

# Supplement

## 1 Mathematical Formulation of Bidding Strategy and Online Bidding Problem

We focus on designing a bid optimization strategy for RTB auctions. When facing a bid request together with an impression, a DSP evaluates specific values of the impression, and then an impression-level bid optimization strategy maps the values, advertisers' requirements, and remnant budget to a real-time bid price.

### 1.1 Technical terms and notations

First, some technical terms and notations are explained; refer to [1,2] for more details.

- (a) RTB-related terms and notations.
  - **Impression**, an opportunity to show an ad in front of users.
  - **Bid request**, a call for pricing an impression together with a message containing the impression feature.
  - **Ad campaign**, advertisers promote their products by setting up an ad campaign containing impression features of interests, KPIs, CPA constraints, etc.
  - **Auction scale**, denoted by  $T \in \mathbb{Z}_+$ . The estimated total number of bid requests an ad campaign receives within a specific interval, usually one-day.
  - **Click-through rate, CTR**, denoted by  $p^{\text{CTR}} \in [0, 1]$ . It is the estimated probability that users click through an ad when they see it.
  - **Conversion rate, CVR**, denoted by  $p^{\text{CVR}} \in [0, 1]$ . It is the estimated probability that users take predefined action after clicking through an ad.
- (b) Impression-related terms and notations.

- **Impression feature**, denoted by a high dimensional vector,  $\mathbf{x} \in \mathcal{X}$ . Various features can describe an impression, such as cookie information, time, location, etc. The feature  $\mathbf{x}$  can be regarded as an integrable random variable with an integrable probability density function (PDF)  $p_{\mathbf{x}}$ , which follows the independent identically distribution (i.i.d) assumption [2,3]. The latter assumption helps us use a simple generic stochastic model to design bid strategies and pay more attention to bid optimization.
- **Impression value**, denoted by  $v \geq 0$ . It is the estimated KPI of an impression if an advertiser wins the auction. It is provided by a feature mapping module that maps  $\mathbf{x}$  to  $v$  [4], that is,  $v = v(\mathbf{x})$ .
- **Impression action value**, denoted by  $v_{a,j} \geq 0$ . It is the estimated amount of action(s) an impression takes if an advertiser wins the auction [2]. There may exist  $k$  interested actions simultaneously that the subscript  $j \in \{1, \dots, k\}$  is used to distinguish them. It is also provided by a feature mapping module, that is,  $v_{a,j} = v_{a,j}(\mathbf{x})$ . It is different from  $v$  because of it is usually used to measure some intermittent performance, such as subscription, add to cart, and so on [3,5].

(c) Bid optimization problem-related terms and notations.

- **Key performance indicator, KPI**. It is a quantitative measurement of advertising performance [3] that advertisers are interested in.
- **Bid price**, denoted by  $b \geq 0$ . The cost that an advertiser wants to pay for an impression being auctioned.
- **Winning price**, denoted by  $w \geq 0$ . For an impression, winning price is the lowest bid price to win its auction and  $w = w(\mathbf{x})$  [2]. For an auction, the relationship among bid price, winning price, and corresponding auction result is,

$$\begin{cases} \text{Advertiser wins,} & b \geq w \\ \text{Advertiser loses,} & b < w \end{cases}.$$

- **Cost**, denoted by  $c \geq 0$ . It is the cost for the impression the advertisers win, which relates to  $\mathbf{x}$ , the auction mechanism and the billing method of DSP [2,3,6]. In this paper, we use the second-price auction, a mechanism in which each bidder provably tends to price on their desires in accordance with the truthful value of  $\mathbf{x}$  [7]. Besides, the billing method is another factor affecting the cost [1]. In this paper, we assume the billing method is click-based, that is, For example, in click ads, advertisers only pay for clicked-through ads [5]; while in display ads, billing has nothing to do with whether an ad has been clicked or not [8]. Thus, we define  $c = c(\mathbf{x}, *)$ , where  $*$  means the billing method.

- **Budget constraint**, denoted by  $B \geq 0$ . It is the total cost the advertiser can use for an ad campaign during its lifetime.
- **CPA constraint**, denoted by  $C \geq 0$ . It is the maximum average cost the advertiser can pay for a specified action observed on the delivered impression and usually is one-day-based. For example, CPC constraint means the maximum average cost the advertiser pays when a user clicks on the delivered ad of the won impression.

## 1.2 Problem of Interests

In the bid optimization problem, we aim to find a strategy  $\pi$  that, by mapping the impression feature  $\mathbf{x}$ , and the remnant budget  $\Delta B$  (equivalently, the consumed budget) to a bid price  $b$ , optimizes the expected KPI and holds CPA constraints. However,  $\mathbf{x}$  is a constructed, high-dimensional vector, e.g., the methodology proposed in [6], and objectives and constraints vary for advertisers. These complicate the direct use of  $\mathbf{x}$  and the corresponding bidding strategy generalization. Therefore, there exists a supplied module called feature mapping, which maps  $\mathbf{x}$  to specific predefined values [4, 9, 10], i.e., the impression value  $v(\mathbf{x})$  and the impression action value  $v_a(\mathbf{x})$ . Feature mapping is another vital problem in RTB and due to the limited space, we assumed it is prior to this work. Then, we write  $b_t = \pi_t(B_t, \mathbf{x}_t)$ , where  $B_t$  is the consumed budget at time  $t$ , and model the bid optimization as

$$\begin{aligned}
& \arg \max_{\pi_t, t \in \{1, \dots, T\}} \mathbb{E} \left[ \sum_{t=1}^T v(\mathbf{x}_t) \mathbf{1}(w(\mathbf{x}_t), \pi_t(B_t, \mathbf{x}_t)) \right] \\
& \text{s.t.} \quad \sum_{t=1}^T c(\mathbf{x}_t, *) \mathbf{1}(w(\mathbf{x}_t), \pi_t(B_t, \mathbf{x}_t)) \leq B \\
& \quad \frac{\sum_{t=1}^T \mathbb{E} [c(\mathbf{x}_t, *) \mathbf{1}(w(\mathbf{x}_t), \pi_t(B_t, \mathbf{x}_t))]}{\sum_{t=1}^T \mathbb{E} [v_{a,j}(\mathbf{x}_t) \mathbf{1}(w(\mathbf{x}_t), \pi_t(B_t, \mathbf{x}_t))]} \leq C_j,
\end{aligned} \tag{1}$$

where  $j = 1, \dots, k$ , indicator function  $\mathbf{1}(w(\mathbf{x}_t), \pi_t(B_t, \mathbf{x}_t))$  means the  $t$ -th auction result, if  $\pi_t(B_t, \mathbf{x}_t) \geq w(\mathbf{x}_t)$ , then  $\mathbf{1}(w(\mathbf{x}_t), \pi_t(B_t, \mathbf{x}_t)) = 1$ ; otherwise, it equals to 0. The winning probability  $g(b_t, \mathbf{x}_t)$  is the expectation of  $\mathbf{1}(w(\mathbf{x}_t), \pi_t(B_t, \mathbf{x}_t))$  with respect to  $\mathbf{x}_t$ . The objective in the above is to find  $\pi_t, t \in \{1, \dots, T\}$  maximizing the total amount of impression value  $v$ . Nevertheless,  $\mathbf{x}$  is a random variable with  $p_{\mathbf{x}}$ , the objective is taken in the expected sense. The first constraint means the consumed budget is limited by a predefined  $B$ ; we should bid carefully for each auction because we may miss desired impressions when the budget depleted. The second constraint measures the cost per actions determined by advertisers (see Section ??), such as cost per click, cost per conversion. etc.

### 1.2.1 A real-world example

As an example, we present a real-world business service [5] provided by *Taobao.com* to illustrate the functional optimization problem (1). Its main settings contain:

- **Maximizing objective**,  $v(\mathbf{x}) = p^{\text{CTR}}(\mathbf{x})p^{\text{CVR}}(\mathbf{x})$ , is the expected total number of transactions.
- **Constraints**, advertisers set a total budget constraint  $B$ , and there is only a cost per click constraint  $C$  for a day.
- **Billing method**, advertisers only pay for those clicked-through ads they won, abbreviated as ‘click’. In this case, the cost is  $c(\mathbf{x}, *) = c(\mathbf{x}, \text{‘click’})$ .
- **Others**. The only considered impression action value  $v_a(\mathbf{x}) = p^{\text{CTR}}(\mathbf{x})$ . The second-price auction is used; cost equals the winning price, that is, the second-highest bid price.

To sum up, its bid optimization problem is

$$\begin{aligned}
& \arg \max_{\pi_t, t \in \{1, \dots, T\}} \mathbb{E} \left[ \sum_{t=1}^T p^{\text{CTR}}(\mathbf{x}_t) p^{\text{CVR}}(\mathbf{x}_t) \mathbf{1}(w(\mathbf{x}_t), \pi_t(B_t, \mathbf{x}_t)) \right] \\
& \text{s.t.} \quad \sum_{t=1}^T c(\mathbf{x}_t, \text{‘click’}) \mathbf{1}(w(\mathbf{x}_t), \pi_t(B_t, \mathbf{x}_t)) \leq B \\
& \quad \frac{\sum_{t=1}^T \mathbb{E} [c(\mathbf{x}_t, \text{‘click’}) \mathbf{1}(w(\mathbf{x}_t), \pi_t(B_t, \mathbf{x}_t))]}{\sum_{t=1}^T \mathbb{E} [p^{\text{CTR}}(\mathbf{x}_t) \mathbf{1}(w(\mathbf{x}_t), \pi_t(B_t, \mathbf{x}_t))]} \leq C.
\end{aligned} \tag{2}$$

### 1.3 Markov Decision Process (MDP) Modelling

The problem (1) is a typical sequential decision making problem. A remarkable class of well-studied sequential decision-making problems is the Markov decision process (MDP) or reinforcement learning, where MDPs are applied to more complex tasks and solved in a data-driven manner. In this part, we will formulate (1) into an MDP problem. To do so, we need to change the formulation of (1) by redefining the objective with its constraints as punishment terms. The reason for doing so is that whether a constraint of (1) is violated or not cannot be tested from real-time accessed data. In this part, we consider a finite time-horizon Markov decision process  $\{\mathcal{S}, \mathcal{B}, p, r, T\}$ :

- **State set  $\mathcal{S}$** . For each start of the  $t$ -th auction, the consumed budget  $B_t$  and the current impression feature  $\mathbf{x}_t$  are the state  $s_t = [B_t \ \mathbf{x}_t^\top]^\top$ .
- **Bid price set  $\mathcal{B}$** . For the  $t$ -th auction, bid price is  $b_t$ .

- **State transition probability**,  $p = \mathbb{P}(s_{t+1}|s_t, b_t)$ . For a state  $B_t$  and bidding  $b_t$ , the transition of  $B_{t+1}$  is:

$$\begin{cases} B_{t+1} = B_t + c(\mathbf{x}_t, *)\mathbf{1}(w(\mathbf{x}_t), b_t), & \text{if } B_t < B \\ B_{t+1} = B_t, & \text{if } B_t \geq B \end{cases}.$$

The sequence of  $\mathbf{x}_t$  is an i.i.d. random process (see Section 1.1 ), that is, the transition of  $\mathbf{x}_{t+1}$  is not affected by action  $b_t$  and is fully described by  $p_{\mathbf{x}}$ . The above two transition rules together characterize the state transition probability.

- **Reward**,  $r = r(s_t, b_t)$ , which is slightly complicated because of the punishment of constraints,

$$r_t = v(\mathbf{x}_t)\mathbf{1}(w(\mathbf{x}_t), b_t) - \sum_{j=1}^k \lambda_j c(\mathbf{x}_t, *)\mathbf{1}(w(\mathbf{x}_t), b_t) - \lambda_j C_j v_{a,j}(\mathbf{x}_t)\mathbf{1}(w(\mathbf{x}_t), b_t)$$

where  $\lambda_j > 0, j \in \{1, \dots, k\}$  is the Lagrange multiplier that remakes the objective with the second constraint in (1) as punishment terms.

- **Time horizon**. The auction scale  $T \in \mathbb{Z}_+$  is the time horizon.

Furthermore, we assume that functions  $v, v_a, g$  and  $p_{\mathbf{x}}$  are known by the DSP prior to solving the bid optimization problem. In a real DSP, these functions mainly come from feature mapping and value estimation modules, such as CTR prediction [11, 12], CVR prediction [13, 14], winning price and winning rate prediction [15–18], landscape prediction [19], feature construction [6], and so on. They are important in RTB, however their prediction is beyond the scope of this paper. The bid optimization under the MDP model is to solve:

$$\arg \max_{\pi_t, t=1, \dots, T} \sum_{t=1}^T \mathbb{E}[r_t] \quad (3)$$

Note that the budget constraint has been encoded in the system dynamics, which does not explicitly appear in the above problem. Besides, in most previous optimization-based studies [2, 5, 8, 20, 21], researchers do not consider the influence of the dynamic consumed budget and thus solve a constraint optimization problem only with a fixed predefined total budget and CPA constraints. These methods can hold optimality but are fragile in the face real uncertainty and disturbance. On the other hand, in some machine-learning-based studies [22–26], budget dynamics have been considered in MDP modelling. Nevertheless, they need enormous resources and data to train the networks. Besides, the question about the optimality still exists. Compared to these two research approaches, we focus on solving the bid optimization problem using the MDP modelling together with optimization methods.

## 2 Rollout Mechanism and Structure of Bidding Functions

Though we can use the classical dynamic programming (DP) methodology to solve (3), it is still challenging because:

- (a) The MDP problem (3) is not completely equivalent to (1).
- (b) The impression feature  $\mathbf{x}_t$  is usually of high dimensions and sparse. It is inefficient to work out and then store the bidding strategy mapping the feature, among others, to the bid price.
- (c) On Taobao, half of the ad campaign received more than  $10^4$  bid requests a day; in other words, they bid every 10 seconds on average. Moreover, if we consider the distribution of impressions within a day, the frequency will be higher. High decision frequency and large-scale data place high demand on infrastructure and algorithms.
- (d) The rewards and transition of the consumed budget are time-delayed [3]. As the example in Section 1.2.1, an advertiser wins auctions and pushes ads in front of users. Whether the latter trades or not is usually not immediately known.

The above is in two parts: methodology and implementation. The methodology part is discussed in this section, which requires us to fill the model gaps and design an efficient solution. The implementation trade-off is discussed in the next section, which mainly comes from realistic constraints, performance, and practicality.

### 2.1 A Rollout Mechanism

Exact value iteration, policy iteration methods, or their variants tend to lead to overwhelming computation requirements as the state space is extremely large. Compromise formulas are various approximates of DPs (ADPs) [27–32]. We introduce an approximate method known as the open-loop feedback control (OLFC), which is a rollout mechanism [27, 32], to solve (3). The core concept of OLFC is to ignore in part the availability of online information and target a more tractable computation. Back to our bid optimization problem, we use the impression feature accessed from online as feedback information to adjust the current and future bid prices in real time and ignore the real budget remnant information. In particular, for the  $m$ -th auction, we solve the remaining  $(T - m + 1)$  number of bid price decision rules  $\bar{\pi}_t(\mathbf{x}_t)$  through solving the following

problem:

$$\begin{aligned}
& \arg \max_{\bar{\pi}_t, t=m, \dots, T} \sum_{t=m}^T \mathbb{E} [v(\mathbf{x}_t) \mathbf{1}(w(\mathbf{x}_t), \bar{\pi}_t(\mathbf{x}_t))] \\
& \quad - \sum_{j=1}^k \lambda_j \left\{ \mathbb{E} \left[ \sum_{t=m}^T c(\mathbf{x}_t, *) \mathbf{1}(w(\mathbf{x}_t), \bar{\pi}_t(\mathbf{x}_t)) \right] \right. \\
& \quad \left. - C_j \sum_{t=m}^T \mathbb{E} [v_{a,j}(\mathbf{x}_t) \mathbf{1}(w(\mathbf{x}_t), \bar{\pi}_t(\mathbf{x}_t))] \right\} \\
& \text{s.t. } \mathbb{E} \left[ \sum_{t=m}^T c(\mathbf{x}_t, *) \mathbf{1}(w(\mathbf{x}_t), \bar{\pi}_t(\mathbf{x}_t)) \right] \leq R_m.
\end{aligned} \tag{4}$$

Based on OLFC, remnant budget  $R_m = B - B_m$ , where  $B_m$  is the cost that has been paid for the past  $(m-1)$  auctions, for auction  $m$  to auction  $T$  is averaged with respect to  $x_t$  and randomness in auction outcomes and then constrained. Thus the altered strategy only relates to the impression feature  $\mathbf{x}$ . At the same time, (4) with different  $m$  and  $B_m$  have the same form. At each iteration  $m$ , we solve (4) for a specific  $R_m$ , and then evaluate the bid price through  $b_m = \bar{\pi}(\mathbf{x}_m)$ . We further simplify the formulation of (4) via the following operations. First, we consider taking expectation with respect to the randomness of auctions and with respect to  $\mathbf{x}_t$ 's separately. Consequently, given  $\mathbf{x}_t$ , we have

$$\mathbb{E}[\mathbf{1}(w(\mathbf{x}_t), \bar{\pi}_t(\mathbf{x}_t))] = \mathbb{P}(\text{Advertiser wins} | \mathbf{x}_t, \bar{\pi}_t) := q_t(\mathbf{x}_t, \bar{\pi}_t).$$

When  $\bar{\pi}_t$  is clearly known from context, we drop it, writing  $q_t(\mathbf{x}_t) := q_t(\mathbf{x}_t, \bar{\pi})$  for short. Letting  $q_t(\mathbf{x}_t)$  be the function to be determined, we come up with the following optimization problem.

$$\begin{aligned}
& \arg \max_{q_t, t=m, \dots, T} \sum_{t=m}^T \mathbb{E} [v(\mathbf{x}_t) q_t(\mathbf{x}_t)] \\
& \text{s.t. } \mathbb{E}[c(\mathbf{x}_t, *) q_t(\mathbf{x}_t)] \leq R_m \\
& \quad \frac{\sum_{t=m}^T \mathbb{E}[c(\mathbf{x}_t, *) q_t(\mathbf{x}_t)]}{\sum_{t=m}^T \mathbb{E}[v(\mathbf{x}_t) q_t(\mathbf{x}_t)]} \leq C_j, \quad j = 1, \dots, k.
\end{aligned} \tag{5}$$

Besides, the impression feature  $\mathbf{x}$  is i.i.d and draws from the same set  $\mathcal{X}$  subject to a PDF  $p_{\mathbf{x}}$ . It indicates that the quantity of the remnant budget  $B_m$  affects the strategy of each auction. If we could solve (5), we would reversely obtain  $\bar{\pi}(\mathbf{x}_t)$  from  $q_t(\mathbf{x}_t)$ . However it is extremely difficult to solve  $q_t(\mathbf{x}_t)$  from (5) as  $\mathcal{X}$  in general is a continuous set, solving  $q_t(\mathbf{x}_t)$  amounts to finding an optimal mapping from  $\mathcal{X}$  to  $[0, 1]$  for (5). To numerically solve (5) in a tractable way, we resort to a sampling method.

### 3 Implementation trade-off in solving the bid optimization

Though we know the structure of the bidding function (5), it is still difficult to solve. The reasons are not only the high frequency and time-delay mentioned in the implementation challenges, but also the difficulties of lacking the exact knowledge about  $w(\mathbf{x})$ ,  $p_{\mathbf{x}}$ , and other functions in (5). In this section, we solve the bid pricing rule numerically using historical auction logs, which transforms the problem of finding an optimal mapping  $q_t$  of  $\mathbf{x}_t$  for (5) into a linear programming problem with decision variable being a set of function values for  $q_t$  evaluated at archived impression features in an auction log. Further, we derive an approximate bid price decision rule to the one obtained from solving (4). This is done by analysing the structure of solutions to the linear programming and its dual. In particular, the pricing rule can be easily parameterised using an optimal solution to the dual linear programming problem. It is renewed over time by solving a sequence of receding dual linear programming problems. Finally, practical issues regarding balancing computational complexity and accuracy are also be discussed.

#### 3.1 Solution Analytics from Linear Programming Dual Theory

To come up with decision rules  $\bar{\pi}$ , we need to have the functions of  $p_{\mathbf{x}}$ ,  $v(\mathbf{x})$ , and other variables in (5) first. However, they are inaccessible at the start of the day because impressions would not have appeared yet. We have to assume that the impression features follow the same PDF  $p_{\mathbf{x}}$  and set  $\mathcal{X}$  in different days, which is widely used in previous studies [2, 5, 8, 24]. Thus, we can approximately identify the variables from the logs. Nevertheless, it is challenging to identify  $p_{\mathbf{x}}$  due to the construction of  $\mathbf{x}$ , a high-dimensional vector containing lots of information (refer to Section 1.1 and [3, 6]). Some machine learning methods can be used to learn  $p_{\mathbf{x}}$  and others [22, 25, 26], it may not be suitable for such a huge scale of ad campaigns due to the limited resources in practice as well as the online implementation.

An auction log record is a description of an auctioned impression in the past. Inspired by [5, 8], every day's auction features recorded in logs can be seen as a collection of samples from  $\mathcal{X}$ . Thus we replace the expectations in (5) with



their estimates, generated from logs and solve :

$$\begin{aligned}
& \arg \max_{\hat{q}_t(\hat{\mathbf{x}}_t), t=m, \dots, \hat{T}} \sum_{t=m}^{\hat{T}} v(\hat{\mathbf{x}}_t) \hat{q}_t(\hat{\mathbf{x}}_t) \\
& \text{s.t.} \quad \sum_{t=m}^{\hat{T}} \hat{c}_t \hat{q}_t(\hat{\mathbf{x}}_t) \leq R_m \\
& \quad \frac{\sum_{t=m}^{\hat{T}} \hat{c}_t \hat{q}_t(\hat{\mathbf{x}}_t)}{\sum_{t=m}^{\hat{T}} v_{a,j} \hat{q}_t(\hat{\mathbf{x}}_t)} \leq C_j, \quad j = 1, \dots, k,
\end{aligned} \tag{6}$$

where the superscript  $\wedge$  denotes the value recorded in logs. The auction log for an ad campaign is one-day based. All the auction records in a day form an auction log. Note that the total impression  $T$  for the current trading day may be different from logged value  $\hat{T}$ . Besides, as we mentioned in Section 1.1, the billing method used in this paper is click-based. It means that  $\hat{c}_t = \hat{w}_t \hat{p}^{\text{CTR}}(\hat{\mathbf{x}}_t)$  when we evaluate  $\mathbb{E}[c(\mathbf{x}_t, *) | \mathbf{x}_t]$  with logs, where  $\hat{p}^{\text{CTR}}(\hat{\mathbf{x}}_t)$  is the predicted click-through rate of impression  $\hat{\mathbf{x}}_t$  recorded in logs.

We analyse the structure of the solution to (6) and have the following lemma.

**Lemma 1** *If (6) is feasible, then there exists an optimal solution  $\hat{q}_t^*(\mathbf{x}_t), t = m, \dots, \hat{T}$ , of (6), which has at least  $(\hat{T} - m - k)$  number of 0 and 1.*

**Proof 1** *We rewrite (6) to an equivalent constrained linear programming by introducing slack variables to transform the inequality constraints into linear ones:*

$$\begin{aligned}
& \arg \max_{\hat{q}_t(\hat{\mathbf{x}}_t), a_t, d_j, g} \sum_{t=m}^{\hat{T}} v(\mathbf{x}_t) \hat{q}_t(\hat{\mathbf{x}}_t) \\
& \text{s.t.} \quad g + \sum_{t=m}^{\hat{T}} \hat{c}_t \hat{q}_t(\hat{\mathbf{x}}_t) = R_m \\
& \quad \hat{q}_t(\hat{\mathbf{x}}_t) + a_t = 1, \quad t = m, \dots, \hat{T} \\
& \quad d_j + \sum_{t=m}^{\hat{T}} \hat{q}_t(\hat{\mathbf{x}}_t) (\hat{c}_t - C_j v_{a,j}(\hat{\mathbf{x}}_t)) = 0, \quad j = 1, \dots, k \\
& \quad \hat{q}_t(\hat{\mathbf{x}}_t) \geq 0, a_t \geq 0, \quad t = m, \dots, \hat{T} \\
& \quad d_j \geq 0, \quad j = 1, \dots, k \\
& \quad g \geq 0,
\end{aligned} \tag{7}$$

which contains  $(2\hat{T} - 2m + k + 3)$  variables and  $(\hat{T} - m + k + 2)$  equality constraints. We can define the so-called basic feasible solutions (BFS) for the form of LP problems like (7). The BFSs of (7) have at least  $(\hat{T} - m + 1)$  number of 0. From [33, Theorem 2.7], if (6) is feasible, then (7) is feasible and there exists

an optimal solution that is a BFS for (7). Furthermore, note that, due to the constraint  $\hat{q}_t(\hat{\mathbf{x}}_t) + a_t = 1, \forall t = m, \dots, \hat{T}$ , if  $a_t = 0$  then  $\hat{q}_t(\hat{\mathbf{x}}_t) = 1$ . Thus, when we only consider  $\hat{q}_m, \dots, \hat{q}_{\hat{T}}$  of an optimal BFS as an optimal solution to (6), it has at least  $(\hat{T} - m - k)$  number of 0 and 1.  $\square$

We continue to investigate (6) from the dual theory. Slack  $\hat{q}$  to  $[0, 1]$  and consider the dual problem of (6),

$$\begin{aligned} \arg \min_{\alpha, \beta_j, \gamma_t} & R_m \alpha + \sum_{t=m}^{\hat{T}} \gamma_t \\ \text{s.t.} & \hat{c}_t \alpha + \sum_{j=1}^k (\hat{c}_t - C_j v_{a,j}(\hat{\mathbf{x}}_t)) \beta_j - v(\hat{\mathbf{x}}_t) + \gamma_t \geq 0, \\ & t = m, \dots, \hat{T}, \\ & \gamma_t \geq 0, t = m, \dots, \hat{T}, \\ & \beta_j \geq 0, j = 1, \dots, k, \\ & \alpha \geq 0. \end{aligned} \tag{8}$$

Let  $\alpha^*, \beta_1^*, \dots, \beta_k^*$  be the optimal to (8). We denote

$$\mu_t = \frac{1}{\hat{p}^{\text{CTR}}(\hat{\mathbf{x}}_t)} \frac{v(\hat{\mathbf{x}}_t) + \sum_{j=1}^k \beta_j^* C_j v_{a,j}(\hat{\mathbf{x}}_t)}{\alpha^* + \sum_{j=1}^k \beta_j^*} \tag{9}$$

**Lemma 2** Suppose that (6) is feasible. Consider the optimal solution  $\hat{q}_t^*(\mathbf{x}_t), t = m, \dots, \hat{T}$ , studied in Lemma 1. Then  $\hat{q}_t^*(\hat{\mathbf{x}}_t) = 0$  if  $\mu_t < \hat{w}_t$ , and  $\hat{q}_t^*(\hat{\mathbf{x}}_t) > 0$  if  $\mu_t > \hat{w}_t$ .

**Proof 2** Based on the complementary slackness [8, 34], the optimal solutions of (8) and (7) satisfy  $\forall t = m, \dots, \hat{T}$ :

$$\hat{q}_t^*(\hat{\mathbf{x}}_t) \left\{ \hat{c}_t \left( \alpha^* + \sum_{j=1}^k \beta_j^* \right) - v(\hat{\mathbf{x}}_t) - \sum_{j=1}^k C_j v_{a,j}(\hat{\mathbf{x}}_t) \beta_j^* + \gamma_t^* \right\} = 0, \tag{10}$$

$$(\hat{q}_t^* - 1) \gamma_t^* = 0, \tag{11}$$

where (11) comes from the complementary slackness of  $\hat{q} \in [0, 1]$ . Then substituting (9) and  $\hat{c}_t = \hat{w}_t \hat{p}_t^{\text{CTR}}$  into (10),

$$\hat{q}_t^*(\hat{\mathbf{x}}_t) \left\{ \hat{p}^{\text{CTR}}(\hat{\mathbf{x}}_t) (\hat{w}_t - \mu_t) \left( \alpha^* + \sum_{j=1}^k \beta_j^* \right) + \gamma_t^* \right\} = 0.$$

If  $\mu_t < \hat{w}_t$ , then  $\hat{p}^{\text{CTR}}(\hat{\mathbf{x}}_t) (\hat{w}_t - \mu_t) \left( \alpha^* + \sum_{j=1}^k \beta_j^* \right) + \gamma_t^* > 0$ , which further implies that  $\hat{q}_t^*(\hat{\mathbf{x}}_t) = 0$  by (10). Otherwise, if  $\hat{q}_t^*(\hat{\mathbf{x}}_t) = 0$ , then  $\gamma_t^* = 0$  from (11). Hence  $\mu_t \leq \hat{w}_t$  according to feasibility of (8), which completes the proof.  $\square$

Finally, we have the following theorem.

**Theorem 1** *Suppose that  $\hat{T} \gg k$  and that (6) is feasible. The following bidding rule of (6)*

$$\bar{\pi}_t(\hat{\mathbf{x}}_t) = \frac{1}{\hat{p}^{\text{CTR}}(\hat{\mathbf{x}}_t)} \frac{v(\hat{\mathbf{x}}_t) + \sum_{j=1}^k \beta_j^* C_j v_{a,j}(\hat{\mathbf{x}}_t)}{\alpha^* + \sum_{j=1}^k \beta_j^*}. \quad (12)$$

*is optimal in hindsight (It is in hindsight since at impression  $t$  the realizations  $\hat{\mathbf{x}}_t$  and  $\hat{w}_t$  from  $t = m$  to  $t = \hat{T}$  have already been known for bidding decision making.) with respect to the second-price auction mechanism.*

**Proof 3** *Consider the bidding strategy (12). If  $\bar{\pi}_t(\hat{\mathbf{x}}_t) > \hat{w}_t$ , then the bidder will win the impression according to the second-price auction mechanism. On the other hand,  $\hat{q}_t^*(\hat{\mathbf{x}}_t) > 0$  by Lemma 2. If  $\bar{\pi}_t(\hat{\mathbf{x}}_t) < \hat{w}_t$ , then the bidder will lose the impression by the second-price auction mechanism. In this case  $\hat{q}_t^*(\hat{\mathbf{x}}_t) = 0$  by Lemma 2. In virtue of Lemma 1, when  $\hat{q}_t^*(\hat{\mathbf{x}}_t) > 0$  there are at most  $(k+1)$  number of  $\hat{q}_t^*(\hat{\mathbf{x}}_t)$  that is not 1. Noticing that  $\hat{T} \gg k$ , it concludes the proof.  $\square$*

The assumption  $\hat{T} \gg k$  is practical in most cases as the number of real ad campaigns is mostly around or greater than the order of  $10^4$ . In contrast, the number of CPA constraints is at most one or two in most cases [1, 2, 5, 8, 9]. Thus, for a batch of  $\hat{T}$  impressions, when we compare auction results that come from an LP solver 1 and Theorem 1, the number of different results would not exceed  $k+1$ , which is an extremely small quantity compared to the total number of auctions  $\hat{T}$ .

We use the rule (12) as an approximate to the bid pricing rule obtained from solving (4) for the ongoing impression auction. Thus a bid price for the current auction can be proposed by valuing the current bid price using the current  $\mathbf{x}_t$  via (12), that is,

$$\bar{\pi}_t(\mathbf{x}_t) = \frac{1}{p_t^{\text{CTR}}(\mathbf{x}_t)} \frac{v(\mathbf{x}_t) + \sum_{j=1}^k \beta_j^* C_j v_{a,j}(\mathbf{x}_t)}{\alpha^* + \sum_{j=1}^k \beta_j^*}. \quad (13)$$

### 3.2 Auctions and Bid Requests Aggregation

In practice we may not exactly know the remnant budget at the start of each auction since the expense for each auction is in connection with users' clicks, which causes time delay [2, 8]. A possible solution is to aggregate impressions within a period of time and then fix  $\bar{\pi}$  in each period [2, 5, 8]. At the end of the period, the remnant budget is re-evaluated via the billing module. After that, we can renew  $\bar{\pi}$  in preparation for the next aggregation of impressions.

## References

- [1] Google.com, "Google ads help," 2022. [Online]. Available: <https://support.google.com/google-ads>

- [2] W. Zhang, “Optimal real-time bidding for display advertising,” Ph.D. dissertation, UCL (University College London), 2016.
- [3] S. Yuan, J. Wang, and X. Zhao, “Real-time bidding for online advertising: measurement and analysis,” in *Proceedings of the seventh international workshop on data mining for online advertising*, 2013, pp. 1–8.
- [4] W. Zhang, J. Qin, W. Guo, R. Tang, and X. He, “Deep learning for click-through rate estimation,” *arXiv preprint arXiv:2104.10584*, 2021.
- [5] Z. Han, J. Jin, T. Chang, P. Fei, and K. Gai, “Optimized cost per click in taobao display advertising,” in *the 23rd ACM SIGKDD International Conference*, 2017.
- [6] W. Zhang and J. Xu, “Learning, prediction and optimisation in rtb display advertising,” in *Proceedings of the 25th Information and Knowledge Management*, 2014, pp. 1077–1086.
- [7] W. Vickrey, “Counterspeculation, auctions, and competitive sealed tenders,” *The Journal of finance*, vol. 16, no. 1, pp. 8–37, 1961.
- [8] X. Yang, Y. Li, H. Wang, D. Wu, Q. Tan, J. Xu, and K. Gai, “Bid optimization by multivariable control in display advertising,” in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 1966–1974.
- [9] Y. Yuan, F. Wang, J. Li, and R. Qin, “A survey on real time bidding advertising,” in *Proceedings of 2014 IEEE International Conference on Service Operations and Logistics, and Informatics*. IEEE, 2014, pp. 418–423.
- [10] F. Lyu, X. Tang, H. Guo, R. Tang, X. He, R. Zhang, and X. Liu, “Memorize, factorize, or be naïve: Learning optimal feature interaction methods for ctr prediction,” in *2022 IEEE 38th International Conference on Data Engineering (ICDE)*. IEEE, 2022, pp. 1450–1462.
- [11] H. Guo, R. Tang, Y. Ye, Z. Li, and X. He, “Deepfm: a factorization-machine based neural network for ctr prediction,” *arXiv preprint arXiv:1703.04247*, 2017.
- [12] G. Zhou, X. Zhu, C. Song, Y. Fan, H. Zhu, X. Ma, Y. Yan, J. Jin, H. Li, and K. Gai, “Deep interest network for click-through rate prediction,” in *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, 2018, pp. 1059–1068.
- [13] H. Wen, J. Zhang, F. Lv, W. Bao, T. Wang, and Z. Chen, “Hierarchically modeling micro and macro behaviors via multi-task learning for conversion rate prediction,” in *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2021, pp. 2187–2191.

- [14] H. Wen, J. Zhang, Q. Lin, K. Yang, and P. Huang, “Multi-level deep cascade trees for conversion rate prediction in recommendation system,” in *proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 338–345.
- [15] W. Wu, M.-Y. Yeh, and M.-S. Chen, “Deep censored learning of the winning price in the real time bidding,” in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 2526–2535.
- [16] X. Li and D. Guan, “Programmatic buying bidding strategies with win rate and winning price estimation in real time mobile advertising,” in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 2014, pp. 447–460.
- [17] W.-Y. Shih, Y.-S. Lu, H.-P. Tsai, and J.-L. Huang, “An expected win rate-based real time bidding strategy for branding campaign by the model-free reinforcement learning model,” *IEEE Access*, vol. 8, pp. 151 952–151 967, 2020.
- [18] W.-Y. Zhu, W.-Y. Shih, Y.-H. Lee, W.-C. Peng, and J.-L. Huang, “A gamma-based regression for winning price estimation in real-time bidding advertising,” in *2017 IEEE International Conference on Big Data (Big Data)*. IEEE, 2017, pp. 1610–1619.
- [19] Y. Wang, K. Ren, W. Zhang, J. Wang, and Y. Yu, “Functional bid landscape forecasting for display advertising,” in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2016, pp. 115–131.
- [20] N. Karlsson, “Feedback control in programmatic advertising: The frontier of optimization in real-time bidding,” *IEEE control systems*, vol. 40, no. 5, pp. 40–77, 2020.
- [21] S. Tunuguntla and P. R. Hoban, “A near-optimal bidding strategy for real-time display advertising auctions,” *Journal of Marketing Research*, vol. 58, no. 1, pp. 1–21, 2021.
- [22] H. Cai, K. Ren, W. Zhang, K. Malialis, J. Wang, Y. Yu, and D. Guo, “Real-time bidding by reinforcement learning in display advertising,” in *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*, 2017, pp. 661–670.
- [23] N. Grislain, N. Perrin, and A. Thabault, “Recurrent neural networks for stochastic control in real-time bidding,” in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 2801–2809.

- [24] Y. He, X. Chen, D. Wu, J. Pan, Q. Tan, C. Yu, J. Xu, and X. Zhu, “A unified solution to constrained bidding in online display advertising,” in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021, pp. 2993–3001.
- [25] J. Jin, C. Song, H. Li, K. Gai, J. Wang, and W. Zhang, “Real-time bidding with multi-agent reinforcement learning in display advertising,” in *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. New York, NY, USA: Association for Computing Machinery, 2018, p. 2193–2201.
- [26] K. Ren, W. Zhang, K. Chang, Y. Rong, Y. Yu, and J. Wang, “Bidding machine: Learning to bid for directly optimizing profits in display advertising,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 30, no. 4, pp. 645–659, 2017.
- [27] D. P. Bertsekas, “Dynamic programming and optimal control,” *Athena Scientific*, 1995.
- [28] —, “Dynamic programming and stochastic control.” 1976.
- [29] G. Tesauro and G. Galperin, “On-line policy improvement using monte-carlo search,” *Advances in Neural Information Processing Systems*, vol. 9, 1996.
- [30] S. Boyd, M. T. Mueller, B. O’Donoghue, Y. Wang *et al.*, “Performance bounds and suboptimal policies for multi-period investment,” *Foundations and Trends® in Optimization*, vol. 1, no. 1, pp. 1–72, 2013.
- [31] M. W. Ulmer, *Approximate dynamic programming for dynamic vehicle routing*. Springer, 2017, vol. 1.
- [32] S. E. Dreyfus, “Dynamic programming and the calculus of variations,” *Journal of mathematical analysis and applications*, vol. 1, no. 2, pp. 228–239, 1960.
- [33] A. Schrijver, *Theory of linear and integer programming*. John Wiley & Sons, 1998.
- [34] D. Bertsimas and J. N. Tsitsiklis, *Introduction to linear optimization*. Athena Scientific Belmont, MA, 1997, vol. 6.