under CPM,

$$\begin{aligned} Q_C(s, \hat{a}_t = 1) &= \mathbb{E}[c_t + \sum_{k=t+1}^{T_i} c_k | s, \hat{a}_t = 1, \pi_i] \\ &= \text{bid}_t^{\text{2nd}} + \mathbb{E}[\sum_{k=t+1}^{T_i} c_k | s, \hat{a}_t = 1, \pi_i] \\ Q_C(s, \hat{a}_t = 0) &= \mathbf{0} + \mathbb{E}[\sum_{k=t+1}^{T_i} c_k | s, \hat{a}_t = 0, \pi_i] \end{aligned} \quad (11)$$

Notice that the second highest bid  $\text{bid}_t^{\text{2nd}}$  is unknown until the current auction is finished. **Substituting Equation (11)**

into Equation (10), we acquire

$$\text{bid}_t^{\text{2nd}} < \left[ \left( \frac{Q_G(s, \hat{a}_t = 1)}{\text{CPR}_{\text{thr}}} - Q_C^{\text{next}}(s, \hat{a}_t = 1) \right) - \left( \frac{Q_G(s, \hat{a}_t = 0)}{\text{CPR}_{\text{thr}}} - Q_C^{\text{next}}(s, \hat{a}_t = 0) \right) \right] \quad (12)$$

where  $Q_C^{\text{next}}(s, \hat{a}_t) = \mathbb{E}[\sum_{k=t+1}^{T_j} c_k | s, \hat{a}_t, \pi_i]$ . We denote the term on the right of the ' $<$ ' in Equation (12) as  $\mathbf{b}_t^*$ . And we conclude that the bidding agent can always set the bid price  $a_t = \mathbf{b}_t^*$  during the online bidding phase, which is the optimal action without any loss of accuracy. Refer to Section B.2 of the Appendix for proof. For CPC or CPS, the optimal bid formula  $\mathbf{b}_t^*$  can be easily acquired by substituting the corresponding  $\mathbb{E}[c_t | \hat{a}_t = 1]$  into Equation 11. Here, we reaffirm that our action space reduction technique is a generalized design and is applicable to different pricing models.

## 4. Empirical Evaluation: Simulations

We start with designing simulation experiments to shed light on the contributions of the proposed framework MSBCB under more controlled settings. Similar to the simulation settings of (Ie et al., 2019), we assume there are a set of users  $\{i = 1, \dots, N\}$ , a set of ads  $\mathcal{D}$  and a set of commodity categories  $\mathcal{T}$ . Each ad  $d \in \mathcal{D}$  has an associated category. Each user  $i$  has various degrees of interests in commodity categories, which is influenced by the displayed ad. When user  $i$  consumes ad  $d$ , his interest in category  $T(d)$  is nudged stochastically, biased slightly towards increasing his interest, but allows some chance of decreasing his interest. We set  $N = 10000$ ,  $|\mathcal{D}| = 2000$  and  $|\mathcal{T}| = 20$  in the following experiments. Detailed settings of the simulation environment can be found in Section D of the Appendix.

### 4.1. Baselines

We compare our MCBCB with following baseline strategies:

- **Myopic Approaches:** (1) Manual Bid is a strategy that the agent continuously bids at the same price initialized by the advertiser. (2) Contextual Bandit (Zhang et al., 2014) aims at maximizing the accumulated short-term value of each request based on the Greedy framework.
- **Greedy with maximized CPR:** This approach is similar to our method under the Greedy framework except that each  $\pi_i$  is optimized by maximizing the long-term CPR. In the offline simulation, we enumerate all policies for each user and select the one which could maximize its CPR. This approach is named as Greedy+maxCPR.
- **Greedy with state-of-the-art RL approaches:** These baselines, i.e., Greedy+DQN, Greedy+DDPG and

Greedy+PPO, utilize the same reward function with our MSBCB to optimize the lower-level optimization of II. The difference is that our MSBCB leverages the action space reduction technique. For DQN and PPO, we discretize the bid action space  $[0, \text{bid}_{\text{max}}]$  evenly into 11 real numbers as the valid actions.

- **Undecomposed Optimization:** These baselines are RL approaches (DQN, DDPG and PPO) based on the Constrained Markov Decision Process (CMDP). They are named as Constrained+DQN, Constrained+DDPG, Constrained+PPO respectively. We follow the CMDP design and settings in (Wu et al., 2018).
- **Offline Optimal:** The optimal solution of the Dynamic Knapsack Problem can be computed by dynamic programming in offline simulation because we could enumerate all possible policies to get the corresponding long-term values and costs for each user. Note that since users' request sequences are unknown beforehand and there is only one chance for the ad to bid for each request in the online advertising systems, the optimal solution can only be obtained in offline simulation.

### 4.2. Experimental Results

We conduct extensive analysis of our MSBCB in the following 5 aspects. All approaches aim to maximize the advertiser's cumulative revenue under a fixed budget constraint. All experimental results are averaged over 10 runs. The hyperparameters for each algorithm are set to the best we found after grid-search optimization.

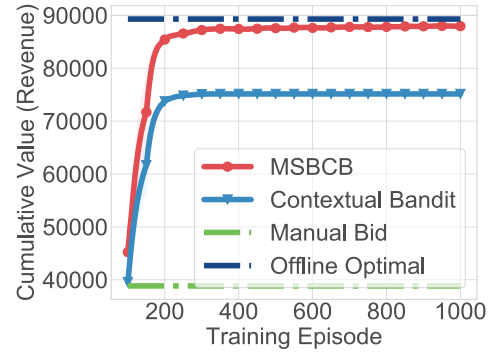


Figure 5. Values comparisons (learning curves) of the myopic approaches with non-myopic approaches and the offline optimal.

**Myopic vs Non-myopic.** To show the benefits of upgrading the myopic advertising system into a farsighted one, we compare the cumulative revenue achieved by our MSBCB with two other myopic baselines. The learning curves and results are shown in Figure 5 and Table 1. We see that MSBCB outperforms the Manual Bid and the Contextual Bandit by a large margin, which indicates that taking account of

the long-term effect of each ad exposure could significantly improve the cumulative advertising results.

**MSBCB vs the Offline Optimal.** In Figure 5, we also compare our MSBCB with the Offline Optimal, which is computed by a modified dynamic programming algorithm. We see that as the training continues, our MSBCB gradually achieves an approximately optimal solution. Detailed results are summarized in Table 1. Our MSBCB empirically achieves an approximation ratio of 98.53% ( $\pm 0.36\%$ ).

**MSBCB vs Greedy with maximized CPR.** As discussed in Section 3.1, under the Greedy framework, maximizing each user’s  $CPR_i$  cannot guarantee that the greedy solution of the Dynamic Knapsack Problem (2) could be maximized. The optimal advertising policy  $\pi_i$  for each user is given by Theorem 1. To experimentally verify the correctness of Theorem 1, we compare the cumulative revenue achieved by MSBCB and the Greedy with maximized CPR. As shown in Figure 6 and Table 1, MSBCB outperforms Greedy with maximized CPR and achieves a +5.11% improvement.

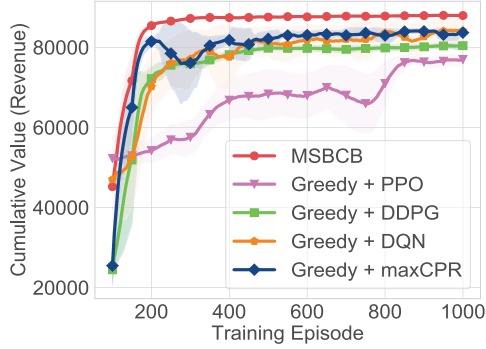


Figure 6. Value comparisons of MSBCB with the Greedy with maximized CPR and the Greedy with state-of-the-art RL.

**MSBCB vs Greedy with state-of-the-art RL approaches.** Besides, to show the effectiveness of the action-space reduction proposed in Section 3.5, we compare MSBCB with the state-of-the-art DRL approaches under the Greedy framework. As shown in Figure 6 and Table 1, MSBCB outperforms Greedy+DQN, Greedy+DDPG and Greedy+PPO both in the cumulative revenue and the **convergence speed, which shows that the action space reduction effectively improves the sample efficiency of RL approaches.**

**Decomposed MSBCB vs Undecomposed optimization.** Similar to (Wu et al., 2018), the undecomposed optimization baselines consider all users requests as a whole and model the budget allocations among all request as a CMDP. As shown in Figure 7 and Table 1, MSBCB outperforms the CMDP based RL approaches by a large margin. The reason of the poor performance in CMDP-based approaches is that these methods model all users’ requests as a whole sequence and thus the learning process is particularly inefficient. In

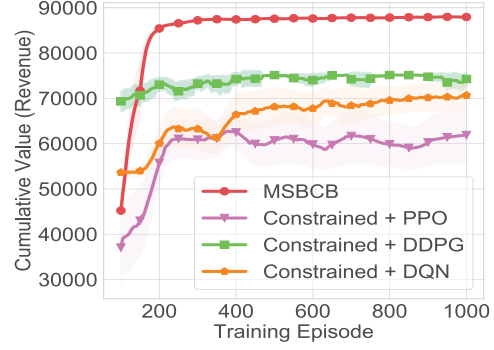


Figure 7. Values comparison (learning curves) of MSBCB and state-of-the-art CMDP based RL approaches.

contrast, our MSBCB decomposes the whole sequence optimization into an efficient two-level optimization process, thus can achieve better performance more easily.

Table 1. Cumulative values, costs, value improvements (over Contextual Bandit) and the approximation ratio of all approaches.

Method	Revenue	Cost	Revenue Impro	Approximation Ratio
Manual Bid	38838.28	11995.10	-48.31%	43.5%
Contextual Bandit	75137.30	11995.46	0%	84.15%
Constrained + PPO	61890.92	11954.07	-17.63 $\pm$ 16.11%	69.31 $\pm$ 13.56%
Constrained + DDPG	74259.12	11996.12	-1.19 $\pm$ 3.66%	83.17 $\pm$ 3.08%
Constrained + DQN	70662.65	11881.12	-5.96 $\pm$ 7.83%	79.14 $\pm$ 6.59%
Greedy + maxCPR	83668.70	11914.12	11.35 $\pm$ 2.84%	93.70 $\pm$ 2.36%
Greedy + PPO	76970.35	11825.59	2.44 $\pm$ 3.52%	86.20 $\pm$ 2.93%
Greedy + DDPG	80424.69	11841.28	7.04 $\pm$ 1.13%	90.07 $\pm$ 0.92%
Greedy + DQN	84117.09	11794.24	11.95 $\pm$ 4.96%	94.21 $\pm$ 4.14%
<b>MSBCB</b>	<b>87947.99</b>	<b>11957.57</b>	<b>17.95<math>\pm</math>0.42%</b>	<b>98.50<math>\pm</math>0.33%</b>
<b>MSBCB (enum)</b>	<b>89251.77</b>	<b>11988.36</b>	<b>18.78%</b>	<b>99.96%</b>
Offline Optimal	89291.11	11999.23	18.84%	100.00%

The complete comparisons of all approaches are shown in Table 1. The budget constraint  $B$  is set to 12000 for all experiments. In Table 1, we also add an MSBCB (enum), which is the theoretical upper bound of our MSBCB. The difference between MSBCB (enum) and MSBCB is that: the MSBCB (enum) computes the optimal advertising policy  $\pi_i^*$  for each user  $i$  by enumerating all possible policies. Instead of utilizing the RL approach, MSBCB (enum) could find the one which maximizes  $V_G(i|\pi_i) - CPR_{thr} * V_C(i|\pi_i)$ . We see MSBCB (enum) is very close to the optimal solution and reaches an approximation ratio of 99.96%.

#### 4.3. Effectiveness of Action Space Reduction

As shown in Table 2, MSBCB achieves a revenue of 75000 in only 61 epochs, reducing more than 60% samples compared with the state-of-the-art RL baselines without using the action-space reduction technique. As for learning process, our MSBCB achieves the same revenue (80000) more than 10 times faster than the baselines, reducing more than 90% samples and finally reaches the highest revenue. Thus,



with the action space reduction technique, our **MSBCB** could reach a higher performance with a faster speed and significantly improve the sample efficiency. More analysis of our MSBCB, e.g., the convergence of  $\Pi^*$  and  $\text{CPR}_{\text{thr}}^*$ , and the hyperparameter settings of the offline experiments are shown in Section D of the Appendix.

Table 2. The training epochs and the number of samples needed by different approaches when achieving the same revenue level.

Revenue Method	75000		80000		85000	
	#Epoch	#Samples	#Epoch	#Samples	#Epoch	#Samples
Greedy+PPO	817	4183040	-	-	-	-
Greedy+DDPG	154	788480	853	4362240	-	-
Greedy+DQN	373	1909760	754	3855360	-	-
<b>MSBCB</b>	<b>61</b>	<b>312320</b>	<b>71</b>	<b>363520</b>	<b>104</b>	<b>532480</b>

## 5. Empirical Evaluation: Online A/B Testing

We deployed MSBCB on one of the world’s largest E-commerce platforms, Taobao. Our platform is authorized by the advertisers to dynamically adjust their bid prices for each user request according its value in the real-time auction. In the online experiments, we compare MSBCB with two models widely used in the industry.

- Cross Entropy Method (CEM), which is a deployed production model, whose target is to optimize the immediate rewards. We consider CEM as the control group in the following evaluations.
- Contextual Bandit, which has been explained in previous section and is reserved as a contrast test.

The experiment involves 135,858,118 users and 72,147 ad items from 186 advertisers. For fair comparison, we control the consumers and the advertisers involved in the A/B testing to be homogeneous. In detail, the 135,858,118 users are randomly and evenly divided into 3 groups. For users in group #1, all 186 advertisers adopt the CEM algorithm. For users in group #2, all 186 advertisers adopt the Contextual Bandit algorithm. For users in group #3, all 186 advertisers adopt our MSBCB. Table 3 summarises the effects of the Contextual Bandit and our MSBCB compared to the Cross Entropy Method from Dec.10 to Dec.20 in 2019. From Table 3, we see that our MSBCB achieves a +10.08% improvement in revenue and a +10.31% improvement in ROI with almost the same cost (-0.20%). The results indicate that upgrading the myopic advertising strategy into a farsighted one could significantly improves the cumulative revenue. Besides, as shown in Figure 8, the daily ROI improvement also demonstrates the effectiveness of our MSBCB compared with the Contextual Bandit.

Given that there are only 186 advertisers take part in our online experiment, one frequently asked question is “How

Table 3. The overall performance comparisons of the A/B testing. CVR represents the Conversion Rate of the users. #PV represents the number of page views.  $\text{ROI} = \frac{\text{Revenue}}{\text{Cost}}$  means Return On Investment. (Notice that CEM is the control group and the improvements of Contextual Bandit and MSBCB are compared over CEM.)

Method	Revenue	Cost	CVR	#PV	ROI
Contextual Bandit	+0.91%	-3.26%	+4.78%	+4.62%	+4.31%
<b>MSBCB</b>	<b>+10.08%</b>	<b>-0.20%</b>	<b>+6.04%</b>	<b>+15.37%</b>	<b>+10.31%</b>

**does the MSBCB work across all ads?”** Since 186 is relatively small compared with the total number of advertisers, their policy updates would not cause dramatic changes to the RTB environment. In other words, the RTB environment is still approximately stationary from a single-ad perspective. This setting also works well with our practical business model-providing better service for VIP advertisers (about 0.2% of all the advertisers). In the case that the majority of the advertisers adopt MSBCB, the system cannot be estimated as being stationary from any single-ads perspective and explicit multi-agent modeling and coordination should be incorporated. Detailed analysis of the improvement in revenue for each advertiser is presented in Table 7 and Figure 19 of the Appendix. More details about the deployment and experimental results (e.g., the online model architecture) can also be found in Section C and E of the Appendix.

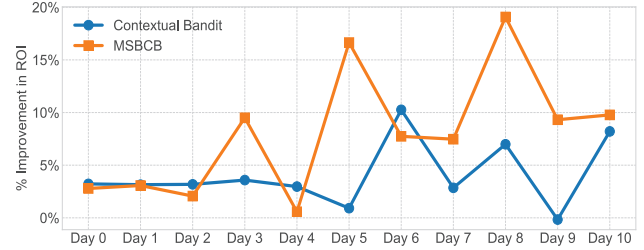


Figure 8. Daily ROI improvement comparisons of Contextual Bandit and MSBCB over Cross Entropy Method.

## 6. Conclusion

We formulate the multi-channel sequential advertising problem as a Dynamic Knapsack Problem, whose target is to maximize the long-term cumulative revenue over a period of time under a budget constraint. We decompose the original problem into an easier bilevel optimization, which significantly reduces the solution space. For the lower-level optimization, we derive an optimal reward function with theoretical guarantees and design an action space reduction technique to improve the sample efficiency. Extensive offline experimental analysis and online A/B testing demonstrate the superior performance of our MSBCB over the state-of-the-art baselines in terms of cumulative revenue.

## Acknowledgements

The work is supported by the National Natural Science Foundation of China (Grant Nos.: 61702362, U1836214), the Special Program of Artificial Intelligence and the Special Program of Artificial Intelligence of Tianjin Municipal Science and Technology Commission (No.: 569 17ZXRGX00150) and the Alibaba Group through Alibaba Innovative Research Program. We deeply appreciate all teammates from Alibaba group for the significant supports for the online experiments.

## References

- Altman, E. *Constrained Markov decision processes*, volume 7. CRC Press, 1999.
- Åström, K. J. and Häggglund, T. *PID controllers: theory, design, and tuning*, volume 2. Instrument society of America Research Triangle Park, NC, 1995.
- Boutillier, C. and Lu, T. Budget allocation using weakly coupled, constrained markov decision processes. 2016.
- Cai, H., Ren, K., Zhang, W., Malialis, K., Wang, J., Yu, Y., and Guo, D. Real-time bidding by reinforcement learning in display advertising. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*, pp. 661–670. ACM, 2017.
- Dantzig, G. B. Discrete-variable extremum problems. *Operations research*, 5(2):266–288, 1957.
- Du, M., Sassioui, R., Varisteas, G., Brorsson, M., Cherkaoui, O., et al. Improving real-time bidding using a constrained markov decision process. In *International Conference on Advanced Data Mining and Applications*, pp. 711–726. Springer, 2017.
- Du, R., Zhong, Y., Nair, H., Cui, B., and Shou, R. Causally driven incremental multi touch attribution using a recurrent neural network. *arXiv preprint arXiv:1902.00215*, 2019.
- Edelman, B., Ostrovsky, M., and Schwarz, M. Internet advertising and the generalized second-price auction: Selling billions of dollars worth of keywords. *American economic review*, 97(1):242–259, 2007.
- Ie, E., Jain, V., Wang, J., Navrekar, S., Agarwal, R., Wu, R., Cheng, H.-T., Lustman, M., Gatto, V., Covington, P., et al. Reinforcement learning for slate-based recommender systems: A tractable decomposition and practical methodology. *arXiv preprint arXiv:1905.12767*, 2019.
- Ji, W. and Wang, X. Additional multi-touch attribution for online advertising. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- Jin, J., Song, C., Li, H., Gai, K., Wang, J., and Zhang, W. Real-time bidding with multi-agent reinforcement learning in display advertising. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pp. 2193–2201. ACM, 2018.
- Li, P., Hawbani, A., et al. An efficient budget allocation algorithm for multi-channel advertising. In *2018 24th International Conference on Pattern Recognition (ICPR)*, pp. 886–891. IEEE, 2018.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Martello, S., Pisinger, D., and Toth, P. Dynamic programming and strong bounds for the 0-1 knapsack problem. *Management Science*, 45(3):414–424, 1999.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Nuara, A., Sosio, N., Trov , F., Zaccardi, M. C., Gatti, N., and Restelli, M. Dealing with interdependencies and uncertainty in multi-channel advertising campaigns optimization. In *The World Wide Web Conference*, pp. 1376–1386. ACM, 2019.
- Ren, K., Zhang, W., Chang, K., Rong, Y., Yu, Y., and Wang, J. Bidding machine: Learning to bid for directly optimizing profits in display advertising. *IEEE Transactions on Knowledge and Data Engineering*, 30(4):645–659, 2017.
- Ren, K., Fang, Y., Zhang, W., Liu, S., Li, J., Zhang, Y., Yu, Y., and Wang, J. Learning multi-touch conversion attribution with dual-attention mechanisms for online advertising. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pp. 1433–1442. ACM, 2018.
- Ren, K., Qin, J., Zheng, L., Yang, Z., Zhang, W., and Yu, Y. Deep landscape forecasting for real-time bidding advertising. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 363–372. ACM, 2019.
- Roberge, M. *The Sales Acceleration Formula: Using Data, Technology, and Inbound Selling to go from 0to100 Million*. John Wiley & Sons, 2015.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Wu, D., Chen, X., Yang, X., Wang, H., Tan, Q., Zhang, X., Xu, J., and Gai, K. Budget constrained bidding by model-free reinforcement learning in display advertising. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pp. 1443–1451. ACM, 2018.
- Yuan, S., Wang, J., and Zhao, X. Real-time bidding for online advertising: measurement and analysis. In *Proceedings of the Seventh International Workshop on Data Mining for Online Advertising*, pp. 1–8, 2013.
- Zhang, W., Yuan, S., and Wang, J. Optimal real-time bidding for display advertising. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 1077–1086. ACM, 2014.
- Zhang, W., Ren, K., and Wang, J. Optimal real-time bidding frameworks discussion. *arXiv preprint arXiv:1602.01007*, 2016.
- Zhou, G., Zhu, X., Song, C., Fan, Y., Zhu, H., Ma, X., Yan, Y., Jin, J., Li, H., and Gai, K. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 1059–1068. ACM, 2018.
- Zhu, H., Jin, J., Tan, C., Pan, F., Zeng, Y., Li, H., and Gai, K. Optimized cost per click in taobao display advertising. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 2191–2200. ACM, 2017.

---

## Appendix

---

### A. Background of Online Advertising

Online advertising is a marketing strategy involving the use of *advertising platform* as a medium to obtain website traffics and targets, and deliver marketing messages of *advertisers* to the suitable *customers*.

*Platform.* Advertising platform plays an important role in connecting consumers and advertisers. For consumers, it provides multiple advertising channels, e.g., channels on news media, social media, E-commerce websites and apps to explore. For advertisers, it provides automated bidding strategies to compete for consumers in all channels under the setting of real-time bidding (RTB), in which advertisers bid for ad exposures and the exposures opportunities go to the highest bidder with a cost which equals to the second-highest bid in the auction.

*Consumers.* Consumers explore multiple channels during the several visits to the platform within a couple of days. A consumer's final purchase of an item is usually a gradually changing process, which often includes the phases of Awareness, Interest, Desire, and Action (AIDA) (Roberge, 2015). The consumer's decision to convert (purchase a product) is usually and has to be driven by multiple touchpoints (exposures) with ads. Each advertising exposure during the sequentially multiple interactions could influence the consumers mind (preferences and interests) and therefore contribute to the final conversion.

*Advertisers.* The goal of advertisers is to cultivate the consumer's awareness, interest and finally driving purchase. As different ad strategies can affect consumers' AIDA, an advertiser should develop a competitive strategy to win the ad exposures in RTB setting. When the ad is displayed to a consumer, in Cost Per Click (CPC) setting, the advertisers should pay commission to the platform after the consumer clicking the ad. When the consumer purchases the advertised item, the advertiser will get the corresponding revenue.

The objective of an advertiser is usually to optimize the accumulated revenue within a time period under a budget constraint. A strategy that maximizes short-term revenue of each ad exposure on different channels independently is obviously unreasonable, since the final purchase is a result of long-term ad-consumer sequential interactions and the consumer's visits between different channels are interdependent. Therefore, the advertiser must develop a strategy to overcome following two key challenges: (1) Find the optimal interaction sequence including interaction times, channels selection and channel orders for a targeted consumer; (2) Choose targeted consumers and allocate predefined limited budget to them in multiple interaction sequences.

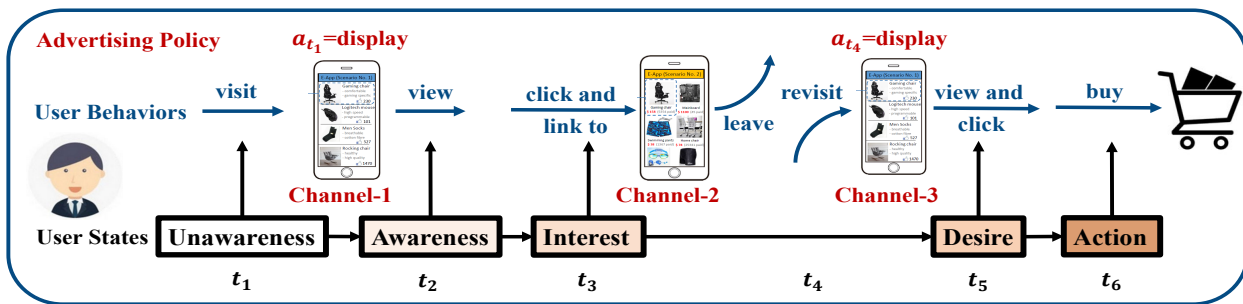


Figure 9. An illustration of the sequential multiple interactions (across different channels) between a user and an ad. Each ad exposure has long-term influence on the user's final purchase decision.

An example of a user's shopping journey is shown in Figure 9. At time  $t_1$ , a user visits the news media channel and triggers an advertising exposure opportunity. Then, the advertising agent executes a display action and leaves an exposure on the user. After that, the user becomes aware of and is interested in the commodity, so he clicks the hyperlink. Quickly, the user is induced into the landing (detail) page of the commodity in the shopping app. After fully understanding the product information, the user leaves the shopping app. After a period of time, the user comes back to the shopping app at time



$t_4$  and triggers an exposure opportunity of banner advertising. The advertising agent executes a display action as well. Consequently, the users desire is stimulated. At time  $t_6$ , the user makes a purchase. In this example, the ad exposure at time  $t_1$  influences the users mind and contributes to the ad exposure at time  $t_4$  and the delayed purchase, which means the ad exposure on one channel would influence the users preferences and interests, and therefore contributes to the final conversion. Thus, the goal of advertising should maximize the total cumulative revenue over a period of time instead of simply maximizing the immediate revenue.

## B. Proof and Analysis

### B.1. Knapsack Problem in Online Advertising Settings

**Theorem 2.** *The greedy solution to the proposed dynamic knapsack problem of online advertising is  $\lambda$  approximately optimal where  $\lambda > 99.9\%$*

*Proof.* In the proposed online advertising problem, each user is with value  $V_G$  (i.e. the profit of advertiser when the user purchase the commodity) and weight  $V_C$  (i.e. the total budget consumption for the target user in the real-time bidding to reach the final purchase). As the item (i.e. user) is non-splittable, the proposed dynamic knapsack problem is essentially a 0-1 knapsack problem which aims to maximize the total value of the knapsack given a fixed capacity  $B$ . For each item, we can calculate the Cost-Performance Ratio (CPR) as  $V_G/V_C$ . Sort all items in descending order of CPR, i.e.  $(V_{G_1}, V_{C_1}), (V_{G_2}, V_{C_2}), \dots, (V_{G_n}, V_{C_n})$  where  $\text{CPR}_i \geq \text{CPR}_j, \forall i \leq j \leq n$ . For  $V_C > 0, V_G > 0$  and  $B > 0$ , we first define that this 0-1 knapsack problem has optimal solution  $K^*(V_C, V_G, B)$  and greedy solution  $K(V_C, V_G, B)$  where  $K^*$  and  $K$  represent the total value of the knapsack.

Assume  $B_{end}$  is the remaining budget after greedy algorithm, the following inequality holds:

$$\frac{B - B_{end}}{B} K^*(V_C, V_G, B) \leq K(V_C, V_G, B) \leq K^*(V_C, V_G, B) \quad (13)$$

This is because:

- 1) If the knapsack can hold all the items after the greedy algorithm, that is, the optimal solution is equal to the greedy solution. As  $B_{end} \geq 0$ , we have  $\frac{B - B_{end}}{B} K^*(V_C, V_G, B) \leq K^*(V_C, V_G, B) = K(V_C, V_G, B)$
- 2) If the knapsack cannot hold all the items after the greedy algorithm, as  $\frac{V_{G_1}}{V_{C_1}} \geq \frac{V_{G_2}}{V_{C_2}} \geq \dots \geq \frac{V_{G_l}}{V_{C_l}}$ , we have  $V_{G_l} \sum_{j=1}^{l-1} V_{C_j} \leq V_{C_l} \sum_{j=1}^{l-1} V_{G_j} \Leftrightarrow V_{G_l} (B - B_{end}) \leq V_{C_l} K(V_C, V_G, B) \Leftrightarrow \frac{V_{G_l}}{V_{C_l}} \leq \frac{K(V_C, V_G, B)}{B - B_{end}} \Leftrightarrow K(V_C, V_G, B) \geq K^*(V_C, V_G, B) - \frac{B_{end} K(V_C, V_G, B)}{B - B_{end}}$  where  $l$  is the index of last item picked by greedy algorithm. This derivation can be simplified to  $K(V_C, V_G, B) \geq \frac{B - B_{end}}{B} K^*(V_C, V_G, B)$ .

In online advertising settings, the budget spent on a single user is much smaller than the advertiser's total budget. We conduct statistics on one of the world's largest E-commerce platforms to prove it. On Feb 3rd of 2020, a total of 1136149 ads result in 983414548 user-ad sequences (a user sequence consists of multiple interactions of the same user with the same ad), with an average of 865 user sequences per ad. Interactions with users of each ad forms a knapsack problem, where each user sequence is an item in the knapsack. The average maximum budget consumed by each user sequence accounts for 0.07068% of the total budget capacity of the advertisers. We also list details of 5 ads with largest budget consumption in Table 4, where the maximum budget consumed by each user sequence is much smaller than 1/1000 (smaller than 3/10000 specifically) of the total budget of each ad.

As proposed in Dantzig (1957),  $\forall i \in 1, 2, \dots, n, V_{C_i} \leq (1 - \lambda)B, 0 \leq \lambda \leq 1$ , the greedy algorithm achieves an approximation guarantee of  $\lambda$ . We can conclude from above statistics that  $\max_i \frac{V_{C_i}}{B} \leq \frac{1}{1000}$ , which means  $\lambda$  is much greater than  $1 - \frac{1}{1000}$ .

The thesis above can be further proved:

- 1) If the knapsack can hold all the items after the greedy algorithm, that is, the greedy solution is obviously equal to the optimal solution, which is also the  $\lambda$  approximately optimal solution.

## Appendix

Ad	#Users Sequences	Budget	Avg Cost	(Avg Cost)/Budget	Max Cost	(Max Cost)/Budget
Ad 1	2460976	119352.51	0.048498039	0.0000406343%	20.04	0.0167905979%
Ad 2	2674738	114388.54	0.04276626	0.000037388%	26.22	0.0229218766%
Ad 3	2848816	90113.08	0.031631766	0.0000351023%	15.29	0.0169675701%
Ad 4	2107497	82951.82	0.03936035	0.0000474497%	5.6	0.0067509067%
Ad 5	1087011	77140.49	0.070965694	0.0000919954%	19.32	0.0250452130%

Table 4. Detailed Comparison between an ad's total budget and cost on a user sequence.

- 2) If the knapsack cannot hold all the items after the greedy algorithm, we have  $V_{C_l} > B_{end}$ , that is,  $B_{end} < V_{C_l} \leq (1-\lambda)B$ . According to Formula 13, we have

$$\begin{aligned}
 K(V_C, V_G, B) &\geq \frac{B - B_{end}}{B} K^*(V_C, V_G, B) \\
 &> \frac{B - (1-\lambda)B}{B} K^*(V_C, V_G, B) \\
 &= \lambda K^*(V_C, V_G, B)
 \end{aligned} \tag{14}$$

Therefore, in theory, the greedy solution in our online advertising settings is  $\lambda$  approximately optimal and the  $\lambda$  is much greater than 99.9% in our case.

### B.2. Regretless Optimal Bidding Strategy $b_t^*$

**Theorem 3.** During the online bidding phase, the bidding agent can always set the bid price as:

$$\begin{aligned}
 \mathbf{b}_t^* &= \left[ \left( \frac{Q_G(s, \hat{a}_t = 1)}{\text{CPR}_{\text{thr}}^*} - Q_C^{\text{next}}(s, \hat{a}_t = 1) \right) \right. \\
 &\quad \left. - \left( \frac{Q_G(s, \hat{a}_t = 0)}{\text{CPR}_{\text{thr}}^*} - Q_C^{\text{next}}(s, \hat{a}_t = 0) \right) \right]
 \end{aligned} \tag{15}$$

where  $Q_C^{\text{next}}(s, \hat{a}_t) = \mathbb{E}[\sum_{k=t+1}^{T_j} c_k | s, \hat{a}_t, \pi_j]$ .  **$\mathbf{b}_t^*$  is a regretless optimal bidding strategy without any loss of accuracy.**

*Proof.* Since  $\mathbf{bid}_t^{2\text{nd}}$  is unknown until the current auction is finished, we prove the regretless of  $\mathbf{b}_t^*$  from the following two cases:

- 1) If  $\mathbf{b}_t^* > \mathbf{bid}_t^{2\text{nd}}$ :  $\mathbf{b}_t^* > \mathbf{bid}_t^{2\text{nd}} \Leftrightarrow Q(s, \hat{a}_t = 1) > Q(s, \hat{a}_t = 0)$ , which means the agent should take action  $\hat{a}_t = 1$  in this case. Exactly,  $\mathbf{b}_t^*$  is greater than the second highest price  $\mathbf{bid}_t^{2\text{nd}}$  based on the condition for entering the current branch. Thus, the agent will always win the auction and the executed action is indeed  $\hat{a}_t = 1$ .
- 2) If  $\mathbf{b}_t^* \leq \mathbf{bid}_t^{2\text{nd}}$ :  $\mathbf{b}_t^* \leq \mathbf{bid}_t^{2\text{nd}} \Leftrightarrow Q(s, \hat{a}_t = 1) \leq Q(s, \hat{a}_t = 0)$ , which means the agent should take action  $\hat{a}_t = 0$  in this case. Exactly,  $\mathbf{b}_t^*$  is less than the second highest price  $\mathbf{bid}_t^{2\text{nd}}$  according to the condition. Thus, the agent will always lose the auction and the executed action is indeed  $\hat{a}_t = 0$ .

Thus, we complete the proof.

### B.3. Convergence Analysis of MSBCB

The overall framework of MSBCB can be described as follows:

- (1) Let the budget constraint of an advertiser be  $B$ . Given a  $\text{CPR}_{\text{thr}}$ , we can use reinforcement learning algorithms to ensure that each user  $i$  is optimized according to  $\pi_i^* := \arg\max_{\pi_i} [V_G(i|\pi_i) - \text{CPR}_{\text{thr}} * V_C(i|\pi_i)]$  and converges to the optimal policy  $\pi_i^*$  under the current  $\text{CPR}_{\text{thr}}$ . Further, picking all users whose  $\text{CPR}_i \geq \text{CPR}_{\text{thr}}$  will result in a total cost of  $B'$  (i.e., the advertiser spends a budget  $B'$ ).

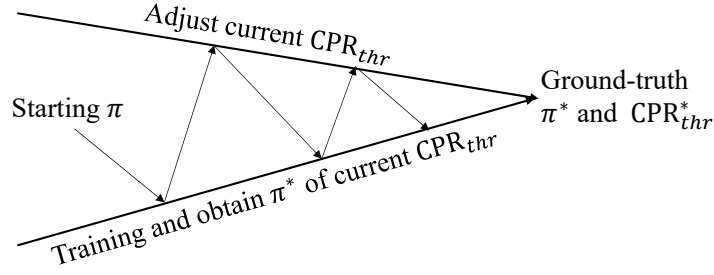


Figure 10. Convergence demonstration of MSBCB

- (2) As the current estimated threshold  $CPR_{thr}$  might have some bias from the optimal  $CPR_{thr}^*$ ,  $B'$  may not equal to the budget  $B$ . Thus, we design a PID controller to dynamically adjust the estimated  $CPR_{thr}$  so as to minimize the gap between the budget constraint  $B$  and the actual feedback of the daily cost  $B'$ .

As described in Figure 10, *MSBCB* repeats the above two steps iteratively. Given an updated  $CPR_{thr}$ , each  $\pi$  will be optimized by the lower-level reinforcement learning algorithms and  $\pi$  will move towards the optimal  $\pi^*$ . As a result, users whose optimized  $CPR_i \geq CPR_{thr}$  will be selected and we get the daily cost  $B'$ . Then, the current  $CPR_{thr}$  will be updated so that the gap between the cost  $B'$  and the budget  $B$  will be further minimized. Thus,  $CPR_{thr}$  will move towards the optimal  $CPR_{thr}^*$  gradually. As long as the learning rates of  $\pi$  and  $CPR_{thr}$  are small enough, the overall iterations will finally converge. In this paper, we also validate the convergence of our *MSBCB* in the experiments. As shown in Section 4.2 of the paper, our method converges quickly and finally reaches an approximation ratio of 98.53%.

## C. Deployment

Here we give the online deployment details of our *MSBCB*.

### C.1. Myopic to Non-Myopic Advertising System Upgrade Solution

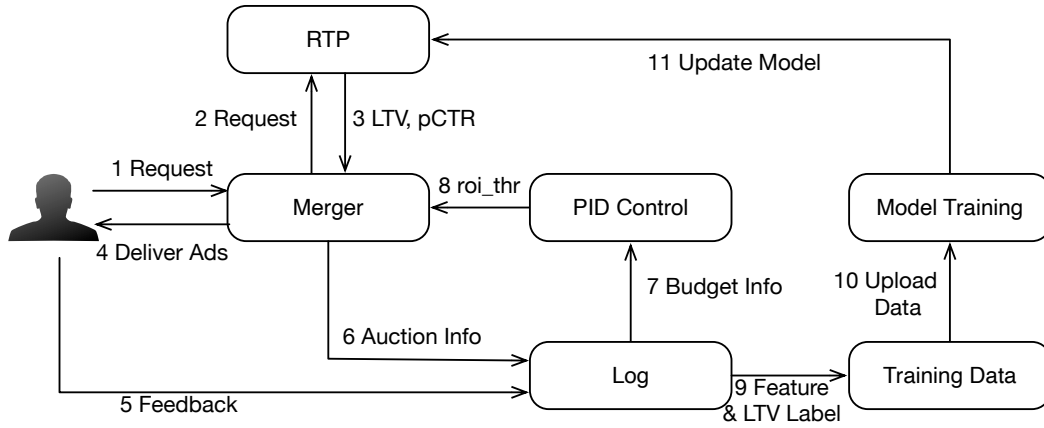


Figure 11. Online System

A myopic advertising system includes several key components as Figure 11 shows: (1) Log module collects auction information and user feedback. (2) Training data are constructed based on log followed by model training with offline evaluation. (3) Real-time prediction (RTP) module provides service for myopic value prediction of user-ad pairs. RTP periodically pulls newly trained models. (4) Merger module receives the user visit, requests RTP for myopic value with which ad bid adjustment ratios and ranking scores are calculated (In advertising, ranking score is  $ecpm = pCTR * bid$  where  $pCTR$  is predicted Click Through Rate and  $bid$  is the bidding price). Finally, top-scored ads are delivered to the user. Above myopic advertising system can upgrade to a non-myopic system by considering the following key changes.

(1) Log module needs to keep long-term auction information and users' feedback, and these data are used to construct features and long-term labels for training. Besides, logged data have to track each advertised item's budget and current cost data which are fed to a PID control module to compute  $CPR_{thr}$  for users selection in Merger. (2) Model training can use Monte Carlo (MC) or Temporal Difference (TD) methods. For MC, the long-term labels are cumulative rewards of a sequence and the training becomes a supervised regression problem. For TD, one-step or multi-step rewards are used to compute a bootstrapped long-term value using a separate network for training. (3) RTP module should periodically pull both myopic and non-myopic newly trained models and provide corresponding value prediction service. (4) Merger maintains an  $\langle \text{item}, CPR_{thr} \rangle$  table which is updated periodically from PID module. When a user visit comes, Merger requests RTP for both  $pCTR$  and long-term values (long-term  $GMV$  i.e.  $V_G$  and  $cost$  i.e.  $V_C$  in our paper), and with  $CPR_{thr}$  decides the selection of current user and bid adjustment.

## C.2. Long-Term Value Prediction Model

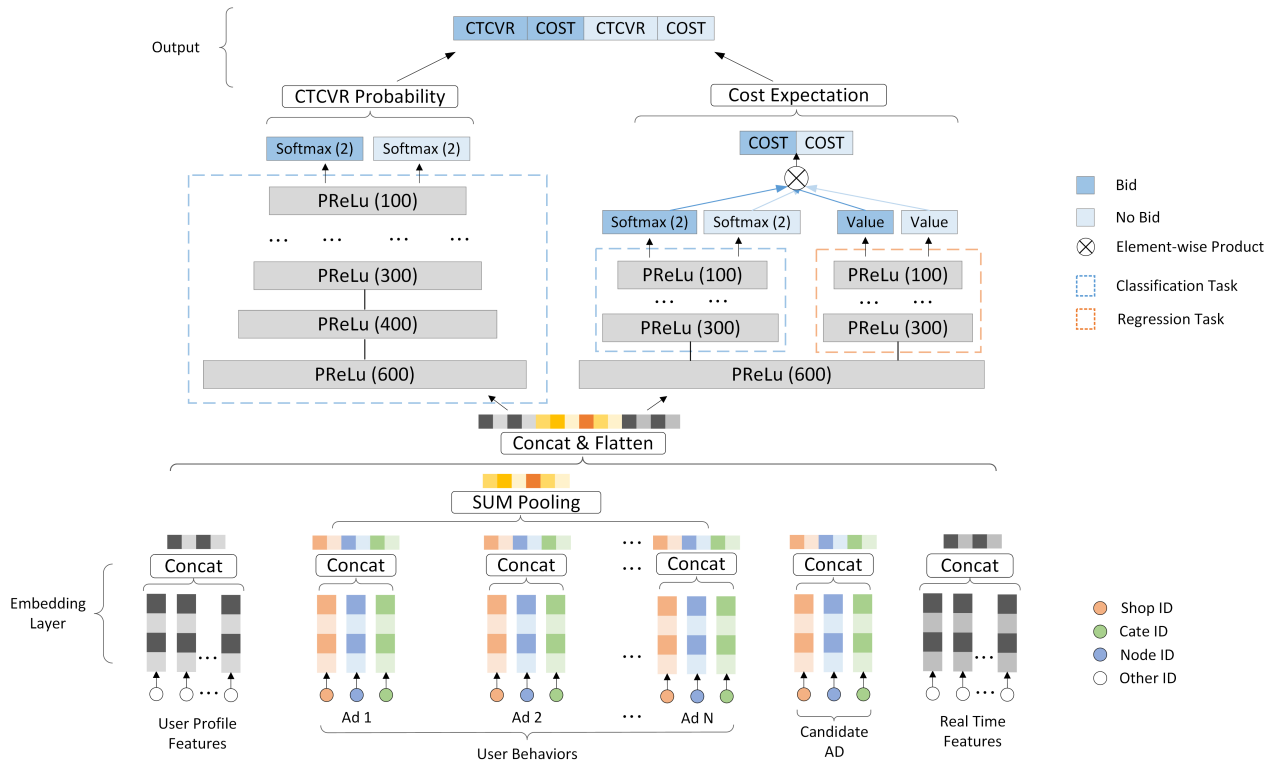


Figure 12. Long-Term Value Prediction Model

### C.2.1. FEATURES AND LABELS

Features for long-term value prediction should contain sufficient user's static profile and historical behavior information. Most myopic advertising systems already have a sound feature system which can summarize user-oriented, ad-oriented and user-ad interactive history very well. Besides, due to the large amount of data collected by the online advertising system, these features are able to generalize across large number users where each user-ad pair's interaction is considered as a separate MDP, thus, help the prediction model learning. To be specific, the state  $s_t$  at step  $t$  includes: 1) user profile features; 2) user behavior features; 3) real-time user behavior features; 4) context features; 5) user-ad interaction histories; 6) user feedback before current step  $t$  and so on. Features are constructed based on past 7-14 days data before user visit time  $t$ . For the MC training method, labels are constructed using the following 7 days data after user visit time  $t$ . For the TD method, labels are the instant rewards at time  $t$  and the long-term labels are constructed using a bootstrap method.



### C.2.2. MODEL ARCHITECTURE

The long-term value model architecture is shown in Figure 12, where the model takes the features as input and output long-term value of  $GMV$  (i.e.  $V_G$  in our formulation) and  $cost$  (i.e.  $V_C$  in our formulation) for both action=1 (display the ad) and action = 0 (do not display the ad).

We use one model to output multiple long-term values ( $GMV$  and  $cost$  for action=1 and action=0). Multiple prediction tasks share the same bottom layers because we consider the underlying knowledge of the user’s sequence behaviors such as opening the app, jumping across channels, turning off the phone and revisiting the app should be learned together and shared. The shared layer converts input features to embeddings and embeddings in the same group are concatenated. The user-behavior group embeddings are then pooled with sum operation. User-profile embeddings, user-behavior embeddings, candidate ad embeddings, and real-time features are finally concatenated and flattened as the output of the bottom layers.

Following the shared bottom layers, the network is split into two forward-pass branches where one is for long-term  $GMV$  prediction and one for long-term cost prediction. We find this two-branch design can reduce the influences among different tasks and stabilize the learning. For the long-term  $GMV$  prediction, since each user usually buys a commodity only once, we only have to predict  $P(buy > 0|feature)$  denoted as  $CTCVR$ . In the online inference phase, the long-term  $GMV$  is computed with  $GMV = P(buy > 0|feature) * item\_price$  where  $item\_price$  is the price of the commodity. For the long-term cost prediction, in CPC (Cost-Per-Click) advertising, a user usually clicks several times before buys and the cost per click along with each click varies, thus, the long-term  $cost$  prediction cannot be decomposed as the long-term  $GMV$  prediction and the only way is to regress the long-term cost value. However, as most sequences’ costs are zero, the direct regression learning process will be very noisy. Therefore, we design an additional hidden layer to compute  $P(cost > 0|feature)$ ,  $P(cost = 0|feature)$  and  $E(cost|cost > 0, feature)$ . Then, the predicted long-term cost is computed as  $pcost = P(cost > 0|feature) * E(cost|cost > 0, feature) + P(cost = 0|feature) * 0 = P(cost > 0|feature) * E(cost|cost > 0, feature)$  where  $P(cost > 0|feature)$  and  $P(cost = 0|feature)$  are learned using logistic regression loss and  $pcost$  is learned using mean-square error loss  $(pcost - cost)^2$ . We find the above designs help improve the model’s prediction performance in practice. For  $CTCVR$  and  $P(cost > 0|feature)$ ,  $P(cost = 0|feature)$ , we use GAUC (Zhou et al., 2018) as metric, and for  $pcost$  regression, we use mean-square error and reverse order metrics.

## D. Empirical Evaluation: Supplementary of Offline Experiments

### D.1. Experiments Settings.

Considering the potential losses of assets and money, it’s usually forbidden to do a lot of trial and error and thoroughly comparisons between available baselines in a live advertising system. Thus we implement a fairly general simulation environment so that we could make extensive analyses of our approach. All experiments are conducted on an Intel(R) Xeon(R) E5-2682 v4 processor based Red Hat Enterprise Linux Server, which consists of two processors (each with 16 cores), running at 2.50GHz (16 cores in total) with 32KB of L1, 256 KB of L2, 40MB of unified L3 cache, and 128 GB of memory and 2 Tesla M40 GPUs.

### D.2. Simulation Environment.

Here, we give the detail of the simulation environment. Similar to (Ie et al., 2019), the simulation environment includes the following 5 modules:

- *Advertisements and Topic Model*: We assume a set of documents  $\mathcal{D}$  representing the content available for advertising. We also assume a set of topics (or users interests)  $\mathcal{T}$  that capture fundamental characteristics of interest to users; we assume topics are indexed  $1, 2, \dots, |\mathcal{T}|$ . Each commodity  $d \in \mathcal{D}$  has an associated topic vector  $\mathbf{d} \in [0, 1]^{|\mathcal{T}|}$ , where  $d_j$  is the degree to which  $d$  reflects topic  $j$ . Each document  $d \in \mathcal{D}$  also have an inherent quality  $Q_d \in [0, 1]$ , representing the topic-independent attractiveness to the average user.
- *Consumer Interest and Satisfaction Model*: Each user  $i$  has various degrees of interests in topics, ranging from 0 (completely uninterested) to 1 (fully interested), with each user  $i$  associated with an interest vector  $\mathbf{u} \in [0, 1]^{|\mathcal{T}|}$ . Consumer  $i$ ’s interest in advertisement  $d$  is given by the dot product  $I(u, d) = \mathbf{u}\mathbf{d}$ . We assume some prior distribution  $P_u$  over user interest vectors, but user  $i$ ’s interest vector is dynamic, i.e., influenced by their advertisement consumption (see

below). Besides, a user’s satisfaction  $S(u, d)$  with a consumed (viewed) advertisement  $d$  is a function  $f(I(u, d), Q_d)$  of user  $i$ ’s interest and ad  $d$ ’s quality. Here, we assume a simple convex combination  $S(u, d) = (1 - \alpha)I(u, d) + \alpha Q_d$ . Satisfaction influences user dynamics as we discuss below.

- **Consumer Choice Model:** The user’s Click-Through Rate (CTR) and Conversion Rate (cvr) are represented by  $I(u, d)$  and  $S(u, d)$  respectively. Each user has the probability of clicking and buying an advertising commodity according the CTR and CVR.
- **Consumer Dynamics:** We assume that a user’s interest evolves as a function of the documents consumed (viewed). When user  $i$  consumes document  $d$ , her interest in topic  $T(d)$  is nudged stochastically, biased slightly towards increasing her interest, but allows some chance of decreasing her interest. In this paper, we set  $\mathbf{u} \leftarrow \gamma \mathbf{u} + \beta * S(u, d) * \mathbf{d}$ , where  $\gamma$  is the interest decay rate and  $\beta \in [-1, 1]$  is a user independent parameter.
- **Consumer Visiting Model and Advertising System Dynamics:** The users’ request sequence are generated from a stable distribution  $P_{req}$ . For each user’s request, all advertisements  $d \in \mathcal{D}$  give a bid and competes with other bidders in real-time. The winner has the privilege to display its ad to the user, which could further influence the user’s interest and behavior.

### D.3. Codes and Datasets.

The codes and datasets to reproduce our offline experiments are provided in another supplementary material.

### D.4. Cost Comparison.

The consumption of budget during the training process is shown in Figure 13. As we can see, the costs of all approaches converge to about 12000, which is exactly equal to the budget we set in experiments. Specific costs of each approach can be found in Table 1 of paper.

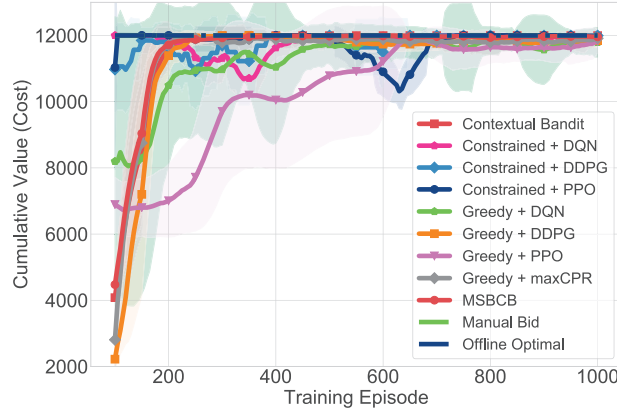


Figure 13. The learning curves of costs of our *MSBCB* and the other baseline approaches.

### D.5. Convergence Analyses

#### D.5.1. CONVERGENCE OF EACH $\pi_i^*$ GIVEN ANY $CPR_{THR}$ .

As shown in Figure 14, given a  $CPR_{thr}$ , the learned advertising policy  $\pi$  of our *MSBCB* converges to the optimal  $\pi_j^*$ . In Figure 14, the x-axis denotes the cumulative cost, the y-axis denotes the cumulative value and the dots in blue represent the cumulative values and costs of all possible policies for each user. The red line represents  $y = CPR_{thr} * x$ , whose slope is  $CPR_{thr}$ . The orange point represents the optimal policy  $\pi_i^*$  computed by enumerating all possible policies (blue points) and finding the one which maximize  $V_G(i|\pi_i) - CPR_{thr} * V_C(i|\pi_i)$  according to **Theorem 1**. The green point denotes the learned policy of *MSBCB*. In theory, the point of the optimal policy is the one whose  $CPR > CPR_{thr}^*$  and vertical distance is the farthest from the red line. A proof is provided in the **Theorem 4** in the later part. We present 3 convergence examples of different types in Figure 14. In Figure 14 (a) and (b), the learned  $\pi$  by the RL algorithm is exactly the same with the optimal  $\pi^*$ . In Figure 14 (b), the optimal policy is do not advertise to this user. In Figure 14 (c), the learned  $\pi$  is

approximately optimal. Detail convergence statistics on the proportion of users whose policies converged to the optimal

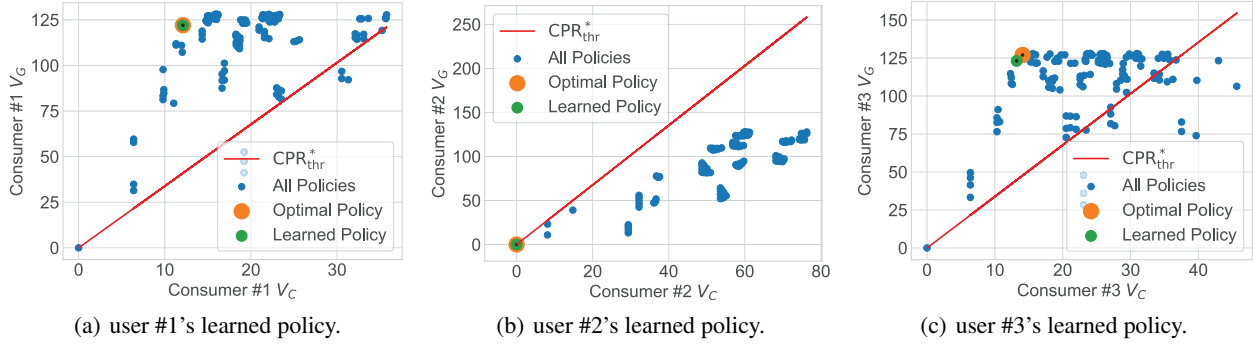


Figure 14. Three examples of the convergence of each  $\pi_i^*$  given a fixed  $CPR_{thr}$ .

ones among all users are shown in Table 5. For each user, we denote the vertical distance of the learned policy to the  $CPR_{thr}^*$  line as  $dis_{learned}^*$  and the vertical distance of the optimal policy  $\pi^*$  to the  $CPR_{thr}^*$  line as  $dis_{optimal}^*$ . We denote  $R_{opt}^* = dis_{learned}^*/dis_{optimal}^*$  as the approximation ratio. According to **Theorem 4**, if the  $R_{opt}^*$  is 100%, then the learned strategy is exactly the optimal strategy. Otherwise, we denote that the learned strategy is the  $R_{opt}^*$ -approximation strategy. As shown in Table 5, there are 74.9% policies achieve more than 90%-approximation ratios and 53.3% policies achieve exactly the optimal.

Table 5. Optimal types of each  $\pi_i^*$  of 10000 users

$R_{opt}^*$	100%	[90%, 100%)	[0%, 90%)
Percentage	53.3%	21.6%	25.3%

**Theorem 4.** The point of the optimal policy is the one whose  $CPR > CPR_{thr}^*$  and vertical distance to  $CPR_{thr}^*$  line (red line) is the farthest among all policy dots in Figure 14.

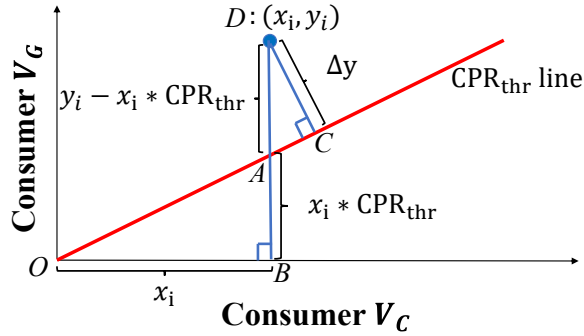


Figure 15. Proof of Optimal Policy Dot

*Proof.* Here we give a simple proof of **Theorem 4**. As we can see in Figure 15, the blue dot  $D : (x_i, y_i)$  is an arbitrary policy  $i$  in Figure 14. Suppose the vertical distance of  $D$  to  $CPR_{thr}^*$  line (red line) is  $\Delta y_i$  (segment  $DC$  in the figure). We then draw a vertical line of x-axis from dot  $D$  to dot  $B$ . We can then calculate the length of segments:  $OB = x_i$ ,  $BA = x_i * CPR_{thr}$ ,  $DA = y_i - x_i * CPR_{thr}$ . It's evident that  $\triangle OAB \sim \triangle DAC$ , which means  $\frac{DC}{OB} = \frac{DA}{OA} = \frac{DA}{\sqrt{(OB)^2 + (AB)^2}}$ . We can derive that

$$\frac{\Delta y_i}{x_i} = \frac{y_i - x_i * CPR_{thr}}{\sqrt{x_i^2 + (x_i * CPR_{thr})^2}} \quad (16)$$

As  $x_i > 0$ , we can further derive that

$$\Delta y_i = \frac{y_i - x_i * \text{CPR}_{thr}}{\sqrt{1 + \text{CPR}_{thr}^2}} \quad (17)$$

Suppose the dot of a policy is  $(x^*, y^*)$ , which has farthest vertical distance  $\Delta y^*$  from the  $\text{CPR}_{thr}^*$  line, that is, for a dot of arbitrary policy  $i$ , we have  $\Delta y_i \leq \Delta y^*$ . According to Equation 17, we have

$$\frac{y_i - x_i * \text{CPR}_{thr}}{\sqrt{1 + \text{CPR}_{thr}^2}} \leq \frac{y^* - x^* * \text{CPR}_{thr}}{\sqrt{1 + \text{CPR}_{thr}^2}} \quad (18)$$

Then we get  $y_i - x_i * \text{CPR}_{thr} \leq y^* - x^* * \text{CPR}_{thr}$ , which means  $(x^*, y^*)$  is the dot of the optimal policy. Thus, we complete the proof.

#### D.5.2. CONVERGENCE OF $\text{CPR}_{thr}^*$ .

In Figure 16, we plot the learning curves of the  $\text{CPR}_{thr}$  of our *MSBCB* as well as 3 RL approaches. The dotted blue line denotes the optimal  $\text{CPR}_{thr}^*$  computed by the *MSBCB (enum)* of Table 1 of paper. Figure 16 shows that the learned  $\text{CPR}_{thr}$  of our *MSBCB* could gradually converge to the optimal  $\text{CPR}_{thr}^*$  approximately, which is much better than the other 3 RL approaches.

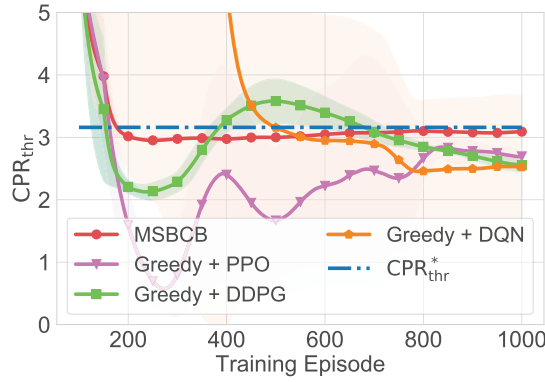


Figure 16. The convergence of  $\text{CPR}_{thr}^*$

#### D.6. Gap to Market Second Price.

Figure 17 shows average gaps between the bid of the agent of different approaches and the second price in the auction. Results indicate that the bid prices given by the *MSBCB* agent are closer to the second price in the auction, which can reduce the risk of economic loss when the market price fluctuates.

#### D.7. Effectiveness of Action Space Reduction.

Here we give a more detailed comparison of *MSBCB* and RL baselines to demonstrate the effectiveness of action space reduction. As shown in Figure 18 and Table 6, *MSBCB* (with action space reduction) can reach exactly the same cumulative value much more quickly than the other 3 RL baselines. *MSBCB* can reach a cumulative value of 85000 in only 104 epochs, which proves that action space reduction can effectively improve the sample utilization to converge to higher performance with faster speed.

### E. Empirical Evaluation: Supplementary of Online A/B Testing

In online A/B Testing, we conduct further analyses to verify the effectiveness of our *MSBCB* and find out whether our approach could benefit most advertisers.



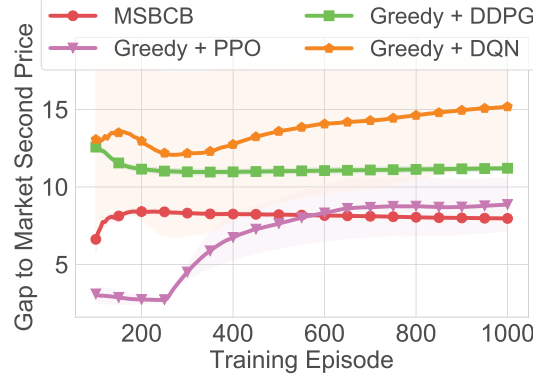


Figure 17. The average gaps of the bids to the second prices in the auction by different methods.

Table 6. The training epochs and the number of samples needed by different approaches when achieving the same revenue level.

Cumulative Value Method	60000		65000		70000		75000		80000		85000	
	#Epoch	#Samples	#Epoch	#Samples	#Epoch	#Samples	#Epoch	#Samples	#Epoch	#Samples	#Epoch	#Samples
Greedy + PPO	251	1280000	299	1530880	776	3973120	817	4183040	-	-	-	-
Greedy + DDPG	68	343040	76	389120	92	471040	154	788480	853	4362240	-	-
Greedy + DQN	90	455680	109	558080	153	783360	373	1909760	754	3855360	-	-
<b>MSBCB</b>	<b>22</b>	<b>112640</b>	<b>33</b>	<b>163840</b>	<b>48</b>	<b>245760</b>	<b>61</b>	<b>312320</b>	<b>71</b>	<b>363520</b>	<b>104</b>	<b>532480</b>

Firstly, we analyze the performance of our *MSBCB* for each advertiser. To guarantee the statistical significance, only the advertisers with more than 100 conversions in a week are included. The detail results of top-10 advertisers with the largest costs are shown in Table 7. In Table 7, under the same budget constraint, our *MSBCB* can increase the Revenues and ROIs of most advertisers compared with the myopic *Contextual Bandit* approach. Although the ROI of advertiser 7 drops slightly, our *MSBCB* contributes to much more PVs (Page Views).

Besides, in Figure 19, we give the detail proportions of advertisers whose ROIs are improved. Among all advertisers, 85.1% advertisers obtain positive ROI improvements while the rest of 14.9% advertisers are in the so-called quantity and quality exchange situations: their PV increments are larger than the ROI drops. We say that its also acceptable for some advertisers because the PV increments might lead to secondary exposures to an advertiser and thus lower the ROI within the current time period. But the increase in PV may leave deeper impressions to the users and contribute to the long-term future revenues. In addition, Figure 19 demonstrates that our *MSBCB* can be well applied to the multi-agent setting (which involves multiple advertisers) in the real-world auction environment, which could increase the overall revenue for most advertisers.

In order to highlight the advantage of our method in long-term revenue optimization, we compared the average number of

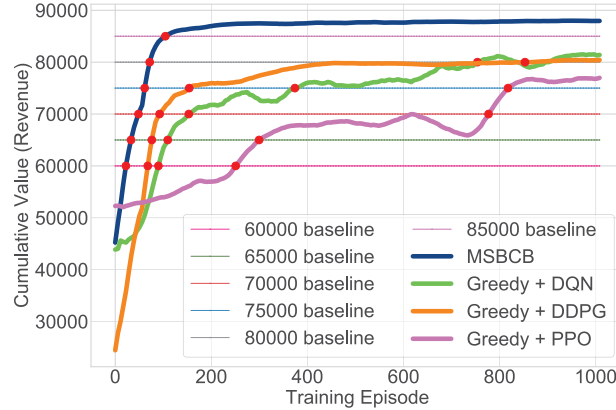


Figure 18. The comparison of the number of training episodes needed by different approaches when achieving the same revenue level.

## Appendix

	Revenue	Cost	CVR	PV	ROI
Advertiser 1	5.1%	-6.3%	17.2%	9.6%	12.2%
Advertiser 2	7.5%	2.1%	5.2%	12.2%	5.3%
Advertiser 3	48.6%	10.9%	27.6%	28.9%	33.9%
Advertiser 4	3.1%	2.8%	1.1%	9.6%	0.3%
Advertiser 5	12.7%	1.7%	12.9%	17.8%	10.8%
Advertiser 6	10.8%	2.2%	4.4%	13.8%	8.4%
Advertiser 7	1.9%	3.8%	4.6%	31.5%	-1.8%
Advertiser 8	5.6%	-4.8%	2.9%	10.7%	11.1%
Advertiser 9	6.7%	-2.4%	6.3%	21.0%	9.4%
Advertiser 10	5.8%	-0.8%	2.5%	8.0%	6.7%

Table 7. The improvements in Revenue, CVR, PV and ROI of our *MSBCB* compared with the myopic *Contextual Bandit* method.

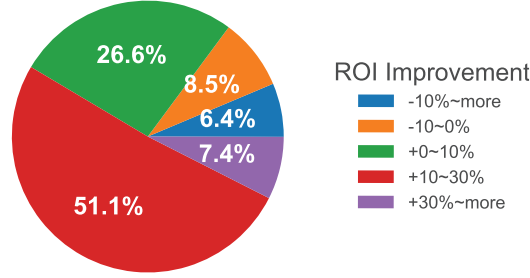


Figure 19. The distribution of ROI improvements for all advertisers of our *MSBCB* compared with the myopic *Contextual Bandit* method.

times (we call the sequence length) that a user contact with an advertisement under different approaches. Figure 20 shows the extent of *MSBCB*'s improvement relative to *Contextual Bandit* in the proportion of the user sequence length. The results show that our *MSBCB* can increase the proportion of the sequences with larger sequence length. Especially, the ratio of sequence length of 7 is increased by nearly 30%. It shows that our method can promote longer user behavior sequences, and longer user behavior sequence means more opportunities to affect the user's mentality towards an advertisement, thereby improving the long-term revenue for an advertisement.

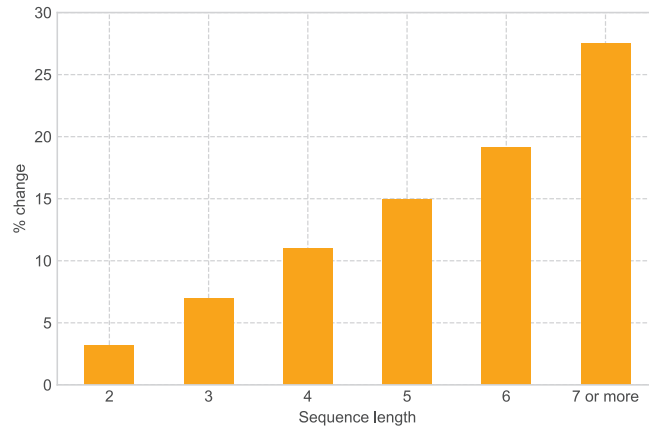


Figure 20. The proportion improvements in the sequence length of our *MSBCB* compared with the myopic *Contextual Bandit* method.

Further, we also analyze the ROI performances of the compared 3 algorithms (i.e., *CEM*, *Contextual Bandit* and our *MSBCB*) in different channels. Figure 21 shows the budget allocation distributions of all approaches among 6 channels and the corresponding ROIs. The left axis represents the ROI, and the ROI performances of each algorithm among different channels are given by the corresponding bar charts. The right axis represents the increments or decrements of the actual costs of *MSBCB* and *Contextual Bandit* relative to *CEM*, which are indicated by the line charts. In Figure 21, we observe the

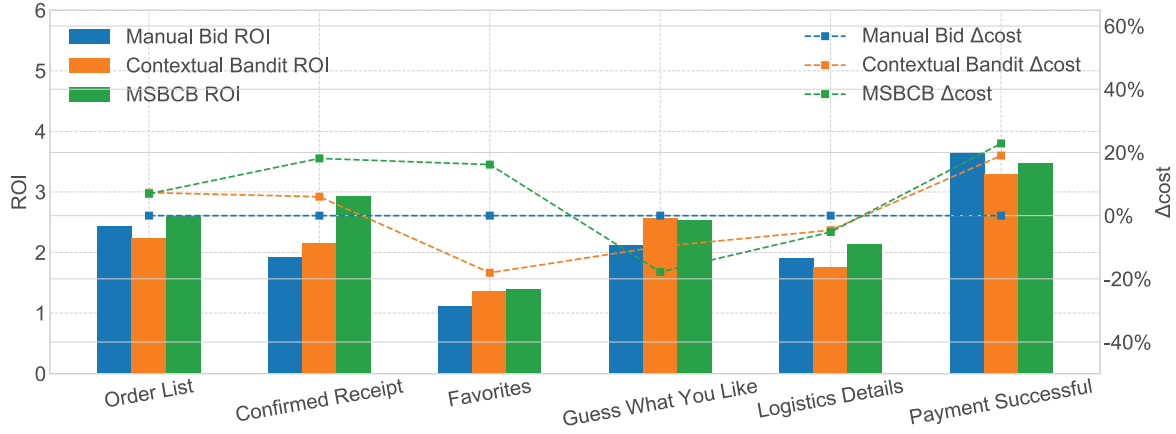


Figure 21. ROI and budget allocation among different channels.

following two phenomena:

- 1) *MSBCB* and *Contextual Bandit* both spend more budgets on channels with higher ROIs, especially on the *Payment Successful* channel, where the average ROI is much higher.
- 2) Compared with *Contextual Bandit*, *MSBCB* allocates more budget from the *Guess What You Like* channel to other channels, especially the *Favorites* channel, *Confirmed Receipt* channel and the *Payment Successful* channel.

These phenomena show that our *MSBCB* can reasonably allocate budgets among different channels and spend more budgets in channels with higher ROIs. In addition, compared with the myopic method *Contextual Bandit*, our long-term *MSBCB* is more optimistic about channels during and after purchasing, which shows that our *MSBCB* prefers a longer interaction sequence to optimize cumulative long-term values.