

# MNL-Bandit: A Dynamic Learning Approach to Assortment Selection

Shipra Agrawal

Industrial Engineering and Operations Research, Columbia University, New York, NY. sa3305@columbia.edu

Vashist Avadhanula

Decision Risk and Operations, Columbia Business School, New York, NY. vavadhanula18@gsb.columbia.edu

Vineet Goyal

Industrial Engineering and Operations Research, Columbia University, New York, NY. vg2277@columbia.edu

Assaf Zeevi

Decision Risk and Operations, Columbia Business School, New York, NY. assaf@gsb.columbia.edu

We consider a dynamic assortment selection problem, where in every round the retailer offers a subset (assortment) of  $N$  substitutable products to a consumer, who selects one of these products according to a multinomial logit (MNL) choice model. The retailer observes this choice and the objective is to dynamically learn the model parameters, while optimizing cumulative revenues over a selling horizon of length  $T$ . We refer to this exploration-exploitation formulation as the *MNL-Bandit problem*. Existing methods for this problem follow an *explore-then-exploit* approach, which estimate parameters to a desired accuracy and then, treating these estimates as if they are the correct parameter values, offers the optimal assortment based on these estimates. These approaches require certain a priori knowledge of “separability,” determined by the true parameters of the underlying MNL model, and this in turn is critical in determining the length of the exploration period. (Separability refers to the distinguishability of the true optimal assortment from the other sub-optimal alternatives.) In this paper, we give an efficient algorithm that *simultaneously* explores and exploits, achieving performance independent of the underlying parameters. The algorithm can be implemented in a fully online manner, without knowledge of the horizon length  $T$ . Furthermore, the algorithm is adaptive in the sense that its performance is near-optimal in both the “well separated” case, as well as the general parameter setting where this separation need not hold.

*Key words:* Exploration-Exploitation, assortment optimization, upper confidence bound, multinomial logit

## 1. Introduction

### 1.1. Overview of the problem.

Assortment optimization problems arise widely in many industries including retailing and online advertising where the seller needs to select a subset of items to offer from a universe of substitutable items such that the expected revenue is maximized. Choice models capture substitution effects

among products by specifying the probability that a consumer selects a product given the offered set. Traditionally, the assortment decisions are made at the start of the selling period based on the estimated choice model from historical data; see Kök and Fisher (2007) for a detailed review.

In this work, we focus on the dynamic version of the problem where the retailer needs to simultaneously learn consumer preferences and maximize revenue. In many business applications such as fast fashion and online retail, new products can be introduced or removed from the offered assortments in a fairly frictionless manner and the selling horizon for a particular product can be short. Therefore, the traditional approach of first estimating the choice model and then using a static assortment based on the estimates, is not practical in such settings. Rather, it is essential to experiment with different assortments to learn consumer preferences, while simultaneously attempting to maximize immediate revenues. Suitable balancing of this exploration-exploitation tradeoff is the focal point of this paper.

We consider a stylized dynamic optimization problem that captures some salient features of this application domain, where our goal is to develop an exploration-exploitation policy that simultaneously learns from current observations and exploits this information gain for future decisions. In particular, we consider a constrained assortment selection problem under the Multinomial logit (MNL) model with  $N$  substitutable products and a “no purchase” option. Our goal is to offer a sequence of assortments,  $S_1, \dots, S_T$ , where  $T$  is the planning horizon, such that the cumulative expected revenues over said horizon is maximized, or alternatively, minimizing the gap between the performance of a proposed policy and that of an oracle that knows instance parameters a priori, a quantity referred to as the *regret*.

**Related literature.** The Multinomial Logit model (MNL), owing primarily to its tractability, is the most widely used choice model for assortment selection problems. (The model was introduced independently by Luce (1959) and Plackett (1975), see also Train (2009), McFadden (1978), Ben-Akiva and Lerman (1985) for further discussion and survey of other commonly used choice models.) If the consumer preferences (MNL parameters in our setting) are known a priori, then the problem of computing the optimal assortment, which we refer to as the *static assortment optimization problem*, is well studied. Talluri and van Ryzin (2004) consider the unconstrained assortment planning problem under the MNL model and present a greedy approach to obtain the optimal assortment. Recent works of Davis et al. (2013) and Désir and Goyal (2014) consider assortment planning problems under MNL with various constraints. Other choice models such as Nested Logit (Williams (1977), Davis et al. (2011), Gallego and Topaloglu (2014) and Li et al. (2015)), Markov Chain (Blanchet et al. (2016) and Désir et al. (2015)) and more general models (Farias et al. (2013) and Gallego et al. (2014)) are also considered in the literature.

Most closely related to our work are the papers of Caro and Gallien (2007), Rusmevichientong et al. (2010) and Sauré and Zeevi (2013), where information on consumer preferences is not known and needs to be learned over the course of the selling horizon. Caro and Gallien (2007) consider the setting under which demand for products is independent of each other. Rusmevichientong et al. (2010) and Sauré and Zeevi (2013) consider the problem of minimizing regret under the MNL choice model and present an “explore first and exploit later” approach. In particular, a selected set of assortments are explored until parameters can be estimated to a desired accuracy and then the optimal assortment corresponding to the estimated parameters is offered for the remaining selling horizon. The exploration period depends on certain a priori knowledge about instance parameters. Assuming that the optimal and next-best assortment are “well separated,” they show an asymptotic  $O(N \log T)$  regret bound. However, their algorithm relies crucially on the a priori knowledge of certain instance parameters which is not readily available in practice. Furthermore, their policies also require a priori knowledge of the length of the planning horizon. In this work, we focus on approaches that simultaneously explore and exploit demand information and do not require any such a priori knowledge or assumptions; thereby, making our approach more universal in its scope.

Our problem is closely related to the multi-armed bandit (MAB) paradigm. A naive mapping to that setting would consider every assortment as an arm, and as such, would lead to exponentially many arms. Popular extensions of MAB for large scale problems include the linear bandit (e.g., Auer (2003), Rusmevichientong and Tsitsiklis (2010)) and generalized linear bandit (Filippi et al. (2010)) problems. However, these do not apply directly to our problem, since the revenue corresponding to an assortment is nonlinear in problem parameters. Other works (see Chen et al. (2013)) have considered versions of MAB where one can play a subset of arms in each round and the expected reward is a function of rewards for the arms played. However, this approach assumes that the reward for each arm is generated independently of the other arms in the subset. This is not the case typically in retail settings, and in particular, in the MNL choice model where purchase decisions depend on the assortment of products offered in a time step. In this work, we use the structural properties of the MNL model, along with techniques from MAB literature, to optimally explore and exploit in the presence of a large number of alternatives (assortments).

## 1.2. Contributions

**Parameter independent online algorithm and regret bounds.** We give an efficient online algorithm that judiciously balances the exploration and exploitation trade-off intrinsic to our problem and achieves a worst-case regret bound of  $O(\sqrt{NT \log T})$  under a mild assumption, namely that the no-purchase is the most “frequent” outcome. Our algorithm is online in the sense that it

does not require any prior knowledge of the instance parameters or the time horizon,  $T$ . Moreover, the regret bound we present for this algorithm is non-asymptotic and parameter independent. The “big-oh” notation is used for brevity and only hides absolute constants. To the best of our knowledge, this is the first such policy to have a parameter independent regret bound for the MNL choice model. The assumption regarding no-purchase is quite natural and a norm in online retailing for example. Furthermore, we can establish a similar regret bound when this assumption is relaxed.

We also show that for “well separated” instances, the regret of our policy is bounded by  $O\left(\min\left(N^2 \log T / \Delta, \sqrt{NT}\right)\right)$  where  $\Delta$  is the “separability” parameter. This is comparable to the bounds established in Sauré and Zeevi (2013) and Rusmevichientong et al. (2010), even though we do not require any prior information on  $\Delta$  unlike the aforementioned work. It is also interesting to note that the regret bounds hold true for a large class of constraints, e.g., we can handle matroid constraints such as assignment, partition and more general totally unimodular constraints (see Davis et al. (2013)). Our algorithm is predicated on upper confidence bound (UCB) type logic, originally developed to balance the aforementioned exploration-exploitation trade-off in the context of the multi-armed bandit (MAB) problem (cf. Lai and Robbins (1985) and Auer et al. (2002)). In this paper the UCB approach, also known as optimism in the face of uncertainty, is customized to the assortment optimization problem under the MNL model.

**Lower bounds.** We establish a non-asymptotic lower bound for the online assortment optimization problem under the MNL model. In particular, we show that for the cardinality constrained problem under the MNL model, any algorithm must incur a regret of  $\Omega(\sqrt{NT/K})$ , where  $K$  is the bound on the number of products that can be offered in an assortment. This bound is derived via a reduction from a parametric multi-armed bandit problem, for which such lower bounds are constructed by means of information theoretic arguments. This result establishes that our online algorithm is nearly optimal, the upper bound being within a factor of  $\sqrt{K}$  of the lower bound.

**Computational study.** We present a computational study that highlights several salient features of our algorithm. In particular, we test the performance of our algorithm over instances with varying degrees of separability between optimal and sub-optimal solutions and observe that the performance is bounded irrespective of the “separability parameter”. In contrast, the approach of Sauré and Zeevi (2013) “breaks down” and results in linear regret for some values of the “separability parameter.” We also present results from a simulated study on a real world data set.

**Outline.** In Section 2, we give the precise problem formulation. In Section 3, we present our algorithm for the MNL-Bandit problem, and in Section 4, we prove the worst-case regret bound of  $\tilde{O}(\sqrt{NT})$  for our policy. We present the modified algorithm that works for a more general class of

MNL parameters and establish  $\tilde{O}(\sqrt{BNT})$  regret bounds in Section 5. In Section 6, we present the logarithmic regret bound for our policy for the case of “well separated” instances. In Section 7, we present our non-asymptotic lower bound on regret for any algorithm for MNL-Bandit. In Section 8, we present results from our computational study.

## 2. Problem formulation

**The basic assortment problem.** In our problem, at every time instance  $t$ , the seller selects an assortment  $S_t \subset \{1, \dots, N\}$  and observes the customer purchase decision  $C_t \in S_t \cup \{0\}$ . As noted earlier, we assume consumer preferences are modeled using a multinomial logit (MNL) model. Under this model, the probability that a consumer purchases product  $i$  at time  $t$  when offered an assortment  $S_t = S \subset \{1, \dots, N\}$  is given by,

$$p_i(S) := \mathbb{P}(C_t = i | S_t = S) = \begin{cases} \frac{v_i}{v_0 + \sum_{j \in S} v_j}, & \text{if } i \in S \cup \{0\} \\ 0, & \text{otherwise,} \end{cases} \quad (2.1)$$

for all  $t$ , where  $v_i$  is the *attraction parameter* for product  $i$  in the MNL model. The random variables  $\{C_t : t = 1, 2, \dots\}$  are conditionally independent, namely,  $C_t$  conditioned on the event  $\{S_t = S\}$  is independent of  $C_1, \dots, C_{t-1}$ . Without loss of generality, we can assume that  $v_0 = 1$ . It is important to note that the parameters of the MNL model  $v_i$ , are not known to the seller. From (2.1), the expected revenue when assortment  $S$  is offered and the MNL parameters are denoted by the vector  $\mathbf{v}$  is given by

$$R(S, \mathbf{v}) = \mathbb{E} \left[ \sum_{i \in S} r_i \mathbb{1}\{C_t = i | S_t = S\} \right] = \sum_{i \in S} \frac{r_i v_i}{1 + \sum_{j \in S} v_j}. \quad (2.2)$$

We consider several naturally arising constraints over the assortments that the retailer can offer. These include cardinality constraints (where there is an upper bound on the number of products that can be offered in the assortment), partition matroid constraints (where the products are partitioned into segments and the retailer can select at most a specified number of products from each segment) and joint display and assortment constraints (where the retailer needs to decide both the assortment as well as the display segment of each product in the assortment and there is an upper bound on the number of products in each display segment). More generally, we consider the set of totally unimodular (TU) constraints on the assortments. Let  $\mathbf{x}(S) \in \{0, 1\}^N$  be the incidence vector for assortment  $S \subseteq \{1, \dots, N\}$ , i.e.,  $x_i(S) = 1$  if product  $i \in S$  and 0 otherwise. We consider constraints of the form

$$\mathcal{S} = \{S \subseteq \{1, \dots, N\} \mid \mathbf{A} \mathbf{x}(S) \leq \mathbf{b}, \mathbf{0} \leq \mathbf{x} \leq \mathbf{1}\}, \quad (2.3)$$

where  $\mathbf{A}$  is a totally unimodular matrix. The totally unimodular constraints model a rich class of practical assortment planning problems including the examples discussed above. We refer the

reader to Davis et al. (2013) for a detailed discussion on assortment and pricing optimization problems that can be formulated under the TU constraints.

**Admissible Policies.** To define the set of policies that can be used by the seller, let  $U$  be a random variable, which encodes any additional sources of randomization and  $(\mathbb{U}, \mathcal{U}, \mathbb{P}_u)$  be the corresponding probability space. We define  $\{\pi_t, t = 1, 2, \dots\}$  to be a measurable mapping as follows:

$$\begin{aligned}\pi_1 &: \mathbb{U} \rightarrow \mathcal{S} \\ \pi_t &: \mathbb{U} \times \mathcal{S}^{t-1} \times \{0, \dots, N\}^{t-1} \rightarrow \mathcal{S}, \text{ for each } t = 2, 3, \dots\end{aligned}$$

where  $\mathcal{S}$  is as defined in (2.3). Then the assortment selection for the seller at time  $t$  is given by

$$S_t = \begin{cases} \pi_1(U), & t = 1 \\ \pi_t(U, C_1, \dots, C_{t-1}, S_1, \dots, S_{t-1}), & t = 2, 3, \dots, \end{cases} \quad (2.4)$$

For further reference, let  $\{\mathcal{H}_t : t = 1, 2, \dots\}$  denote the filtration associated with the policy  $\pi = (\pi_1, \pi_2, \dots, \pi_t, \dots)$ . Specifically,

$$\begin{aligned}\mathcal{H}_1 &= \sigma(U) \\ \mathcal{H}_t &= \sigma(U, C_1, \dots, C_{t-1}, S_1, \dots, S_{t-1}), \text{ for each } t = 2, 3, \dots\end{aligned}$$

We denote by  $\mathbb{P}_\pi\{\cdot\}$  and  $\mathbb{E}_\pi\{\cdot\}$  the probability distribution and expectation value over path space induced by the policy  $\pi$ .

**The online assortment optimization problem.** The objective is to design a policy  $\pi = (\pi_1, \dots, \pi_T)$  that selects a sequence of history dependent assortments  $(S_1, S_2, \dots, S_T)$  so as to maximize the cumulative expected revenue,

$$\mathbb{E}_\pi \left( \sum_{t=1}^T R(S_t, \mathbf{v}) \right), \quad (2.5)$$

where  $R(S, \mathbf{v})$  is defined as in (2.2). Direct analysis of (2.5) is not tractable given that the parameters  $\{v_i, i = 1, \dots, N\}$  are not known to the seller a priori. Instead we propose to measure the performance of a policy  $\pi$  via its *regret*. The objective then is to design a policy that approximately minimizes the *regret* defined as

$$Reg_\pi(T, \mathbf{v}) = \sum_{t=1}^T R(S^*, \mathbf{v}) - \mathbb{E}_\pi[R(S_t, \mathbf{v})], \quad (2.6)$$

where  $S^*$  is the optimal assortment for (2.2), namely,  $S^* = \operatorname{argmax}_{S \in \mathcal{S}} R(S, \mathbf{v})$ . This exploration-exploitation problem, which we refer to as **MNL-Bandit**, is the focus of this paper.

### 3. The proposed policy

In this section, we describe our proposed policy for the MNL-Bandit problem. The policy is designed using the characteristics of the MNL model based on the principle of optimism under uncertainty.

### 3.1. Challenges and overview

A key difficulty in applying standard multi-armed bandit techniques to this problem is that the response observed on offering a product  $i$  is *not* independent of other products in assortment  $S$ . Therefore, the  $N$  products cannot be directly treated as  $N$  independent arms. As mentioned before, a naive extension of MAB algorithms for this problem would treat each of the feasible assortments as an arm, leading to a computationally inefficient algorithm with exponential regret. Our policy utilizes the specific properties of the dependence structure in MNL model to obtain an efficient algorithm with order  $\sqrt{NT}$  regret.

Our policy is based on a non-trivial extension of the UCB algorithm Auer et al. (2002). It uses the past observations to maintain increasingly accurate upper confidence bounds for the MNL parameters  $\{v_i, i = 1, \dots, N\}$ , and uses these to (implicitly) maintain an estimate of expected revenue  $R(S)$  for every feasible assortment  $S$ . In every round, our policy picks the assortment  $S$  with the highest optimistic revenue. There are two main challenges in implementing this scheme. First, the customer response to being offered an assortment  $S$  depends on the entire set  $S$ , and does not directly provide an unbiased sample of demand for a product  $i \in S$ . In order to obtain unbiased estimates of  $v_i$  for all  $i \in S$ , we offer a set  $S$  multiple times: specifically, it is offered repeatedly until a no-purchase occurs. We show that proceeding in this manner, the average number of times a product  $i$  is purchased provides an unbiased estimate of the parameter  $v_i$ . The second difficulty is the computational complexity of maintaining and optimizing revenue estimates for each of the exponentially many assortments. To this end, we use the structure of the MNL model and define our revenue estimates such that the assortment with maximum estimated revenue can be efficiently found by solving a simple optimization problem. This optimization problem turns out to be a static assortment optimization problem with upper confidence bounds for  $v_i$ 's as the MNL parameters, for which efficient solution methods are available.

### 3.2. Details of the policy

We divide the time horizon into epochs, where in each epoch we offer an assortment repeatedly until a no purchase outcome occurs. Specifically, in each epoch  $\ell$ , we offer an assortment  $S_\ell$  repeatedly. Let  $\mathcal{E}_\ell$  denote the set of consecutive time steps in epoch  $\ell$ .  $\mathcal{E}_\ell$  contains all time steps after the end of epoch  $\ell - 1$ , until a no-purchase happens in response to offering  $S_\ell$ , including the time step at which no-purchase happens. The length of an epoch  $|\mathcal{E}_\ell|$  conditioned on  $S_\ell$  is a geometric random variable with success probability defined as the probability of no-purchase in  $S_\ell$ . The total number of epochs  $L$  in time  $T$  is implicitly defined as the minimum number for which  $\sum_{\ell=1}^L |\mathcal{E}_\ell| \geq T$ .

At the end of every epoch  $\ell$ , we update our estimates for the parameters of MNL, which are used in epoch  $\ell + 1$  to choose assortment  $S_{\ell+1}$ . For any time step  $t \in \mathcal{E}_\ell$ , let  $c_t$  denote the consumer's

response to  $S_\ell$ , i.e.,  $c_t = i$  if the consumer purchased product  $i \in S$ , and 0 if no-purchase happened. We define  $\hat{v}_{i,\ell}$  as the number of times a product  $i$  is purchased in epoch  $\ell$ .

$$\hat{v}_{i,\ell} := \sum_{t \in \mathcal{E}_\ell} \mathbb{1}(c_t = i) \quad (3.1)$$

For every product  $i$  and epoch  $\ell \leq L$ , we keep track of the set of epochs before  $\ell$  that offered an assortment containing product  $i$ , and the number of such epochs. We denote the set of epochs by  $\mathcal{T}_i(\ell)$  and the number of epochs by  $T_i(\ell)$ . That is,

$$\mathcal{T}_i(\ell) = \{\tau \leq \ell \mid i \in S_\tau\}, \quad T_i(\ell) = |\mathcal{T}_i(\ell)|. \quad (3.2)$$

We compute  $\bar{v}_{i,\ell}$  as the number of times product  $i$  was purchased per epoch,

$$\bar{v}_{i,\ell} = \frac{1}{T_i(\ell)} \sum_{\tau \in \mathcal{T}_i(\ell)} \hat{v}_{i,\tau}. \quad (3.3)$$

We show that for all  $i \in S_\ell$ ,  $\hat{v}_{i,\ell}$  and  $\bar{v}_{i,\ell}$  are unbiased estimators of the MNL parameter  $v_i$  (see Lemma A.1). Using these estimates, we compute the upper confidence bounds,  $v_{i,\ell}^{\text{UCB}}$  for  $v_i$  as,

$$v_{i,\ell}^{\text{UCB}} := \bar{v}_{i,\ell} + \sqrt{\bar{v}_{i,\ell} \frac{48 \log(\ell+1)}{T_i(\ell)}} + \frac{48 \log(\ell+1)}{T_i(\ell)}. \quad (3.4)$$

We establish that  $v_{i,\ell}^{\text{UCB}}$  is an upper confidence bound on the true parameter  $v_i$ , i.e.,  $v_{i,\ell}^{\text{UCB}} \geq v_i$ , for all  $i, \ell$  with high probability (see Lemma 4.1). The role of the upper confidence bounds is akin to their role in hypothesis testing; they ensure that the likelihood of identifying the parameter value is sufficiently large. We then offer the optimistic assortment in the next epoch, based on the previous updates as follows,

$$S_{\ell+1} := \operatorname{argmax}_{S \in \mathcal{S}} \max \{R(S, \hat{\mathbf{v}}) : \hat{v}_i \leq v_{i,\ell}^{\text{UCB}}\} \quad (3.5)$$

where  $R(S, \hat{\mathbf{v}})$  is as defined in (2.2). We later show that the above optimization problem is equivalent to the following optimization problem (see Lemma A.3).

$$S_{\ell+1} := \operatorname{argmax}_{S \in \mathcal{S}} \tilde{R}_{\ell+1}(S), \quad (3.6)$$

where  $\tilde{R}_{\ell+1}(S)$  is defined as,

$$\tilde{R}_{\ell+1}(S) := \frac{\sum_{i \in S} r_i v_{i,\ell}^{\text{UCB}}}{1 + \sum_{j \in S} v_{j,\ell}^{\text{UCB}}}. \quad (3.7)$$

We summarize the steps in our policy in Algorithm 1. Finally, we may remark on the computational complexity of implementing (3.5). The optimization problem (3.5) is formulated as a static assortment optimization problem under the MNL model with TU constraints, with model parameters being  $v_{i,\ell}^{\text{UCB}}, i = 1, \dots, N$  ((3.6)). There are efficient polynomial time algorithms to solve the static assortment optimization problem under MNL model with known parameters (see Davis et al. (2013), Rusmevichientong et al. (2010)).



---

**Algorithm 1** Exploration-Exploitation algorithm for MNL-Bandit

---

```

1: Initialization:  $v_{i,0}^{\text{UCB}} = 1$  for all  $i = 1, \dots, N$ .
2:  $t = 1$  ;  $\ell = 1$  keeps track of the time steps and total number of epochs respectively
3: while  $t < T$  do
4:   Compute  $S_\ell = \operatorname{argmax}_{S \in \mathcal{S}} \tilde{R}_\ell(S) = \frac{\sum_{i \in S} r_i v_{i,\ell-1}^{\text{UCB}}}{1 + \sum_{j \in S} v_{j,\ell-1}^{\text{UCB}}}$ 
5:   Offer assortment  $S_\ell$ , observe the purchasing decision,  $\mathbf{c}_t$  of the consumer
6:   if  $c_t = 0$  then
7:     compute  $\hat{v}_{i,\ell} = \sum_{t \in \mathcal{E}_\ell} \mathbb{1}(c_t = i)$ , no. of consumers who preferred  $i$  in epoch  $\ell$ , for all  $i \in S_\ell$ .
8:     update  $\mathcal{T}_i(\ell) = \{\tau \leq \ell \mid i \in S_\tau\}$ ,  $T_i(\ell) = |\mathcal{T}_i(\ell)|$ , no. of epochs until  $\ell$  that offered product  $i$ .
9:     update  $\bar{v}_{i,\ell} = \frac{1}{T_i(\ell)} \sum_{\tau \in \mathcal{T}_i(\ell)} \hat{v}_{i,\tau}$ , sample mean of the estimates
10:    update  $v_{i,\ell}^{\text{UCB}} = \bar{v}_{i,\ell} + \sqrt{\bar{v}_{i,\ell} \frac{48 \log(\ell+1)}{T_i(\ell)}} + \frac{48 \log(\ell+1)}{T_i(\ell)}$ ;  $\ell = \ell + 1$ 
11:   else
12:      $\mathcal{E}_\ell = \mathcal{E}_\ell \cup t$ , time indices corresponding to epoch  $\ell$ .
13:   end if
14:    $t = t + 1$ 
15: end while

```

---

## 4. Main results

**Assumption 4.1** We make the following assumptions.

1. The MNL parameter corresponding to any product  $i \in \{1, \dots, N\}$  satisfies  $v_i \leq v_0 = 1$ .
2. The family of assortments  $\mathcal{S}$  is such that  $S \in \mathcal{S}$  and  $Q \subseteq S$  implies that  $Q \in \mathcal{S}$ .

The first assumption is equivalent to the ‘no purchase option’ being preferable to any other product. We note that this holds in many realistic settings, in particular, in online retailing and online display-based advertising. The second assumption implies that removing a product from a feasible assortment preserves feasibility. This holds for most constraints arising in practice including cardinality, and matroid constraints more generally. We would like to note that the first assumption is made for ease in presenting the key results and is not central to deriving bounds on the regret. In section 5, we relax this assumption and derive regret bounds that hold for any parameter instance.

Our main result is the following upper bound on the regret of the policy stated in Algorithm 1.

**Theorem 1** For any instance  $\mathbf{v} = (v_0, \dots, v_N)$  of the MNL-Bandit problem with  $N$  products,  $r_i \in [0, 1]$  under Assumption 4.1, the regret of the policy given by Algorithm 1 at time  $T$  is bounded as,

$$\text{Reg}_\pi(T, \mathbf{v}) \leq C_1 \sqrt{NT \log T} + C_2 N \log^2 T,$$

where  $C_1$  and  $C_2$  are absolute constants (independent of problem parameters).

#### 4.1. Proof Outline

In this section, we provide an outline of different steps involved in proving Theorem 1.

**Confidence intervals.** The first step in our regret analysis is to prove the following two properties of the estimates  $v_{i,\ell}^{UCB}$  computed as in (3.4) for each product  $i$ . Specifically, that  $v_i$  is bounded by  $v_{i,\ell}^{UCB}$  with high probability, and that as a product is offered more and more times, the estimates  $v_{i,\ell}^{UCB}$  converge to the true value with high probability. Intuitively, these properties establish  $v_{i,\ell}^{UCB}$  as upper confidence bounds converging to actual parameters  $v_i$ , akin to the upper confidence bounds used in the UCB algorithm for MAB in Auer et al. (2002). We provide the precise statements for the above mentioned properties in Lemma 4.1. These properties follow from an observation that is conceptually equivalent to the IIA (Independence of Irrelevant Alternatives) property of MNL. This observation shows that in each epoch  $\tau$ ,  $\hat{v}_{i,\tau}$  (the number of purchases of product  $i$ ) provides an independent unbiased estimates of  $v_i$ . Intuitively,  $\hat{v}_{i,\tau}$  is the ratio of probabilities of purchasing product  $i$  to preferring product 0 (no-purchase), which is independent of  $S_\tau$ . This also explains why we choose to offer  $S_\tau$  repeatedly until no-purchase occurs. Given these unbiased i.i.d. estimates from every epoch  $\tau$  before  $\ell$ , we apply a multiplicative Chernoff-Hoeffding bound to prove concentration of  $\bar{v}_{i,\ell}$ .

**Correctness of the optimistic assortment.** The product demand estimates  $v_{i,\ell-1}^{UCB}$  were used in (3.7) to define expected revenue estimates  $\tilde{R}_\ell(S)$  for every set  $S$ . In the beginning of every epoch  $\ell$ , Algorithm 1 computes the optimistic assortment as  $S_\ell = \arg \max_S \tilde{R}_\ell(S)$ , and then offers  $S_\ell$  repeatedly until no-purchase happens. The next step in the regret analysis is to use above properties of  $v_{i,\ell}^{UCB}$  to prove similar, though slightly weaker, properties for the estimates  $\tilde{R}_\ell(S)$ . First, we show that estimated revenue is an upper confidence bound on the optimal revenue, i.e.  $R(S^*, \mathbf{v})$  is bounded by  $\tilde{R}_\ell(S_\ell)$  with high probability. The proof for these properties involves careful use of the structure of MNL model to show that the value of  $\tilde{R}_\ell(S_\ell)$  is equal to the highest expected revenue achievable by any feasible assortment, among *all instances of the problem with parameters in the range*  $[0, v_i^{UCB}]$ ,  $i = 1, \dots, n$ . Since the actual parameters lie in this range with high probability, we have  $\tilde{R}_\ell(S_\ell)$  is at least  $R(S^*, \mathbf{v})$  with high probability. Lemma 4.2 provides the precise statement.

**Bounding the regret.** The final part of our analysis is to bound the regret in each epoch. First, we use the above property to bound the loss due to offering the optimistic assortment  $S_\ell$ . In particular,

we show that the loss is bounded by the difference between optimistic estimated revenue,  $\tilde{R}_\ell(S_\ell)$ , and actual expected revenue,  $R(S_\ell)$ . We then prove a Lipschitz property of the expected revenue function to bound the difference between optimistic estimate and expected revenues in terms of errors in individual product estimates  $|v_{i,\ell}^{\text{UCB}} - v_i|$ . Finally, we leverage the structure of MNL model and the properties of  $v_{i,\ell}^{\text{UCB}}$  to bound the regret in each epoch. Lemma 4.3 provide the precise statements of above properties.

Given the above properties, the rest of the proof is relatively straightforward. In rest of this section, we make the above notions precise. Finally, in Appendix A.3, we utilize these properties to complete the proof of Theorem 1.

## 4.2. Upper confidence bounds

In this section, we will show that the upper confidence bounds  $v_{i,\ell}^{\text{UCB}}$  converge to the true parameters  $v_i$  from above. Specifically, we have the following result.

**Lemma 4.1** *For every  $\ell$ , we have:*

1.  $v_{i,\ell}^{\text{UCB}} \geq v_i$  with probability at least  $1 - \frac{5}{\ell}$  for all  $i = 1, \dots, N$ .
2. There exists constants  $C_1$  and  $C_2$  such that

$$v_{i,\ell}^{\text{UCB}} - v_i \leq C_1 \sqrt{\frac{v_i \log(\ell + 1)}{T_i(\ell)}} + C_2 \frac{\log(\ell + 1)}{T_i(\ell)}$$

with probability at least  $1 - \frac{5}{\ell}$ .

We first establish that the estimates  $\hat{v}_{i,\ell}$ ,  $\ell \leq L$  are unbiased i.i.d estimates of the true parameter  $v_i$  for all products. It is not immediately clear a priori if the estimates  $\hat{v}_{i,\ell}$ ,  $\ell \leq L$  are independent. In our setting, it is possible that the distribution of estimate  $\hat{v}_{i,\ell}$  depends on the offered assortment  $S_\ell$ , which in turn depends on the history and therefore, previous estimates,  $\{\hat{v}_{i,\tau}, \tau = 1, \dots, \ell - 1\}$ . In Lemma A.1, we show that the moment generating of  $\hat{v}_{i,\ell}$  conditioned on  $S_\ell$  only depends on the parameter  $v_i$  and not on the offered assortment  $S_\ell$ , there by establishing that estimates are independent and identically distributed. Using the moment generating function, we show that  $\hat{v}_{i,\ell}$  is a geometric random variable with mean  $v_i$ , i.e.,  $E(\hat{v}_{i,\ell}) = v_i$ . We will use this observation and extend the multiplicative Chernoff-Hoeffding bounds discussed in Mitzenmacher and Upfal (2005) and Babaioff et al. (2015) to geometric random variables and derive the result. The proof is provided in Appendix A.2

### 4.3. Optimistic estimate and convergence rates

In this section, we show that the estimated revenue converges to the optimal expected revenue from above. First, we show that the estimated revenue is an upper confidence bound on the optimal revenue. In particular, we have the following result.

**Lemma 4.2** *Suppose  $S^* \in \mathcal{S}$  is the assortment with highest expected revenue, and Algorithm 1 offers  $S_\ell \in \mathcal{S}$  in each epoch  $\ell$ . Then, for every epoch  $\ell$ , we have*

$$\tilde{R}_\ell(S_\ell) \geq \tilde{R}_\ell(S^*) \geq R(S^*, \mathbf{v}) \text{ with probability at least } 1 - \frac{5}{\ell}.$$

In Lemma A.3, we show that the optimal expected revenue is monotone in the MNL parameters. It is important to note that we do not claim that the expected revenue is in general a monotone function, but only value of the expected revenue corresponding to optimal assortment increases with increase in MNL parameters. The result follows from Lemma 4.1, where we established that  $v_{i,\ell}^{\text{UCB}} > v_i$  with high probability. We provide the detailed proof in Appendix A.2.

The following result provides the convergence rates of the estimate  $\tilde{R}_\ell(S_\ell)$  to the optimal expected revenue.

**Lemma 4.3** *There exists constants  $C_1$  and  $C_2$  such that for every  $\ell$ , we have*

$$(1 + \sum_{j \in S_\ell} v_j)(\tilde{R}_\ell(S_\ell) - R(S_\ell, \mathbf{v})) \leq C_1 \sqrt{\frac{v_i \log(\ell+1)}{|\mathcal{T}_i(\ell)|}} + C_2 \frac{\log(\ell+1)}{|\mathcal{T}_i(\ell)|},$$

*with probability at least  $1 - \frac{5}{\ell}$*

In Lemma A.4, we show that the expected revenue function satisfies a certain kind of Lipschitz condition. Specifically, the difference between estimated,  $\tilde{R}_\ell(S_\ell)$ , and expected revenues,  $R_\ell(S_\ell)$ , is bounded by the sum of errors in parameter estimates for the products,  $|v_{i,\ell}^{\text{UCB}} - v_i|$ . This observation in conjunction with the “optimistic estimates” property will let us bound the regret as an aggregated difference between estimated revenues and expected revenues of the offered assortments. Noting that we have already computed convergence rates between the parameter estimates earlier, we can extend them to show that the estimated revenues converge to the optimal revenue from above. We complete the proof in Appendix A.2.

## 5. Extensions

In this section, we extend our approach (Algorithm 1) to the setting where the assumption that  $v_i \leq v_0$  for all  $i$  is relaxed. The essential ideas in the extension remain the same as our earlier approach,

specifically optimism under uncertainty, and our policy is structurally similar to Algorithm 1. The modified policy requires a small but mandatory initial exploration period. However, unlike the works of Rusmevichientong et al. (2010) and Sauré and Zeevi (2013), the exploratory period does not depend on the specific instance parameters and is constant for all problem instances. Therefore, our algorithm is parameter independent and remains relevant for practical applications. Moreover, our approach continues to simultaneously explore and exploit after the initial exploratory phase. In particular, the initial exploratory phase is to ensure that the estimates converge to the true parameters from above particularly in cases when the attraction parameter,  $v_i$  (frequency of purchase), is large for certain products. We describe our approach in Algorithm 2.

---

**Algorithm 2** Exploration-Exploitation algorithm for MNL-Bandit general parameters

---

```

1: Initialization:  $v_{i,0}^{\text{UCB}} = 1$  for all  $i = 1, \dots, N$ .
2:  $t = 1$  ;  $\ell = 1$  keeps track of the time steps and total number of epochs respectively
3:  $T_i(1) = 0$  for all  $i = 1, \dots, N$ .
4: while  $t < T$  do
5:   Compute  $S_\ell = \operatorname{argmax}_{S \in \mathcal{S}} \tilde{R}_\ell(S) = \frac{\sum_{i \in S} r_i v_{i,\ell-1}^{\text{UCB}}}{1 + \sum_{j \in S} v_{j,\ell-1}^{\text{UCB}}}$ 
6:   if  $T_i(\ell) < 96 \log(\ell + 1)$  for some  $i \in S_\ell$  then
7:     Consider  $\hat{S} = \{i | T_i(\ell) < 48 \log(\ell + 1)\}$ .
8:     Chose  $S_\ell \in \mathcal{S}$  such that  $S_\ell \subset \hat{S}$ .
9:   end if
10:  Offer assortment  $S_\ell$ , observe the purchasing decision,  $\mathbf{c}_t$  of the consumer
11:  if  $c_t = 0$  then
12:    compute  $\hat{v}_{i,\ell} = \sum_{t \in \mathcal{E}_\ell} \mathbb{1}(c_t = i)$ , no. of consumers who preferred  $i$  in epoch  $\ell$ , for all  $i \in S_\ell$ .
13:    update  $\mathcal{T}_i(\ell) = \{\tau \leq \ell | i \in S_\ell\}$ ,  $T_i(\ell) = |\mathcal{T}_i(\ell)|$ , no. of epochs until  $\ell$  that offered product  $i$ .
14:    update  $\bar{v}_{i,\ell} = \frac{1}{T_i(\ell)} \sum_{\tau \in \mathcal{T}_i(\ell)} \hat{v}_{i,\tau}$ , sample mean of the estimates
15:    update  $v_{i,\ell}^{\text{UCB2}} = \bar{v}_{i,\ell} + \max \left\{ \sqrt{\bar{v}_{i,\ell}}, \bar{v}_{i,\ell} \right\} \sqrt{\frac{48 \log(\ell+1)}{T_i(\ell)}} + \frac{48 \log(\ell+1)}{T_i(\ell)}$ 
16:     $\ell = \ell + 1$ 
17:  else
18:     $\mathcal{E}_\ell = \mathcal{E}_\ell \cup t$ , time indices corresponding to epoch  $\ell$ .
19:     $t = t + 1$ 
20:  end if
21: end while

```

---

We can extend the analysis in Section 4 to bound the regret of Algorithm 2 as follows.

**Theorem 2** *For any instance  $\mathbf{v} = (v_0, \dots, v_N)$ , of the MNL-Bandit problem with  $N$  products,  $r_i \in [0, 1]$  for all  $i = 1, \dots, N$ , the regret of the policy corresponding to Algorithm 2 at time  $T$  is bounded as,*

$$\text{Reg}_\pi(T, \mathbf{v}) \leq C_1 \sqrt{BNT \log T} + C_2 N \log^2 T + C_3 NB \log T,$$

where  $C_1$ ,  $C_2$  and  $C_3$  are absolute constants and  $B = \max\{\max_i \frac{v_i}{v_0}, 1\}$ .

**Proof outline.** Note that the Algorithm 2 is very similar to Algorithm 1 except for the initial exploratory phase. Hence, to bound the regret, we first prove that the initial exploratory phase is indeed bounded and then follow the approach discussed in Section 4 to establish the correctness of confidence intervals, the optimistic assortment and finally deriving the convergence rates and regret bounds. Given the above properties, the rest of the proof is relatively straightforward. We make the above notions precise and provide the complete proof in Appendix B.

## 6. Parameter dependent regret bounds

In this section, we derive an  $O(\log T)$  regret bound for Algorithm 1 that is parameter dependent. In Section 4 and 5, we established worst case regret bounds for Algorithm 1 that hold for all problem instances. Although our algorithm ensures that the exploration-exploitation tradeoff is balanced at all times, for problem instances that are “well separated”, our algorithm quickly converges to the optimal solution leading to better regret bounds. More specifically, we consider problem instances, where the optimal assortment and “second best” assortment are sufficiently “separated” and derive a  $O(\log T)$  regret bound that depends on the parameters of the instance. Note that, unlike the parameter-independent bound derived in Section 4, the bound we derive only holds for sufficiently large  $T$  and is dependent on the separation between the revenues corresponding optimal and second best assortments. In particular, let  $\Delta$  denote the difference between the expected revenues of the optimal and second-best assortment, i.e.,

$$\Delta(\mathbf{v}) = \min_{\{S \in \mathcal{S} \mid R(S) \neq R(S^*, \mathbf{v})\}} \{R(S^*, \mathbf{v}) - R(S)\},$$

We have the following result.

**Theorem 3** *For any instance,  $\mathbf{v} = (v_0, \dots, v_N)$  of the MNL-Bandit problem with  $N$  products,  $r_i \in [0, 1]$  and  $v_0 \geq v_i$  for  $i = 1, \dots, N$ , there exists finite constants  $B_1$  and  $B_2$  such that the regret of the policy defined in Algorithm 1 at time  $T$  is bounded as,*

$$\text{Reg}_\pi(T, \mathbf{v}) \leq B_1 \left( \frac{N^2 \log T}{\Delta(\mathbf{v})} \right) + B_2.$$

### 6.1. Proof outline.

In this setting, we analyze the regret by separately considering the epochs that satisfy certain desirable properties and the ones that do not. Specifically, we denote, epoch  $\ell$  as a “good” epoch, if the parameters,  $v_{i,\ell}^{\text{UCB}}$  satisfy the following property,

$$0 \leq v_{i,\ell}^{\text{UCB}} - v_i \leq C_1 \sqrt{\frac{v_i \log(\ell+1)}{T_i(\ell)}} + C_2 \frac{\log(\ell+1)}{T_i(\ell)},$$

and “bad” epoch, otherwise, where  $C_1$  and  $C_2$  are constants as defined in Lemma 4.1. Note that every epoch  $\ell$  is a good epoch with high probability  $(1 - \frac{5}{\ell})$  and we show that regret due to bad epochs is bounded by a constant (see Appendix C). Therefore, we focus on good epochs and show that there exists a constant  $\tau$ , such that after each product has been offered in at least  $\tau$  good epochs, Algorithm 1 finds the optimal assortment with high probability. Based on this result, we can then bound the total number of good epochs in which a sub-optimal assortment can be offered by our algorithm.

Specifically, let

$$\tau = \frac{4NC \log T}{\Delta^2(\mathbf{v})}, \quad (6.1)$$

where  $C = \max\{C_1, C_2\}$ . Then we have the following result.

**Lemma 6.1** *Let  $\ell$  be a good epoch and  $S_\ell$  be the assortment offered by Algorithm 1 in epoch  $\ell$  and if every product in assortment  $S_\ell$  is offered in at least  $\tau$  good epochs, i.e.  $T_i(\ell) \geq \tau$  for all  $i$ , then we have  $R(S_\ell, \mathbf{v}) = R(S^*, \mathbf{v})$ .*

*Proof.* From Lemma 4.3, we have,

$$\begin{aligned} R(S^*, \mathbf{v}) - R(S_\ell, \mathbf{v}) &\leq \frac{1}{V(S_\ell) + 1} \sum_{i \in S_\ell} \left( C_1 \sqrt{\frac{v_i \log(\ell+1)}{T_i(\ell)}} + C_2 \frac{\log(\ell+1)}{T_i(\ell)} \right) \\ &\leq \frac{\Delta(\mathbf{v}) \sum_{i \in S_\ell} \sqrt{v_i}}{2\sqrt{|S_\ell|}(V(S_\ell) + 1)} \leq \frac{\Delta(\mathbf{v})}{2}. \end{aligned}$$

The result follows from the definition of  $\Delta(\mathbf{v})$ .  $\square$

The next step in the analysis is to show that Algorithm 1 will offer a small number of sub-optimal assortments in good epochs. We make this precise in the following result, the proof is a simple counting exercise using Lemma 6.1 and is completed in Appendix C.

**Lemma 6.2** *Algorithm 1 cannot offer sub-optimal assortment in more than  $\frac{N(N-1)}{2} \tau$  good epochs.*

The proof for Theorem 3 follows from the above result. In particular, noting that the number of epochs in which sub-optimal assortment is offered is small, we re-use the regret analysis of Section

4 to bound the regret by  $O(N^2 \log T)$ . We provide the rigorous proof in Appendix C for the sake of completeness. Noting that for the special case of cardinality constraints, we have  $|S_\ell| \leq K$  for every epoch  $\ell$ , we can improve the regret bound to  $O(N^{3/2} K^{1/2} \log T)$ . Specifically,

**Corollary 6.1** *For any instance,  $\mathbf{v} = (v_0, \dots, v_N)$  of the MNL-Bandit problem with  $N$  products and cardinality constraint  $K$ ,  $r_i \in [0, 1]$  and  $v_0 \geq v_i$  for  $i = 1, \dots, N$ , there exists finite constants  $B_1$  and  $B_2$  such that the regret of the policy defined in Algorithm 1 at time  $T$  is bounded as,*

$$\text{Reg}_\pi(T, \mathbf{v}) \leq B_1 \frac{N^{\frac{3}{2}} K^{\frac{1}{2}} \log T}{\Delta(\mathbf{v})} + B_2.$$

It should be noted that the bound obtained in Corollary 6.1 is similar in magnitude to the regret bounds obtained by the approaches of Rusmevichientong et al. (2010) and Sauré and Zeevi (2013) for the cardinality constrained problem. (In fact our algorithm also has improved regret bounds compared to the  $O(N^2 \log^2 T)$  bound established by Rusmevichientong et al. (2010)). It is also important to note that our algorithm is independent of the parameter  $\Delta(\mathbf{v})$  unlike the existing work which requires a conservative estimate of  $\Delta(\mathbf{v})$  to implement a policy.

## 7. Lower bounds and near-optimality of the proposed policy

In this section, we consider the special case of TU constraints, namely, cardinality constrained assortment optimization problem and establish that any policy must incur a regret of  $\Omega(\sqrt{NT/K})$ . More precisely, we prove the following result.

**Theorem 4** *There exists a (randomized) instance of the MNL-Bandit problem with  $v_0 \geq v_i, i = 1, \dots, N$ , such that for any  $N, K < N, T \geq N$ , and any policy  $\pi$  that offers assortment  $S_t^\pi, |S_t^\pi| \leq K$  at time  $t$ , we have*

$$\text{Reg}_\pi(T, \mathbf{v}) := \mathbb{E}_\pi \left( \sum_{t=1}^T R(S^*, \mathbf{v}) - R(S_t^\pi, \mathbf{v}) \right) \geq C \sqrt{\frac{NT}{K}}$$

where  $S^*$  is (at-most)  $K$ -cardinality assortment with maximum expected revenue, and  $C$  is an absolute constant.

### 7.1. Proof overview

We prove Theorem 4 by a reduction to a parametric multi-armed bandit (MAB) problem, for which a lower bound is known.

**Definition 7.1 (MAB instance  $I_{MAB}$ )** *Define  $I_{MAB}$  as a (randomized) instance of MAB problem with  $N \geq 2$  Bernoulli arms, and following parameters (probability of reward 1)*

$$\mu_i = \begin{cases} \alpha, & \text{if } i \neq j, \\ \alpha + \epsilon, & \text{if } i = j, \end{cases} \quad \text{for all } i = 1, \dots, N,$$

where  $j$  is set uniformly at random from  $\{1, \dots, N\}$ ,  $\alpha < 1$  and  $\epsilon = \frac{1}{100} \sqrt{\frac{N\alpha}{T}}$ .



Throughout this section, we will use the terms algorithm and policy interchangeably. An algorithm  $\mathcal{A}$  is referred to as online if it adaptively selects a history dependent  $\mathcal{A}_t \in \{1, \dots, n\}$  at each time  $t$  on the lines of (2.4) for the MAB problem.

**Lemma 7.1** *For any  $N \geq 2$ ,  $\alpha < 1$ ,  $T$  and any online algorithm  $\mathcal{A}$  that plays arm  $\mathcal{A}_t$  at time  $t$ , the expected regret on instance  $I_{MAB}$  is at least  $\frac{\epsilon T}{6}$ . That is,*

$$\text{Reg}_{\mathcal{A}}(T, \boldsymbol{\mu}) := \mathbb{E} \left[ \sum_{t=1}^T (\mu_j - \mu_{\mathcal{A}_t}) \right] \geq \frac{\epsilon T}{6},$$

where, the expectation is both over the randomization in generating the instance (value of  $j$ ), as well as the random outcomes that result from pulled arms.

The proof of Lemma 7.1 is a simple extension of the proof of the  $\Omega(\sqrt{NT})$  lower bound for the Bernoulli instance with parameters  $\frac{1}{2}$  and  $\frac{1}{2} + \epsilon$  (for example, see Bubeck and Cesa-Bianchi (2012)), and has been provided in Appendix E for the sake of completeness. We use the above lower bound for the MAB problem to prove that any algorithm must incur at least  $\Omega(\sqrt{NT/K})$  regret on the following instance of the MNL-Bandit problem.

**Definition 7.2 (MNL-Bandit instance  $I_{MNL}$ )** *Define  $I_{MNL}$  as the following (randomized) instance of MNL-Bandit problem with  $K$ -cardinality constraint,  $\hat{N} = NK$  products, parameters  $v_0 = K$  and for  $i = 1, \dots, \hat{N}$ ,*

$$v_i = \begin{cases} \alpha, & \text{if } \lceil \frac{i}{K} \rceil \neq j, \\ \alpha + \epsilon, & \text{if } \lceil \frac{i}{K} \rceil = j, \end{cases}$$

where  $j$  is set uniformly at random from  $\{1, \dots, N\}$ ,  $\alpha < 1$ , and  $\epsilon = \frac{1}{100} \sqrt{\frac{N\alpha}{T}}$ .

We will show that any MNL-Bandit algorithm has to incur a regret of  $\Omega\left(\sqrt{\frac{NT}{K}}\right)$  on instance  $I_{MNL}$ . The main idea in our reduction is to show that if there exists an algorithm  $\mathcal{A}_{MNL}$  for MNL-Bandit that achieves  $o\left(\sqrt{\frac{NT}{K}}\right)$  regret on instance  $I_{MNL}$ , then we can use  $\mathcal{A}_{MNL}$  as a subroutine to construct an algorithm  $\mathcal{A}_{MAB}$  for the MAB problem that achieves strictly less than  $\frac{\epsilon T}{6}$  regret on instance  $I_{MAB}$  in time  $T$ , thus contradicting the lower bound of Lemma 7.1. This will prove Theorem 4 by contradiction.

## 7.2. Construction of the MAB algorithm using the MNL algorithm

Algorithm 3 provides the exact construction of  $\mathcal{A}_{MAB}$ , which simulates  $\mathcal{A}_{MNL}$  algorithm as a “black-box”. Note that  $\mathcal{A}_{MAB}$  pulls arms at time steps  $t = 1, \dots, T$ . These arm pulls are interleaved by simulations of  $\mathcal{A}_{MNL}$  steps (**Call  $\mathcal{A}_{MNL}$** , **Feedback to  $\mathcal{A}_{MNL}$** ). When step  $\ell$  of  $\mathcal{A}_{MNL}$  is simulated, it uses the feedback from  $1, \dots, \ell - 1$  to suggest an assortment  $S_\ell$ ; and recalls a feedback

**Algorithm 3** Algorithm  $\mathcal{A}_{MAB}$ 


---

```

1: Initialization:  $t = 0, \ell = 0$ .
2: while  $t \leq T$  do
3:   Update  $\ell = \ell + 1$ ;
4:   Call  $\mathcal{A}_{MNL}$ , receive assortment  $S_\ell \subset [\hat{N}]$ ;
5:   Repeat until ‘exit loop’
6:     With probability  $\frac{1}{2}$ , send Feedback to  $\mathcal{A}_{MNL}$  ‘no product was purchased’, exit loop;
7:     Update  $t = t + 1$ ;
8:     Pull arm  $\mathcal{A}_t = \lceil \frac{i}{K} \rceil$ , where  $i \in S_\ell$  is chosen uniformly at random.
9:     If reward is 1, send Feedback to  $\mathcal{A}_{MNL}$  ‘ $i$  was purchased’ and exit loop;
10:  end loop
11: end while

```

---

from  $\mathcal{A}_{MAB}$  on which product (or no product) was purchased out of those offered in  $S_\ell$ , where the probability of purchase of product  $i \in S_\ell$  is be  $v_i / (v_0 + \sum_{i \in S_\ell} v_i)$ . Before showing that the  $\mathcal{A}_{MAB}$  indeed provides the right feedback to  $\mathcal{A}_{MNL}$  in the  $\ell^{th}$  step for each  $\ell$ , we introduce some notation.

Let  $M_\ell$  denote the number of times no arm is pulled or arm  $\lceil \frac{i}{K} \rceil$  is pulled for some  $i \in S_\ell$  by  $\mathcal{A}_{MAB}$  before exiting the loop. For every  $i \in S_\ell \cup 0$ , let  $\zeta_\ell^i$  denote the event that the feedback to  $\mathcal{A}_{MNL}$  from  $\mathcal{A}_{MAB}$  after step  $\ell$  of  $\mathcal{A}_{MNL}$  is “product  $i$  is purchased”. We have,

$$\mathcal{P}(M_\ell = m \cap \zeta_\ell^i) = \frac{v_i}{2K} \left( \frac{1}{2K} \sum_{i \in S_\ell} (1 - v_i) \right)^{m-1} \quad \text{for each } i \in S_\ell \cup \{0\}.$$

Hence, the probability that  $\mathcal{A}_{MAB}$ ’s feedback to  $\mathcal{A}_{MNL}$  is “product  $i$  is purchased” is,

$$p_{S_\ell}(i) = \sum_{m=1}^{\infty} \mathcal{P}(M_\ell = m \cap \zeta_\ell^i) = \frac{v_i}{v_0 + \sum_{q \in S_\ell} v_q}.$$

This establish that  $\mathcal{A}_{MAB}$  provides the appropriate feedback to  $\mathcal{A}_{MNL}$ .

### 7.3. Proof of Theorem 4.

We prove the result by establishing three key results. First, we upper bound the regret for the MAB algorithm,  $\mathcal{A}_{MAB}$ . Then, we prove a lower bound on the regret for the MNL algorithm,  $\mathcal{A}_{MNL}$ . Finally, we relate the regret of  $\mathcal{A}_{MAB}$  and  $\mathcal{A}_{MNL}$  and use the established lower and upper bounds to show a contradiction.

For the rest of this proof, assume that  $L$  is the total number of calls to  $\mathcal{A}_{MNL}$  in  $\mathcal{A}_{MAB}$ . Let  $S^*$  be the optimal assortment for  $I_{MNL}$ . For any instantiation of  $I_{MNL}$ , it is easy to see that the optimal

assortment contains  $K$  items, all with parameter  $\alpha + \epsilon$ , i.e., it contains all  $i$  such that  $\lceil \frac{i}{K} \rceil = j$ . Therefore,  $V(S^*) = K(\alpha + \epsilon) = K\mu_j$ .

**Upper bound for the regret of the MAB algorithm.** The first step in our analysis is to prove an upper bound on the regret of the MAB algorithm,  $\mathcal{A}_{MAB}$  on the instance  $I_{MAB}$ . Let us label the loop following the  $\ell$ th call to  $\mathcal{A}_{MNL}$  in Algorithm 3 as  $\ell$ th loop. Note that the probability of exiting the loop is  $p = E[\frac{1}{2} + \frac{1}{2}\mu_{\mathcal{A}_\ell}] = \frac{1}{2} + \frac{1}{2K}V(S_\ell)$ . In every step of the loop until exited, an arm is pulled with probability  $1/2$ . The optimal strategy would pull the best arm so that the total expected optimal reward in the loop is  $\sum_{r=1}^{\infty} (1-p)^{r-1} \frac{1}{2} \mu_j = \frac{\mu_j}{2p} = \frac{1}{2Kp}V(S^*)$ . Algorithm  $\mathcal{A}_{MAB}$  pulls a random arm from  $S_\ell$ , so total expected algorithm's reward in the loop is  $\sum_{r=1}^{\infty} (1-p)^{r-1} \frac{1}{2K}V(S_\ell) = \frac{1}{2Kp}V(S_\ell)$ . Subtracting the algorithm's reward from the optimal reward, and substituting  $p$ , we obtain that the total expected regret of  $\mathcal{A}_{MAB}$  over the arm pulls in loop  $\ell$  is

$$\frac{V(S^*) - V(S_\ell)}{(K + V(S_\ell))}.$$

Noting that  $V(S_\ell) \geq K\alpha$ , we have the following upper bound on the regret for the MAB algorithm.

$$Reg_{\mathcal{A}_{MAB}}(T, \boldsymbol{\mu}) \leq \frac{1}{(1 + \alpha)} \mathbb{E} \left( \sum_{\ell=1}^L \frac{1}{K} (V(S^*) - V(S_\ell)) \right), \quad (7.1)$$

where the expectation in equation (7.1) is over the random variables  $L$  and  $S_\ell$ .

**Lower bound for the regret of the MNL algorithm.** Here, using simple algebraic manipulations we derive a lower bound on the regret of the MNL algorithm,  $\mathcal{A}_{MNL}$  on the instance  $I_{MNL}$ . Specifically,

$$\begin{aligned} Reg_{\mathcal{A}_{MNL}}(L, \mathbf{v}) &= \mathbb{E} \left[ \sum_{\ell=1}^L \frac{V(S^*)}{v_0 + V(S^*)} - \frac{V(S_\ell)}{v_0 + V(S_\ell)} \right] \\ &\geq \frac{1}{K(1 + \alpha)} \mathbb{E} \left[ \sum_{\ell=1}^L \left( \frac{V(S^*)}{1 + \frac{\epsilon}{1 + \alpha}} - V(S_\ell) \right) \right] \end{aligned}$$

Therefore, it follows that,

$$Reg_{\mathcal{A}_{MNL}}(L, \mathbf{v}) \geq \frac{1}{(1 + \alpha)} \mathbb{E} \left[ \sum_{\ell=1}^L \frac{1}{K} (V(S^*) - V(S_\ell)) - \frac{\epsilon v^* L}{(1 + \alpha)^2} \right], \quad (7.2)$$

where the expectation in equation (7.2) is over the random variables  $L$  and  $S_\ell$ .

**Relating the regret of MNL algorithm and MAB algorithm.** Finally, we relate the regret of MNL algorithm  $\mathcal{A}_{MNL}$  and MAB algorithm  $\mathcal{A}_{MAB}$  to derive a contradiction. The first step in relating the regret involves relating the length of the horizons of  $\mathcal{A}_{MNL}$  and  $\mathcal{A}_{MAB}$ ,  $L$  and  $T$  respectively. Intuitively, after any call to  $\mathcal{A}_{MNL}$  (“Call  $\mathcal{A}_{MNL}$ ” in Algorithm 3), many iterations

of the following loop may be executed; in roughly  $1/2$  of those iterations, an arm is pulled and  $t$  is advanced (with probability  $1/2$ , the loop is exited without advancing  $t$ ). Therefore,  $T$  should be at least a constant fraction of  $L$ . Lemma E.3 in Appendix makes this intuition precise by showing that  $\mathbb{E}(L) \leq 3T$ .

Now we are ready to prove Theorem 4. From (7.1) and (7.2), we have

$$\text{Reg}_{\mathcal{A}_{MAB}}(T, \boldsymbol{\mu}) \leq \mathbb{E} \left( \text{Reg}_{\mathcal{A}_{MNL}}(L, \mathbf{v}) + \frac{\epsilon v^* L}{(1 + \alpha)^2} \right).$$

For the sake of contradiction, suppose that the regret of the  $\mathcal{A}_{MNL}$ ,  $\text{Reg}_{\mathcal{A}_{MNL}}(L, \mathbf{v}) \leq c\sqrt{\frac{\hat{N}L}{K}}$  for a constant  $c$  to be prescribed below. Then, from Jensen's inequality, it follows that,

$$\text{Reg}_{\mathcal{A}_{MAB}}(T, \boldsymbol{\mu}) \leq c\sqrt{\frac{\hat{N}\mathbb{E}(L)}{K}} + \frac{\epsilon v^* \mathbb{E}(L)}{(1 + \alpha)^2}$$

From lemma E.3, we have that  $\mathbb{E}(L) \leq 3T$ . Therefore, we have,  $c\sqrt{\frac{\hat{N}\mathbb{E}(L)}{K}} = c\sqrt{N\mathbb{E}(L)} \leq c\sqrt{3NT} = c\epsilon T\sqrt{\frac{3}{\alpha}} < \frac{\epsilon T}{12}$  on setting  $c < \frac{1}{12}\sqrt{\frac{\alpha}{3}}$ . Also, using  $v^* = \alpha + \epsilon \leq 2\alpha$ , and setting  $\alpha$  to be a small enough constant, we can get that the second term above is also strictly less than  $\frac{\epsilon T}{12}$ . Combining these observations, we have

$$\text{Reg}_{\mathcal{A}_{MAB}}(T, \boldsymbol{\mu}) < \frac{\epsilon T}{12} + \frac{\epsilon T}{12} = \frac{\epsilon T}{6},$$

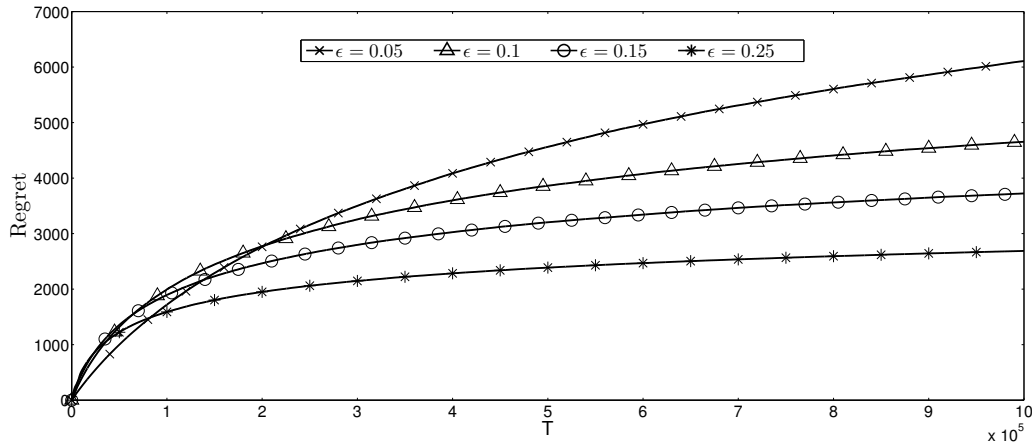
thus arriving at a contradiction.  $\square$

## 8. Computational study

In this section, we present insights from numerical experiments that test the empirical performance of our policy and highlight some of its salient features. We study the performance of Algorithm 1 from the perspective of robustness with respect to the “separability parameter” of the underlying instance. In particular, we consider varying levels of separation between the optimal revenue and the second best revenue and perform a regret analysis numerically. We contrast the performance of Algorithm 1 with the approach in Sauré and Zeevi (2013) for different levels of separation between the optimal and sub-optimal revenues. We observe that when the separation between the optimal assortment and second best assortment is sufficiently small, the approach in Sauré and Zeevi (2013) breaks down, i.e., incurs linear regret, while the regret of Algorithm 1 only grows sub-linearly with respect to the selling horizon. We also present results from a simulated study on a real world data set.

### 8.1. Robustness of Algorithm 1.

Here, we present a study that examines the robustness of Algorithm 1 with respect to the instance separability. We consider a parametric instance (see Example 8.1), where the separation between



**Figure 1** Performance of Algorithm 1 measured as the regret on the parametric instance (8.1). The graphs illustrate the dependence of the regret on  $T$  for  $\epsilon = 0.05, 0.1, 0.15$  and  $0.25$  respectively.

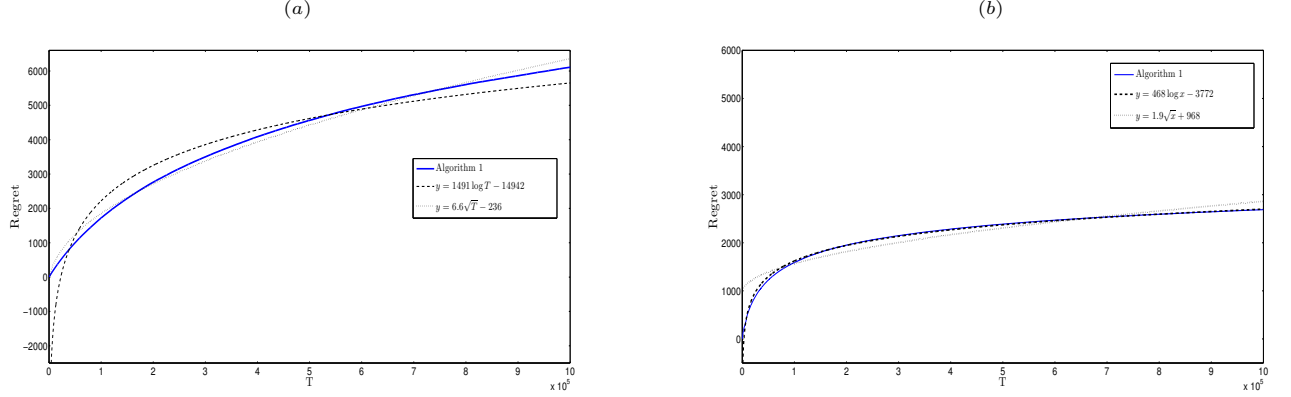
the revenues of the optimal assortment and next best assortment is specified by the parameter  $\epsilon$  and compare the performance of Algorithm 1 for different values of  $\epsilon$ .

**Experimental setup.** We consider the parametric MNL setting with  $N = 10$ ,  $K = 4$ ,  $r_i = 1$  for all  $i$  and utility parameters  $v_0 = 1$  and for  $i = 1, \dots, N$ ,

$$v_i = \begin{cases} 0.25 + \epsilon, & \text{if } i \in \{1, 2, 9, 10\} \\ 0.25, & \text{else,} \end{cases} \quad (8.1)$$

where  $0 < \epsilon < 0.25$ ; recall  $\epsilon$  specifies the difference between revenues corresponding to the optimal assortment and the next best assortment. Note that this problem has a unique optimal assortment,  $\{1, 2, 9, 10\}$  with an expected revenue of  $1 + 4\epsilon/2 + 4\epsilon$  and next best revenue of  $1 + 3\epsilon/2 + 3\epsilon$ . We consider four different values for  $\epsilon$ ,  $\epsilon = \{0.05, 0.1, 0.15, 0.25\}$ , where higher value of  $\epsilon$  corresponding to larger separation, and hence an “easier” problem instance.

**Results.** Figure 1 summarizes the performance of Algorithm 1 for different values of  $\epsilon$ . The results are based on running 100 independent simulations, the standard errors are within 2%. Note that the performance of Algorithm 1 is consistent across different values of  $\epsilon$ ; with a regret that exhibits sub linear growth. Observe that as the value of  $\epsilon$  increases the regret of Algorithm 1 decreases. While not immediately obvious from Figure 1, the regret behavior is fundamentally different in the case of “small”  $\epsilon$  and “large”  $\epsilon$ . To see this, in Figure 2, we focus on the regret for  $\epsilon = 0.05$  and  $\epsilon = 0.25$  and fit to  $\log T$  and  $\sqrt{T}$  respectively. (The parameters of these functions are obtained via respectively simple linear regression of the regret vs  $\log T$  and regret vs  $\sqrt{T}$ ). It can be observed that regret is roughly logarithmic when  $\epsilon = 0.25$ , and in contrast roughly behaves like  $\sqrt{T}$  when  $\epsilon = 0.05$ . This illustrates the theory developed in Section 6, where we showed that the regret grows



**Figure 2** Best fit for the regret of Algorithm 1 on the parametric instance (8.1). The graphs (a), (b) illustrate the dependence of the regret on  $T$  for  $\epsilon = 0.05$ , and  $0.25$  respectively. The best  $y = \beta_1 \log T + \beta_0$  fit and best  $y = \beta_1 \sqrt{T} + \beta_0$  fit are superimposed on the regret curve.

logarithmically with time, if the optimal assortment and next best assortment are “well separated”, while the worst-case regret scales as  $\sqrt{T}$ .

## 8.2. Comparison with existing approaches.

In this section, we present a computational study comparing the performance of our algorithm to that of Sauré and Zeevi (2013). (To the best of our knowledge, Sauré and Zeevi (2013) is currently the best existing approach for our problem setting.) To be implemented, their approach requires certain a priori information of a “separability parameter”; roughly speaking, measuring the degree to which the optimal and next-best assortments are distinct from a revenue standpoint. More specifically, their algorithm follows an *explore-then-exploit* approach, where every product is first required to be offered for a minimum duration of time that is determined by an estimate of said “separability parameter.” After this mandatory exploration phase, the parameters of the choice model are estimated based on the past observations and the optimal assortment corresponding to the estimated parameters is offered for the subsequent consumers. If the optimal assortment and the next best assortment are “well separated,” then the offered assortment is optimal with high probability, otherwise, the algorithm could potentially incur linear regret. Therefore, the knowledge of this “separability parameter” is crucial. For our comparison, we consider the exploration period suggested by Sauré and Zeevi (2013) and compare it with the performance of Algorithm 1 for different values of separation ( $\epsilon$ .) We will show that for any given exploration period, there is an instance where the approach in Sauré and Zeevi (2013) “breaks down” or in other words incurs linear regret, while the regret of Algorithm 1 grows sub-linearly (as  $O(\sqrt{T})$ , more precisely) for all values of  $\epsilon$  as asserted in Theorem 1.

Attribute	Attribute Values
price	Very-high, high, medium, low
maintenance costs	Very-high, high, medium, low
# doors	2, 3, 4, 5 or more
passenger capacity	2, 4, more than 4
luggage capacity	small, medium and big
safety perception	low, medium, high

**Table 1**      Attributional information of cars in the database

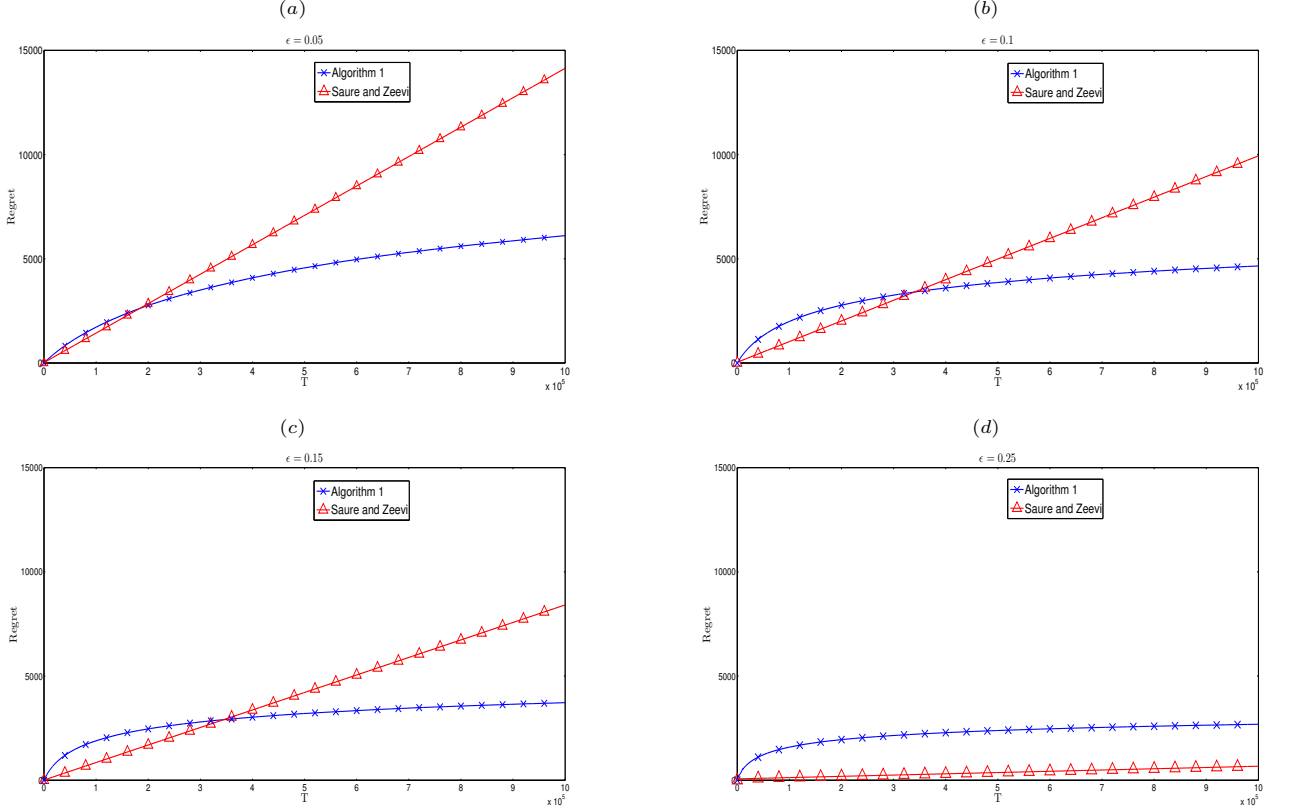
**Experimental setup and results.** We consider the parametric MNL setting as described in (8.1) and for each value of  $\epsilon \in \{0.05, 0.1, 0.15, 0.25\}$ . Since the implementation of the policy in Sauré and Zeevi (2013) requires knowledge of the selling horizon and minimum exploration period a priori, we consider the exploration period to be  $20 \log T$  as suggested in Sauré and Zeevi (2013) and the selling horizon as  $T = 10^6$ . Figure 3 compares the regret of Algorithm 1 with that of Sauré and Zeevi (2013). The results are based on running 100 independent simulations with standard error of 2%. We observe that the regret of the Sauré and Zeevi (2013) is better than the regret of Algorithm 1 when  $\epsilon = 0.25$  but is worse for other values of  $\epsilon$ . This can be attributed to the fact that for the assumed exploration period, Their algorithm fails to identify the optimal assortment within the exploration phase with sufficient probability and hence incurs a linear regret for  $\epsilon = 0.05, 0.1$  and  $0.15$ . Specifically, among the 100 simulations we tested, the algorithm of Sauré and Zeevi (2013) identified the optimal assortment for only 7%, 40%, 61% and 97% cases, when  $\epsilon = 0.05, 0.1, 0.15$ , and  $0.25$  respectively. This highlights the sensitivity to the “separability parameter” and the importance of having a reasonable estimate on the exploration period. Needless to say, such information is typically not available in practice. In contrast, the performance of Algorithm 1 is consistent across different values of  $\epsilon$ , insofar as the regret grows in a sub-linear fashion in all cases.

### 8.3. Performance of Algorithm 1 on real data.

Here, we present the results of a simulated study of Algorithm 1 on a real data set.

**Data description.** We consider the “UCI Car Evaluation Database” (see Lichman (2013)) which contains attributes based information of  $N = 1728$  cars and consumer ratings for each car. The exact details of the attributes are provided in Table 1. Rating for each car is also available. In particular, every car is associated with one of the following four ratings, unacceptable, acceptable, good and very good.

**Assortment optimization framework.** We assume that the consumer choice is modeled by the MNL model, where the mean utility of a product is linear in the values of attributes. More



**Figure 3** Comparison with the algorithm of Sauré and Zeevi (2013). The graphs (a), (b), (c) and (d) compares the performance of Algorithm 1 to that of Sauré and Zeevi (2013) on problem instance (8.1), for  $\epsilon = 0.05, 0.1, 0.15$  and  $0.25$  respectively.

specifically, we convert the categorical attributes described in Table 1 to attributes with binary values by adding dummy attributes (for e.g. “price very high”, “price low” are considered as two different attributes that can take values 1 or 0). Now every car is associated with an attribute vector  $m_i \in \{0, 1\}^{22}$ , which is known a priori and the mean utility of product  $i$  is given by the inner product

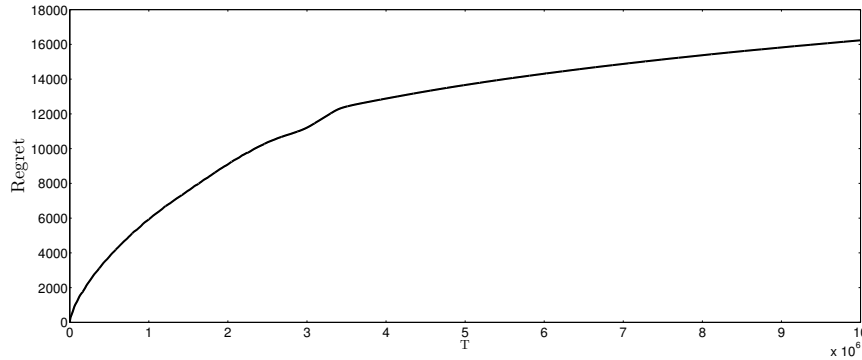
$$\mu_i = \theta \cdot m_i \quad i = 1, \dots, N,$$

where  $\theta \in \mathbb{R}^{22}$  is some fixed, but initially unknown attribute weight vector. Under this model, the probability that a consumer purchases product  $i$  when offered an assortment  $S \subset \{1, \dots, N\}$  is given by,

$$p_i(S) = \begin{cases} \frac{e^{\theta \cdot m_i}}{1 + \sum_{j \in S} e^{\theta \cdot m_j}}, & \text{if } i \in S \cup \{0\} \\ 0, & \text{otherwise,} \end{cases} \quad (8.2)$$

Let  $\mathbf{m} = (m_1, \dots, m_N)$ . Our goal is to offer assortments  $S_1, \dots, S_T$  at times  $1, \dots, T$  respectively such that the cumulative sales are maximized or alternatively, minimize the regret defined as





**Figure 4** Performance of Algorithm 1 on “UCI Car Evaluation Database”. The graph illustrate the dependence of regret on  $T$ .

$$Reg_{\pi}(T, \mathbf{m}) = \sum_{t=1}^T \left( \sum_{i \in S^*} p_i(S) - \sum_{i \in S_t} p_i(S_t) \right), \quad (8.3)$$

where

$$S^* = \arg \max_S \sum_{i \in S} \frac{e^{\theta \cdot m_i}}{1 + \sum_{j \in S} e^{\theta \cdot m_j}}.$$

Note that regret defined in (8.3) is a special case formulation of the regret defined in (2.6) with  $r_i = 1$  and  $v_i = e^{\theta \cdot m_i}$  for all  $i = 1, \dots, N$ .

**Experimental setup and results.** We first estimate a ground truth MNL model as follows. Using the available attribute level data and consumer rating for each car, we estimate a logistic model assuming every car’s rating is independent of the ratings of other cars to estimate the attribute weight vector  $\theta$ . Specifically, under the logistic model, the probability that a consumer will purchase a car whose attributes are defined by the vector  $m \in \{0, 1\}^{22}$  and the attribute weight vector  $\theta$  is given by

$$p_{\text{buy}}(\theta, m) \triangleq \mathbb{P}(\text{buy}|\theta) = \frac{e^{\theta \cdot m}}{1 + e^{\theta \cdot m}}.$$

For the purpose of training the logistic model on the available data, we consider the consumer ratings of “acceptable,” “good,” and “very good” as success or intention to buy and the consumer rating of “unacceptable” as a failure or no intention to buy. We then use the maximum likelihood estimate  $\theta_{\text{MLE}}$  for  $\theta$  to run simulations and study the performance of Algorithm 1 for the realized  $\theta_{\text{MLE}}$ . In particular, we compute  $\theta_{\text{MLE}}$  that maximizes the following regularized log-likelihood

$$\theta_{\text{MLE}} = \arg \max_{\theta} \sum_{i=1}^N \log p_{\text{buy}}(\theta, m_i) - \|\theta\|_2.$$

The objective function in the preceding optimization problem is convex and therefore we can use any of the standard convex optimization techniques to obtain the estimate,  $\theta_{\text{MLE}}$ . It is important to

note that the logistic model is only employed to obtain an estimate for  $\theta$ ,  $\theta_{\text{MLE}}$ . The estimate  $\theta_{\text{MLE}}$  is assumed to be the ground truth MNL model and is used to simulate the feedback of consumer choices for our learning Algorithm 1. We measure the performance of Algorithm 1 in terms of regret as defined in (8.3) with  $\theta = \theta_{\text{MLE}}$ .

Figure 4 plots the regret of Algorithm 1 on the “UCI Car Evaluation Database”, when the selling horizon is  $T = 10^7$  and at each time index, retailer can show at most  $k = 100$  cars. The results are based on running 100 independent simulations and have a standard error of 2%. The regret of Algorithm 1, is seen to grow in a sublinear fashion with respect to the selling horizon. It should be noted that Algorithm 1 did not require any a priori knowledge on the parameters unlike the other existing approaches such as Sauré and Zeevi (2013) and therefore can be applied to a wide range of other settings.

## Acknowledgments

V. Goyal is supported in part by NSF Grants CMMI-1351838 (CAREER) and CMMI-1636046. A. Zeevi is supported in part by NSF Grants NetSE-0964170 and BSF-2010466.

## References

- Angluin, D., L. G. Valiant. 1977. Fast probabilistic algorithms for hamiltonian circuits and matchings. *Proceedings of the Ninth Annual ACM Symposium on Theory of Computing*. STOC '77, 30–41.
- Auer, P. 2003. Using confidence bounds for exploitation-exploration trade-offs. *J. Mach. Learn. Res.* .
- Auer, Peter, Nicolo Cesa-Bianchi, Paul Fischer. 2002. Finite-time analysis of the multiarmed bandit problem. *Machine learning* **47**(2-3) 235–256.
- Babaioff, M., S. Dughmi, R. Kleinberg, A. Slivkins. 2015. Dynamic pricing with limited supply. *ACM Transactions on Economics and Computation* **3**(1) 4.
- Ben-Akiva, M., S. Lerman. 1985. *Discrete choice analysis: theory and application to travel demand*, vol. 9. MIT press.
- Blanchet, Jose, Guillermo Gallego, Vineet Goyal. 2016. A markov chain approximation to choice modeling. *Operations Research* **64**(4) 886–905.
- Bubeck, S., N. Cesa-Bianchi. 2012. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning* .
- Caro, F., J. Gallien. 2007. Dynamic assortment with demand learning for seasonal consumer goods. *Management Science* **53**(2) 276–292.
- Chen, W., Y. Wang, Y. Yuan. 2013. Combinatorial multi-armed bandit: General framework, results and applications. *Proceedings of the 30th international conference on machine learning*. 151–159.

- Davis, J., G. Gallego, H. Topaloglu. 2013. Assortment planning under the multinomial logit model with totally unimodular constraint structures. *Technical Report, Cornell University*. .
- Davis, J.M., G. Gallego, H. Topaloglu. 2011. Assortment optimization under variants of the nested logit model. *Technical report, Cornell University*. .
- Désir, A., V. Goyal. 2014. Near-optimal algorithms for capacity constrained assortment optimization. *Available at SSRN* .
- Désir, A., V. Goyal, D. Segev, C. Ye. 2015. Capacity constrained assortment optimization under the markov chain based choice model. *Under Review* .
- Farias, V., S. Jagabathula, D. Shah. 2013. A nonparametric approach to modeling choice with limited data. *Management Science* **59**(2) 305–322.
- Filippi, S., O. Cappe, A. Garivier, C. Szepesvári. 2010. Parametric bandits: The generalized linear case. *Advances in Neural Information Processing Systems*. 586–594.
- Gallego, G., R. Ratliff, S. Shebalov. 2014. A general attraction model and sales-based linear program for network revenue management under customer choice. *Operations Research* **63**(1) 212–232.
- Gallego, G., H. Topaloglu. 2014. Constrained assortment optimization for the nested logit model. *Management Science* **60**(10) 2583–2601.
- Kleinberg, R., A. Slivkins, E. Upfal. 2008. Multi-armed bandits in metric spaces. *Proceedings of the Fortieth Annual ACM Symposium on Theory of Computing*. STOC '08, 681–690.
- Kök, A. G., M. L. Fisher. 2007. Demand estimation and assortment optimization under substitution: Methodology and application. *Operations Research* **55**(6) 1001–1021.
- Lai, T.L., H. Robbins. 1985. Asymptotically efficient adaptive allocation rules. *Adv. Appl. Math.* **6**(1) 4–22.
- Li, G., P. Rusmevichientong, H. Topaloglu. 2015. The d-level nested logit model: Assortment and price optimization problems. *Operations Research* **63**(2) 325–342.
- Lichman, M. 2013. UCI machine learning repository. URL <http://archive.ics.uci.edu/ml>.
- Luce, R.D. 1959. *Individual choice behavior: A theoretical analysis*. Wiley.
- McFadden, D. 1978. Modeling the choice of residential location. *Transportation Research Record* (673).
- Mitzenmacher, Michael, Eli Upfal. 2005. *Probability and computing: Randomized algorithms and probabilistic analysis*. Cambridge university press.
- Plackett, R. L. 1975. The analysis of permutations. *Applied Statistics* 193–202.
- Rusmevichientong, P., Z. M. Shen, D. B. Shmoys. 2010. Dynamic assortment optimization with a multinomial logit choice model and capacity constraint. *Operations research* **58**(6) 1666–1680.
- Rusmevichientong, P., J. N. Tsitsiklis. 2010. Linearly parameterized bandits. *Math. Oper. Res.* **35**(2) 395–411.

- Sauré, D., A. Zeevi. 2013. Optimal dynamic assortment planning with demand learning. *Manufacturing & Service Operations Management* **15**(3) 387–404.
- Talluri, K., G. van Ryzin. 2004. Revenue management under a general discrete choice model of consumer behavior. *Management Science* **50**(1) 15–33.
- Train, K. E. 2009. *Discrete choice methods with simulation*. Cambridge university press.
- Williams, H.C.W.L. 1977. On the formation of travel demand models and economic evaluation measures of user benefit. *Environment and Planning A* **9**(3) 285–344.

## A. Proof of Theorem 1

In rest of this section, we provide a detailed proof of Theorem 1 following the outline discussed in Section 4.1. The proof is organized as follows. In Section A.1, we complete the proof of Lemma 4.1 and in Section A.2, we prove similar properties for estimates  $\tilde{R}_\ell(S_\ell)$ . Finally, in Section A.3, we utilize these properties to complete the proof of Theorem 1.

### A.1. Properties of estimates $v_{i,\ell}^{\text{UCB}}$

First, we focus on properties of the estimates  $\hat{v}_{i,\ell}$  and  $\bar{v}_{i,\ell}$ , and then extend these properties to establish the necessary properties of  $v_{i,\ell}^{\text{UCB}}$ .

#### Unbiased Estimates.

**Lemma A.1** *The moment generating function of estimate conditioned on  $S_\ell$ ,  $\hat{v}_i$ , is given by,*

$$\mathbb{E}_\pi \left( e^{\theta \hat{v}_{i,\ell}} \middle| S_\ell \right) = \frac{1}{1 - v_i(e^\theta - 1)}, \text{ for all } \theta \leq \log \frac{1 + v_i}{v_i}, \text{ for all } i = 1, \dots, N.$$

*Proof.* From (2.1), we have that probability of no purchase event when assortment  $S_\ell$  is offered is given by

$$p_0(S_\ell) = \frac{1}{1 + \sum_{j \in S_\ell} v_j}.$$

Let  $n_\ell$  be the total number of offerings in epoch  $\ell$  before a no purchased occurred, i.e.,  $n_\ell = |\mathcal{E}_\ell| - 1$ . Therefore,  $n_\ell$  is a geometric random variable with probability of success  $p_0(S_\ell)$ . And, given any fixed value of  $n_\ell$ ,  $\hat{v}_{i,\ell}$  is a binomial random variable with  $n_\ell$  trials and probability of success given by

$$q_i(S_\ell) = \frac{v_i}{\sum_{j \in S_\ell} v_j}.$$

In the calculations below, for brevity we use  $p_0$  and  $q_i$  respectively to denote  $p_0(S_\ell)$  and  $q_i(S_\ell)$ . Hence, we have

$$\mathbb{E}_\pi \left( e^{\theta \hat{v}_{i,\ell}} \right) = E_{n_\ell} \left\{ \mathbb{E}_\pi \left( e^{\theta \hat{v}_{i,\ell}} \middle| n_\ell \right) \right\}.$$

Since the moment generating function for a binomial random variable with parameters  $n, p$  is  $(pe^\theta + 1 - p)^n$ , we have

$$\mathbb{E}_\pi(e^{\theta \hat{v}_{i,\ell}} | n_\ell) = E_{n_\ell} \left\{ (q_i e^\theta + 1 - q_i)^{n_\ell} \right\}.$$

For any  $\alpha$ , such that,  $\alpha(1 - p) < 1$   $n$  is a geometric random variable with parameter  $p$ , we have

$$\mathbb{E}_\pi(\alpha^n) = \frac{p}{1 - \alpha(1 - p)}.$$

Note that for all  $\theta < \log \frac{1+v_i}{v_i}$ , we have  $(q_i e^\theta + (1 - q_i))(1 - p_0) = (1 - p_0) + p_0 v_i (e^\theta - 1) < 1$ . Therefore, we have  $\mathbb{E}_\pi(e^{\theta \hat{v}_{i,\ell}}) = \frac{1}{1 - v_i(e^\theta - 1)}$  for all  $\theta < \log \frac{1+v_i}{v_i}$ .  $\square$

From the moment generating function, we can establish that  $\hat{v}_{i,\ell}$  is a geometric random variable with parameter  $\frac{1}{1+v_i}$ . Thereby also establishing that  $\hat{v}_{i,\ell}$  and  $\bar{v}_{i,\ell}$  are unbiased estimators of  $v_i$ . More specifically, from lemma A.1, we have the following result.

**Corollary A.1** *We have the following results.*

1.  $\hat{v}_{i,\ell}$ ,  $\ell \leq L$  are i.i.d geometrical random variables with parameter  $\frac{1}{1+v_i}$ , i .e. for any  $\ell, i$

$$\mathbb{P}_\pi(\hat{v}_{i,\ell} = m) = \left( \frac{v_i}{1+v_i} \right)^m \left( \frac{1}{1+v_i} \right) \quad \forall m = \{0, 1, 2, \dots\}.$$

2.  $\hat{v}_{i,\ell}$ ,  $\ell \leq L$  are unbiased i.i.d estimates of  $v_i$ , i .e.  $\mathbb{E}_\pi(\hat{v}_{i,\ell}) = v_i \forall \ell, i$ .

**Concentration Bounds.** From Corollary A.1, it follows that  $\hat{v}_{i,\tau}, \tau \in \mathcal{T}_i(\ell)$  are i.i.d geometric random variables with mean  $v_i$ . We will use this observation and extend the multiplicative Chernoff-Hoeffding bounds discussed in Mitzenmacher and Upfal (2005) and Babaioff et al. (2015) to geometric random variables and derive the result.

**Lemma A.2** *If  $v_i \leq v_0$  for all  $i$ , for every epoch  $\ell$ , in Algorithm 1, we have the following concentration bounds.*

1.  $\mathbb{P}_\pi \left( |\bar{v}_{i,\ell} - v_i| > \sqrt{48 \bar{v}_{i,\ell} \frac{\log(\ell+1)}{T_i(\ell)}} + \frac{48 \log(\ell+1)}{T_i(\ell)} \right) \leq \frac{5}{\ell}.$
2.  $\mathbb{P}_\pi \left( |\bar{v}_{i,\ell} - v_i| > \sqrt{24 v_i \frac{\log(\ell+1)}{T_i(\ell)}} + \frac{48 \log(\ell+1)}{T_i(\ell)} \right) \leq \frac{5}{\ell}.$
3.  $\mathbb{P}_\pi \left( |\bar{v}_{i,\ell} - v_i| > \sqrt{\frac{24 \log(\ell+1)}{n}} + \frac{48 \log(\ell+1)}{n} \right) \leq \frac{5}{\ell}$

Note that to apply standard Chernoff-Hoeffding inequality (see p.66 in Mitzenmacher and Upfal (2005)), we must have the individual sample values bounded by some constant, which is not the case with our estimates  $\hat{v}_{i,\tau}$ . However, these estimates are geometric random variables and therefore

have extremely small tails. Hence, we can extend the Chernoff-Hoeffding bounds discussed in Mitzenmacher and Upfal (2005) and Babaioff et al. (2015) to geometric variables and prove the above result. Lemma 4.1 follows directly from Lemma A.2. The proof of Lemma A.2 is long and tedious and in the interest of continuity, we complete the proof in Appendix D. Following the proof of Lemma A.2, we obtain a very similar result that is useful to establish concentration bounds for the general parameter setting.

Lemma 4.1 follows from Corollary A.1 and Lemma A.2 and establishes the necessary properties of  $v_{i,\ell}^{\text{UCB}}$  and  $v_{i,\ell}^{\text{UCB2}}$  as alluded in the proof outline.

## A.2. Properties of estimate $\tilde{R}(S)$

In this section, we show that the estimates  $\tilde{R}_\ell(S_\ell)$  are upper confidence bounds converging “quickly” to the expected revenue corresponding to the optimal assortment,  $R(S^*, \mathbf{v})$ .

**Optimistic Estimates.** First, we establish that  $\tilde{R}_\ell(S)$  is an upper bound estimate of the optimal revenue,  $R(S^*, \mathbf{v})$ . Let  $R(S, \mathbf{w})$  be as defined in (2.2), namely, the expected revenue when assortment  $S$  is offered and if the parameters of the MNL were given by the vector  $\mathbf{w}$ . Then, from definition of  $\tilde{R}_\ell(S)$  (refer to (3.7)), it follows that  $\tilde{R}_\ell(S) = R(S, \mathbf{v}_\ell^{\text{UCB}})$ .

**Lemma A.3** *Assume  $0 \leq w_i \leq v_i^{\text{UCB}}$  for all  $i = 1, \dots, n$ . Suppose  $S$  is an optimal assortment when the MNL parameters are given by  $\mathbf{w}$ . Then,  $R(S, \mathbf{v}^{\text{UCB}}) \geq R(S, \mathbf{w})$ .*

*Proof.* We prove the result by first showing that for any  $j \in S$ , we have  $R(S, \mathbf{w}^j) \geq R(S, \mathbf{w})$ , where  $\mathbf{w}^j$  is vector  $\mathbf{w}$  with the  $j^{\text{th}}$  component increased to  $v_j^{\text{UCB}}$ , i.e.  $w_i^j = w_i$  for all  $i \neq j$  and  $w_j^j = v_j^{\text{UCB}}$ . We can use this result iteratively to argue that increasing each parameter of MNL to the highest possible value increases the value of  $R(S, \mathbf{w})$  to complete the proof.

It is easy to see that if  $r_j < R(S)$  removing the product  $j$  from assortment  $S$  yields higher expected revenue contradicting the optimality of  $S$ . Therefore, we have

$$r_j \geq R(S).$$

Multiplying by  $(v_j^{\text{UCB}} - w_j)(\sum_{i \in S/j} w_i + 1)$  on both sides of the above inequality and re-arranging terms, we can show that  $R(S, \mathbf{w}^j) \geq R(S, \mathbf{w})$ .  $\square$

Let  $\hat{S}, \mathbf{w}^*$  be maximizers of the optimization problem,  $\max_{S \in \mathcal{S}} \max_{0 \leq w_i \leq v_{i,\ell}^{\text{UCB}}} R(S, \mathbf{w})$ . Then applying lemma A.3 on assortment  $\hat{S}$  and parameters  $\mathbf{v}^*$  and noting that  $v_{i,\ell}^{\text{UCB}} > v_i$  with high probability, we have that

$$\tilde{R}_\ell(S_\ell) = \max_{S \in \mathcal{S}} R(S, \mathbf{v}_\ell^{\text{UCB}}) \geq \max_{S \in \mathcal{S}} \max_{0 \leq w_i \leq v_{i,\ell}^{\text{UCB}}} R(S, \mathbf{w}) \geq R(S^*, \mathbf{v}).$$

**Bounding Regret.** Now we will establish the connection between the error on the expected revenues and the error on the estimates of MNL parameters. In particular, we have the following result.

**Lemma A.4** *If  $0 \leq v_i \leq v_{i,\ell}^{\text{UCB}}$  for all  $i \in S_\ell$ , then*

$$\tilde{R}_\ell(S_\ell) - R(S_\ell, \mathbf{v}) \leq \frac{\sum_{j \in S_\ell} (v_{j,\ell}^{\text{UCB}} - v_j)}{1 + \sum_{j \in S_\ell} v_j}.$$

*Proof.* Since  $1 + \sum_{i \in S_\ell} v_{i,\ell}^{\text{UCB}} \geq 1 + \sum_{i \in S_\ell} v_i$ , we have

$$\begin{aligned} \tilde{R}_\ell(S_\ell) - R(S_\ell, \mathbf{v}) &\leq \frac{\sum_{i \in S_\ell} r_i v_{i,\ell}^{\text{UCB}}}{1 + \sum_{j \in S_\ell} v_{j,\ell}^{\text{UCB}}} - \frac{\sum_{i \in S_\ell} r_i v_i}{1 + \sum_{j \in S_\ell} v_{j,\ell}^{\text{UCB}}}, \\ &\leq \frac{\sum_{i \in S_\ell} (v_{i,\ell}^{\text{UCB}} - v_i)}{1 + \sum_{j \in S_\ell} v_{j,\ell}^{\text{UCB}}} \leq \frac{\sum_{i \in S_\ell} (v_{i,\ell}^{\text{UCB}} - v_i)}{1 + \sum_{j \in S_\ell} v_j} \end{aligned}$$

□

Lemma 4.3 follows directly from the above result and Lemma 4.1, while Lemma 4.3 follows directly from the above result and Lemma 4.1.

### A.3. Putting it all together: Proof of Theorem 1

In this section, we formalize the intuition developed in the previous sections and complete the proof of Theorem 1.

Let  $S^*$  denote the optimal assortment, our objective is to minimize the *regret* defined in (2.6), which is same as

$$\text{Reg}_\pi(T, \mathbf{v}) = \mathbb{E}_\pi \left\{ \sum_{\ell=1}^L |\mathcal{E}_\ell| (R(S^*, \mathbf{v}) - R(S_\ell, \mathbf{v})) \right\}, \quad (\text{A.1})$$

Note that  $L$ ,  $\mathcal{E}_\ell$  and  $S_\ell$  are all random variables and the expectation in equation (A.1) is over these random variables. Let  $\mathcal{H}_\ell$  be the filtration (history) associated with the policy upto epoch  $\ell$ . In particular,

$$\mathcal{H}_\ell = \sigma(U, C_1, \dots, C_{t(\ell)}, S_1, \dots, S_{t(\ell)}),$$

where  $t(\ell)$  is the time index corresponding to the end of epoch  $\ell$ . The length of the  $\ell^{\text{th}}$  epoch,  $|\mathcal{E}_\ell|$  conditioned on  $S_\ell$  is a geometric random variable with success probability defined as the probability of no-purchase in  $S_\ell$ , i.e.

$$p_0(S_\ell) = \frac{1}{1 + \sum_{j \in S_\ell} v_j}.$$

Let  $V(S_\ell) = \sum_{j \in S_\ell} v_j$ , then we have  $\mathbb{E}_\pi (|\mathcal{E}_\ell| \mid S_\ell) = 1 + V(S_\ell)$ . Noting that  $S_\ell$  in our policy is determined by  $\mathcal{H}_{\ell-1}$ , we have  $\mathbb{E}_\pi (|\mathcal{E}_\ell| \mid \mathcal{H}_{\ell-1}) = 1 + V(S_\ell)$ . Therefore, by law of conditional expectations, we have

$$\text{Reg}_\pi(T, \mathbf{v}) = \mathbb{E}_\pi \left\{ \sum_{\ell=1}^L \mathbb{E}_\pi [|\mathcal{E}_\ell| (R(S^*, \mathbf{v}) - R(S_\ell, \mathbf{v})) \mid \mathcal{H}_{\ell-1}] \right\},$$

and hence the regret can be reformulated as

$$Reg_\pi(T, \mathbf{v}) = \mathbb{E}_\pi \left\{ \sum_{\ell=1}^L (1 + V(S_\ell)) (R(S^*, \mathbf{v}) - R(S_\ell, \mathbf{v})) \right\} \quad (\text{A.2})$$

the expectation in equation (A.2) is over the random variables  $L$  and  $S_\ell$ . We now provide the proof sketch for Theorem 1 and the complete proof is provided in the full version of the paper.

**Proof of Theorem 1.** For sake of brevity, let  $\Delta R_\ell = (1 + V(S_\ell)) (R(S^*, \mathbf{v}) - R(S_\ell, \mathbf{v}))$ , for all  $\ell = 1, \dots, L$ . Now the regret can be reformulated as

$$Reg_\pi(T, \mathbf{v}) = \mathbb{E}_\pi \left\{ \sum_{\ell=1}^L \Delta R_\ell \right\} \quad (\text{A.3})$$

Let  $T_i$  denote the total number of epochs that offered an assortment containing product  $i$ . Let  $\mathcal{A}_0$  denote the complete set  $\Omega$  and for all  $\ell = 1, \dots, L$ , event  $\mathcal{A}_\ell$  is given by

$$\mathcal{A}_\ell = \left\{ v_{i,\ell}^{\text{UCB}} < v_i \text{ or } v_{i,\ell}^{\text{UCB}} > v_i + C_1 \sqrt{\frac{v_i}{T_i(\ell)}} \log(\ell+1) + C_2 \frac{\log^2(\ell+1)}{T_i(\ell)} \text{ for some } i = 1, \dots, N \right\}.$$

Noting that  $\mathcal{A}_\ell$  is a “low probability” event, we analyze the regret in two scenarios, one when  $\mathcal{A}_\ell$  is true and another when  $\mathcal{A}_\ell^c$  is true. Hence, we have

$$\mathbb{E}_\pi(\Delta R_\ell) = E[\Delta R_\ell \cdot \mathbb{I}(\mathcal{A}_{\ell-1}) + \Delta R_\ell \cdot \mathbb{I}(\mathcal{A}_{\ell-1}^c)]$$

Using the fact that  $R(S^*, \mathbf{v})$  and  $R(S_\ell, \mathbf{v})$  are both bounded by one and  $V(S_\ell) \leq N$ , we have

$$\mathbb{E}_\pi(\Delta R_\ell) \leq (N+1)\mathbb{P}_\pi(\mathcal{A}_{\ell-1}) + \mathbb{E}_\pi[\Delta R_\ell \cdot \mathbb{I}(\mathcal{A}_{\ell-1}^c)].$$

Whenever  $\mathbb{I}(\mathcal{A}_{\ell-1}^c) = 1$ , from Lemma A.3, we have  $\tilde{R}_\ell(S^*) \geq R(S^*, \mathbf{v})$  and by our algorithm design, we have  $\tilde{R}_\ell(S_\ell) \geq \tilde{R}_\ell(S^*)$  for all  $\ell \geq 2$ . Therefore, it follows that

$$\mathbb{E}_\pi\{\Delta R_\ell\} \leq (N+1)\mathbb{P}_\pi(\mathcal{A}_{\ell-1}) + \mathbb{E}_\pi\left\{[(1+V(S_\ell))(\tilde{R}_\ell(S_\ell) - R(S_\ell, \mathbf{v}))] \cdot \mathbb{I}(\mathcal{A}_{\ell-1}^c)\right\}$$

From Lemma 4.3, it follows that

$$[(1+V(S_\ell))(\tilde{R}_\ell(S_\ell) - R(S_\ell, \mathbf{v}))] \cdot \mathbb{I}(\mathcal{A}_{\ell-1}^c) \leq \sum_{i \in S_\ell} \left( C_1 \sqrt{\frac{v_i \log T}{T_i(\ell)}} + \frac{C_2 \log T}{T_i(\ell)} \right)$$

Therefore, we have

$$\mathbb{E}_\pi\{\Delta R_\ell\} \leq (N+1)\mathbb{P}_\pi(\mathcal{A}_{\ell-1}) + C \sum_{i \in S_\ell} \mathbb{E}_\pi \left( \sqrt{\frac{v_i \log T}{T_i(\ell)}} + \frac{\log T}{T_i(\ell)} \right) \quad (\text{A.4})$$

where  $C = \max\{C_1, C_2\}$ . Combining equations (A.2) and (A.4), we have

$$Reg_\pi(T, \mathbf{v}) \leq \mathbb{E}_\pi \left\{ \sum_{\ell=1}^L \left[ (N+1)\mathbb{P}_\pi(\mathcal{A}_{\ell-1}) + C \sum_{i \in S_\ell} \left( \sqrt{\frac{v_i \log T}{T_i(\ell)}} + \frac{\log T}{T_i(\ell)} \right) \right] \right\}.$$



Therefore, from Lemma 4.1, we have

$$\begin{aligned}
 \text{Reg}_\pi(T, \mathbf{v}) &\leq C\mathbb{E}_\pi \left\{ \sum_{\ell=1}^L \frac{N+1}{\ell} + \sum_{i \in S_\ell} \sqrt{\frac{v_i \log T}{T_i(\ell)}} + \sum_{i \in S_\ell} \frac{\log T}{T_i(\ell)} \right\}, \\
 &\stackrel{(a)}{\leq} CN \log T + CN \log^2 T + C\mathbb{E}_\pi \left( \sum_{i=1}^n \sqrt{v_i T_i \log T} \right) \\
 &\stackrel{(b)}{\leq} CN \log T + CN \log^2 T + C \sum_{i=1}^N \sqrt{v_i \log T \mathbb{E}_\pi(T_i)}
 \end{aligned} \tag{A.5}$$

Inequality (a) follows from the observation that  $L \leq T$ ,  $T_i \leq T$ ,  $\sum_{T_i(\ell)=1}^{T_i} \frac{1}{\sqrt{T_i(\ell)}} \leq \sqrt{T_i}$  and

$\sum_{T_i(\ell)=1}^{T_i} \frac{1}{T_i(\ell)} \leq \log T_i$ , while Inequality (b) follows from Jensen's inequality.

For any realization of  $L$ ,  $\mathcal{E}_\ell$ ,  $T_i$ , and  $S_\ell$  in Algorithm 1, we have the following relation  $\sum_{\ell=1}^L n_\ell \leq T$ . Hence, we have  $\mathbb{E}_\pi \left( \sum_{\ell=1}^L n_\ell \right) \leq T$ . Let  $\mathcal{S}$  denote the filtration corresponding to the offered assortments  $S_1, \dots, S_L$ , then by law of total expectation, we have,

$$\begin{aligned}
 \mathbb{E}_\pi \left( \sum_{\ell=1}^L n_\ell \right) &= \mathbb{E}_\pi \left\{ \sum_{\ell=1}^L E_{\mathcal{S}}(n_\ell) \right\} = \mathbb{E}_\pi \left\{ \sum_{\ell=1}^L 1 + \sum_{i \in S_\ell} v_i \right\} \\
 &= \mathbb{E}_\pi \left\{ L + \sum_{i=1}^n v_i T_i \right\} = \mathbb{E}_\pi \{L\} + \sum_{i=1}^n v_i \mathbb{E}_\pi(T_i).
 \end{aligned}$$

Therefore, it follows that

$$\sum v_i \mathbb{E}_\pi(T_i) \leq T. \tag{A.6}$$

To obtain the worst case upper bound, we maximize the bound in equation (A.5) subject to the condition (A.6) and hence, we have  $\text{Reg}_\pi(T, \mathbf{v}) = O(\sqrt{NT \log T} + N \log^2 T)$ .  $\square$

## B. Proof of Theorem 2

The proof for Theorem 2 is very similar to the proof of Theorem 1. Specifically, we first prove that the initial exploratory phase is indeed bounded and then follow the proof of Theorem 1 to establish the correctness of confidence intervals, optimistic assortment and finally deriving the convergence rates and regret bounds.

**Bounding Exploratory Epochs.** We would denote an epoch  $\ell$  as an “exploratory epoch” if the assortment offered in the epoch contains a product that has been offered in less than  $48 \log(\ell + 1)$  epochs. It is easy to see that the number of exploratory epochs is bounded by  $48N \log T$ , where  $T$  is the selling horizon under consideration. We then use the observation that the length of any epoch is a geometric random variable to bound the total expected duration of the exploration phase. Hence, we bound the expected regret due to explorations.

**Lemma B.1** *Let  $L$  be the total number of epochs in Algorithm 2 and let  $\mathcal{E}_L$  denote the set of “exploratory epochs”, i.e.*

$$E_L = \{\ell \mid \exists i \in S_\ell \text{ such that } T_i(\ell) < 48 \log(\ell + 1)\},$$

where  $T_i(\ell)$  is the number of epochs product  $i$  has been offered before epoch  $\ell$ . If  $\mathcal{E}_\ell$  denote the time indices corresponding to epoch  $\ell$  and  $v_i \leq Bv_0$  for all  $i = 1, \dots, N$ , for some  $B \geq 1$ , then we have that,

$$\mathbb{E}_\pi \left( \sum_{\ell \in E_L} |\mathcal{E}_\ell| \right) < 48NB \log T,$$

where the expectation is over all possible outcomes of Algorithm 2.

*Proof.* Consider an  $\ell \in E_L$ , note that  $|\mathcal{E}_\ell|$  is a geometric random variable with parameter  $\frac{1}{V(S_\ell)+1}$ . Since  $v_i \leq Bv_0$ , for all  $i$  and we can assume without loss of generality  $v_0 = 1$ , we have  $|\mathcal{E}_\ell|$  as a geometric random variable with parameter  $p$ , where  $p \geq \frac{1}{B|S_\ell|+1}$ . Therefore, we have the conditional expectation of  $|\mathcal{E}_\ell|$  given that assortment  $S_\ell$  is offered is bounded as,

$$\mathbb{E}_\pi (|\mathcal{E}_\ell| \mid S_\ell) \leq B|S_\ell| + 1.$$

Note that after every product has been offered in at least  $48 \log T$  epochs, then we do not have any exploratory epochs. Therefore, we have that

$$\sum_{\ell \in E_L} |S_\ell| \leq 48BN \log T + 48 \leq 96NB \log T.$$

The required result follows from the preceding two equations.  $\square$ .

**Confidence Intervals.** We will now show a result analogous to Lemma 4.1, that establish the updates in Algorithm 2,  $v_{i,\ell}^{\text{UCB2}}$ , as upper confidence bounds converging to actual parameters  $v_i$ . Specifically, we have the following result.

**Lemma B.2** *For every epoch  $\ell$ , if  $T_i(\ell) \geq 48 \log(\ell + 1)$  for all  $i \in S_\ell$ , then we have,*

1.  $v_{i,\ell}^{\text{UCB2}} \geq v_i$  with probability at least  $1 - \frac{5}{\ell}$  for all  $i = 1, \dots, N$ .
2. There exists constants  $C_1$  and  $C_2$  such that

$$v_{i,\ell}^{\text{UCB2}} - v_i \leq C_1 \max\{\sqrt{v_i}, v_i\} \sqrt{\frac{\log(\ell + 1)}{T_i(\ell)}} + C_2 \frac{\log(\ell + 1)}{T_i(\ell)},$$

with probability at least  $1 - \frac{5}{\ell}$ .

The proof is very similar to the proof of Lemma 4.1, where we first establish the following concentration inequality for the estimates  $\hat{v}_{i,\ell}$ , when  $T_i(\ell) \geq 48 \log(\ell + 1)$  from which the above result follows. The proof of Corollary B.1 is provided in Appendix D.

**Corollary B.1** *If in epoch  $\ell$ ,  $T_i(\ell) \geq 48 \log(\ell + 1)$  for all  $i \in S_\ell$ , then we have the following concentration bounds*

1.  $\mathbb{P}_\pi \left( |\bar{v}_{i,\ell} - v_i| \geq \max \{ \sqrt{\bar{v}_{i,\ell}}, \bar{v}_{i,\ell} \} \sqrt{\frac{48 \log(\ell + 1)}{n}} + \frac{48 \log(\ell + 1)}{n} \right) \leq \frac{5}{\ell}.$
2.  $\mathbb{P}_\pi \left( |\bar{v}_{i,\ell} - v_i| \geq \max \{ \sqrt{v_i}, v_i \} \sqrt{\frac{24 \log(\ell + 1)}{n}} + \frac{48 \log(\ell + 1)}{n} \right) \leq \frac{5}{\ell}.$
3.  $\mathbb{P}_\pi \left( |\bar{v}_{i,\ell} - v_i| > v_i \sqrt{\frac{12 \log(\ell + 1)}{n}} + \sqrt{\frac{6 \log(\ell + 1)}{n}} + \frac{48 \log(\ell + 1)}{n} \right) \leq \frac{5}{\ell}$

**Optimistic Estimate and Convergence Rates:** We will now establish two results analogous to Lemma 4.2 and 4.3, that show that the estimated revenue converges to the optimal expected revenue from above and also specify the convergence rate. In particular, we have the following two results.

**Lemma B.3** *Suppose  $S^* \in \mathcal{S}$  is the assortment with highest expected revenue, and Algorithm 2 offers  $S_\ell \in \mathcal{S}$  in each epoch  $\ell$ . Further, if  $T_i(\ell) \geq 48 \log(\ell + 1)$  for all  $i \in S_\ell$ , then we have,*

$$\tilde{R}_\ell(S_\ell) \geq \tilde{R}_\ell(S^*) \geq R(S^*, \mathbf{v}) \text{ with probability at least } 1 - \frac{5}{\ell}.$$

**Lemma B.4** *For every epoch  $\ell$ , if  $T_i(\ell) \geq 48 \log(\ell + 1)$  for all  $i \in S_\ell$ , then there exists constants  $C_1$  and  $C_2$  such that for every  $\ell$ , we have*

$$(1 + \sum_{j \in S_\ell} v_j)(\tilde{R}_\ell(S_\ell) - R(S_\ell, \mathbf{v})) \leq C_1 \max \{ \sqrt{v_i}, v_i \} \sqrt{\frac{\log(\ell + 1)}{|T_i(\ell)|}} + C_2 \frac{\log(\ell + 1)}{|T_i(\ell)|},$$

*with probability at least  $1 - \frac{5}{\ell}$ .*

### B.1. Putting it all together: Proof of Theorem 2

Proof of Theorem 2 is very similar to the proof of Theorem 1. We use the key results discussed above instead of similar results in Section 4 to complete the proof. Regret can be decomposed as

$$Reg_\pi(T, \mathbf{v}) = \mathbb{E}_\pi \left\{ \sum_{\ell \in E_L} |\mathcal{E}_\ell| (R(S^*, \mathbf{v}) - R(S_\ell, \mathbf{v})) + \sum_{\ell \notin E_L} |\mathcal{E}_\ell| (R(S^*, \mathbf{v}) - R(S_\ell, \mathbf{v})) \right\},$$

From Lemma B.1 and  $R(S, \mathbf{v}) \leq 1$  for any  $S \in \mathcal{S}$ , it follows that

$$Reg_\pi(T, \mathbf{v}) \leq 192NB \log T + \mathbb{E}_\pi \left\{ \sum_{\ell \notin E_L} |\mathcal{E}_\ell| (R(S^*, \mathbf{v}) - R(S_\ell, \mathbf{v})) \right\},$$

For sake of brevity, let  $\Delta R_\ell = (1 + V(S_\ell)) (R(S^*, \mathbf{v}) - R(S_\ell, \mathbf{v}))$ , for all  $\ell \notin E_L$ . Now the regret can be bounded as,

$$Reg_\pi(T, \mathbf{v}) \leq 192NB \log T + \mathbb{E}_\pi \left\{ \sum_{\ell \notin E_L} \Delta R_\ell \right\} \tag{B.1}$$

In the interest of avoiding redundant analysis, we claim that the following inequality can be derived from the proof of Theorem 1 adapted to the general setting.

$$\mathbb{E}_\pi \{\Delta R_\ell\} \leq C \sum_{i \in S_\ell} \mathbb{E}_\pi \left( \max \{\sqrt{v_i}, v_i\} \sqrt{\frac{\log T}{T_i(\ell)}} + \frac{\log T}{T_i(\ell)} \right) \quad (\text{B.2})$$

where  $C = \max\{C_1, C_2\}$ . Combining equations (B.1) and (B.2), we have

$$\text{Reg}_\pi(T, \mathbf{v}) \leq 192NB \log T + C \mathbb{E}_\pi \left\{ \sum_{\ell \notin E_L} \sum_{i \in S_\ell} \left( \max \{\sqrt{v_i}, v_i\} \sqrt{\frac{\log T}{T_i(\ell)}} + \frac{\log T}{T_i(\ell)} \right) \right\}.$$

Define sets  $\mathcal{I} = \{i | v_i \geq 1\}$  and  $\mathcal{D} = \{i | v_i < 1\}$ . Therefore, we have ,

$$\begin{aligned} \text{Reg}_\pi(T, \mathbf{v}) &\leq 192NB \log T + C \mathbb{E}_\pi \left\{ \sum_{\ell \notin E_L} \sum_{i \in S_\ell} \left( \max \{\sqrt{v_i}, v_i\} \sqrt{\frac{\log T}{T_i(\ell)}} + \frac{\log T}{T_i(\ell)} \right) \right\}, \\ &\stackrel{(a)}{\leq} 192NB \log T + CN \log^2 T + C \mathbb{E}_\pi \left( \sum_{i \in \mathcal{D}} \sqrt{v_i T_i \log T} + \sum_{i \in \mathcal{I}} v_i \sqrt{T_i \log T} \right) \quad (\text{B.3}) \\ &\stackrel{(b)}{\leq} 192NB \log T + CN \log^2 T + C \sum_{i \in \mathcal{D}} \sqrt{v_i \mathbb{E}_\pi(T_i) \log T} + \sum_{i \in \mathcal{I}} v_i \sqrt{\mathbb{E}_\pi(T_i) \log T} \end{aligned}$$

Inequality (a) follows from the observation that  $L \leq T$ ,  $T_i \leq T$ ,  $\sum_{T_i(\ell)=1}^{T_i} \frac{1}{\sqrt{T_i(\ell)}} \leq \sqrt{T_i}$  and

$\sum_{T_i(\ell)=1}^{T_i} \frac{1}{T_i(\ell)} \leq \log T_i$ , while Inequality (b) follows from Jensen's inequality.

From (A.6), we have that,

$$\sum v_i \mathbb{E}_\pi(T_i) \leq T.$$

To obtain the worst case upper bound, we maximize the bound in equation (B.3) subject to the above constraint. Noting that the objective in (B.3) is concave, we use the KKT conditions to derive the worst case bound as  $\text{Reg}_\pi(T, \mathbf{v}) = O(\sqrt{BNT \log T} + N \log^2 T + BN \log T)$ .  $\square$

## C. Parameter dependent bounds

**Proof of lemma 6.2** Assume for the sake of contradiction that Algorithm 1 has offered sub-optimal assortments in more than  $\frac{N(N-1)}{2} \tau$  epochs.

Let  $\ell_1$  be the epoch by which we have offered the first  $N\tau$  sub-optimal assortments and  $S_1$  be the subset of products that have been offered in at least  $\tau$  times by epoch  $\ell_1$ , i.e.

$$S_1 = \{i | T_i(\ell_1) \geq \tau\},$$

then by the pigeon hole principle, we have at least one product that has been offered at least  $\tau$  times (else the total number of offerings would not be more than  $N\tau$ ), i.e.

$$|S_1| \geq 1.$$

Similarly, let  $\ell_2$  be the epoch by which we have offered the next  $(N-1)\tau$  suboptimal assortments and  $S_2$  be the subset of products that have been offered in at least  $\tau$  epochs by epoch  $\ell_2$ , i.e.

$$S_2 = \{i | T_i(\ell_2) \geq \tau\}.$$

From lemma 6.1, we have that if any assortment is a subset of  $S_1$  and is offered after the  $\ell_1^{th}$  epoch, then it is an optimal assortment. Hence, among the next set of  $(N-1)\tau$  sub-optimal assortments, there must be at least one product that does not belong to  $S_1$  and therefore, by pigeon hole principle, we will have at least one more product that does not belong to  $S_1$  being offered in  $\tau$  epochs, i.e.

$$|S_2| - |S_1| \geq 1.$$

We can similarly define subsets  $S_3, \dots, S_N$  and prove

$$|S_{k+1} - S_k| \geq 1 \text{ for all } k \leq N-1,$$

establishing that after offering sub-optimal assortments in  $\frac{N(N-1)}{2}\tau$  good epochs, we have that, all the  $N$  products are offered in at least  $\tau$  epochs contradicting our hypothesis.  $\square$

**Proof of Theorem 3** Let  $V(S_\ell) = \sum_{j \in S_\ell} v_j$ , we have that

$$Reg_\pi(T, \mathbf{v}) = \mathbb{E}_\pi \left\{ \sum_{\ell=1}^L (1 + V(S_\ell)) (R(S^*, \mathbf{v}) - R(S_\ell, \mathbf{v})) \right\}$$

For sake of brevity, let  $\Delta R_\ell = (1 + V(S_\ell)) (R(S^*, \mathbf{v}) - R(S_\ell, \mathbf{v}))$ , for all  $\ell = 1, \dots, L$ . Now the regret can be reformulated as

$$Reg_\pi(T, \mathbf{v}) = \mathbb{E}_\pi \left\{ \sum_{\ell=1}^L \Delta R_\ell \right\} \quad (\text{C.1})$$

Let  $T_i$  denote the total number of epochs that offered a sub-optimal assortment containing product  $i$ . From Lemma 6.2, we have that

$$|T_i| \leq \frac{N(N-1)\tau}{2}. \quad (\text{C.2})$$

Let  $\mathcal{A}_0$  denote the complete set  $\Omega$  and for all  $\ell = 1, \dots, L$ , event  $\mathcal{A}_\ell$  is given by

$$\mathcal{A}_\ell = \left\{ v_{i,\ell}^{\text{UCB}} < v_i \text{ or } v_{i,\ell}^{\text{UCB}} > v_i + C_1 \sqrt{\frac{v_i \log(\ell+1)}{T_i(\ell)}} + C_2 \frac{\log(\ell+1)}{T_i(\ell)} \right\}.$$

Noting that  $\mathcal{A}_\ell$  is a “low probability” event, we analyze the regret in two scenarios, one when  $\mathcal{A}_\ell$  is true and another when  $\mathcal{A}_\ell^c$  is true. Hence, we have

$$\mathbb{E}_\pi(\Delta R_\ell) = \mathbb{E}_\pi [\Delta R_\ell \cdot \mathbb{I}(\mathcal{A}_{\ell-1}) + \Delta R_\ell \cdot \mathbb{I}(\mathcal{A}_{\ell-1}^c)]$$

Using the fact that  $R(S^*, \mathbf{v})$  and  $R(S_\ell, \mathbf{v})$  are both bounded by one and  $V(S_\ell) \leq N$ , we have

$$\mathbb{E}_\pi(\Delta R_\ell) \leq (N+1)\mathbb{P}_\pi(\mathcal{A}_{\ell-1}) + \mathbb{E}_\pi[\Delta R_\ell \cdot \mathbb{I}(\mathcal{A}_{\ell-1}^c)].$$

From Lemma 6.2, we have that a sub-optimal assortment cannot be offered for more than  $\frac{N(N-1)}{2}\tau$  good epochs, i.e epochs for which  $\mathbb{I}(\mathcal{A}_{\ell-1}^c) = 1$ . Let  $NO_L$  denote the set of good epochs in which a sub-optimal assortment is offered. From Lemma 6.2, we have that

$$|NO_L| \leq \frac{N(N-1)\tau}{2}.$$

For every  $\ell \in NO_L$ , whenever  $\mathbb{I}(\mathcal{A}_{\ell-1}^c) = 1$ , from Lemma A.3, we have  $\tilde{R}_\ell(S^*) \geq R(S^*, \mathbf{v})$  and by our algorithm design, we have  $\tilde{R}_\ell(S_\ell) \geq \tilde{R}_\ell(S^*)$  for all  $\ell \geq 2$ . Therefore, it follows that

$$\mathbb{E}_\pi\{\Delta R_\ell\} \leq (N+1)\mathbb{P}_\pi(\mathcal{A}_{\ell-1}) + \mathbb{E}_\pi\left\{\left[(1+V(S_\ell))(\tilde{R}_\ell(S_\ell) - R(S_\ell, \mathbf{v}))\right] \cdot \mathbb{I}(\mathcal{A}_{\ell-1}^c)\right\}$$

From Lemma 4.3, it follows that

$$\left[(1+V(S_\ell))(\tilde{R}_\ell(S_\ell) - R(S_\ell, \mathbf{v}))\right] \cdot \mathbb{I}(\mathcal{A}_{\ell-1}^c) \leq \sum_{i \in S_\ell} \left(C_1 \sqrt{\frac{v_i \log T}{T_i(\ell)}} + \frac{C_2 \log T}{T_i(\ell)}\right)$$

Therefore, we have

$$\mathbb{E}_\pi\{\Delta R_\ell\} \leq (N+1)\mathcal{P}(\mathcal{A}_{\ell-1}) + C\mathbb{E}_\pi\left(\sum_{i \in S_\ell} \left(\sqrt{\frac{v_i \log T}{T_i(\ell)}} + \frac{\log T}{T_i(\ell)}\right)\right) \quad (\text{C.3})$$

where  $C = \max\{C_1, C_2\}$ . Combining equations (C.1) and (C.3), we have

$$\text{Reg}_\pi(T, \mathbf{v}) \leq \mathbb{E}_\pi\left\{\sum_{\ell \in NO_L} \left[(N+1)\mathcal{P}(\mathcal{A}_{\ell-1}) + C\sum_{i \in S_\ell} \left(\sqrt{\frac{v_i \log T}{T_i(\ell)}} + \frac{\log T}{T_i(\ell)}\right)\right]\right\}.$$

For the regret to be maximum, we need the sub-optimal assortments as early as possible. Therefore, from Lemma 4.1, we have

$$\begin{aligned} \text{Reg}_\pi(T, \mathbf{v}) &\leq C\mathbb{E}_\pi\left\{\sum_{\ell \in NO_L} \frac{N+1}{\ell} + \sum_{i \in S_\ell} \sqrt{\frac{v_i \log T}{T_i(\ell)}} + \sum_{i \in S_\ell} \frac{\log T}{T_i(\ell)}\right\}, \\ &\stackrel{a}{\leq} CN + CN \log^2 T + C\mathbb{E}_\pi\left(\sum_{i=1}^n \sqrt{v_i T_i \log T}\right) \\ &\stackrel{b}{\leq} CN + CN \log^2 T + C\left(\sum_{i=1}^n \sqrt{v_i \mathbb{E}_\pi(T_i) \log T}\right) \\ &\stackrel{c}{\leq} CN + CN \log^2 T + \left(CN \sqrt{\frac{(N-1)}{2}\tau \log T}\right) \end{aligned} \quad (\text{C.4})$$

Inequality (a) follows from the observation that  $L \leq T$ ,  $T_i \leq T$ ,  $\sum_{T_i(\ell)=1}^{T_i} \frac{1}{\sqrt{T_i(\ell)}} \leq \sqrt{T_i}$  and  $\sum_{T_i(\ell)=1}^{T_i} \frac{1}{T_i(\ell)} \leq \log T_i$ , while Inequality (b) follows from Jensen's inequality. Inequality (c) follows from maximizing the bound in (C.4) subject to the constraint,  $\sum v_i \mathbb{E}_\pi(T_i) \leq \frac{N(N-1)}{2}$  as we have done in the proof of Theorem 1. Finally, the proof follows from the definition of  $\tau$  (See (6.1)).  $\square$

## D. Multiplicative Chernoff Bounds

We will extend the Chernoff bounds as discussed in Mitzenmacher and Upfal (2005)<sup>1</sup> to geometric random variables and establish the following concentration inequality.

**Theorem 5** Consider  $n$  i.i.d geometric random variables  $X_1, \dots, X_n$  with parameter  $p$ , i.e. for any  $i$

$$\Pr(X_i = m) = (1-p)^m p \quad \forall m = \{0, 1, 2, \dots\},$$

and let  $\mu = \mathbb{E}(X_i) = \frac{1-p}{p}$ . We have,

1.

$$\Pr\left(\frac{1}{n} \sum_{i=1}^n X_i > (1+\delta)\mu\right) \leq \begin{cases} \exp\left(-\frac{n\mu\delta^2}{2(1+\delta)(1+\mu)^2}\right) & \text{if } \mu \leq 1, \\ \exp\left(-\frac{n\delta^2\mu^2}{6(1+\mu)^2} \left(3 - \frac{2\delta\mu}{1+\mu}\right)\right) & \text{if } \mu \geq 1 \text{ and } \delta \in (0, 1). \end{cases}$$

and

2.

$$\Pr\left(\frac{1}{n} \sum_{i=1}^n X_i < (1-\delta)\mu\right) \leq \begin{cases} \exp\left(-\frac{n\delta^2\mu}{6(1+\mu)^2} \left(3 - \frac{2\delta\mu}{1+\mu}\right)\right) & \text{if } \mu \leq 1, \\ \exp\left(-\frac{n\delta^2\mu^2}{2(1+\mu)^2}\right) & \text{if } \mu \geq 1. \end{cases}$$

*Proof.* We will first bound  $\Pr\left(\frac{1}{n} \sum_{i=1}^n X_i > (1+\delta)\mu\right)$  and then follow a similar approach for bounding  $\Pr\left(\frac{1}{n} \sum_{i=1}^n X_i < (1-\delta)\mu\right)$  to complete the proof.

**Bounding  $\Pr\left(\frac{1}{n} \sum_{i=1}^n X_i > (1+\delta)\mu\right)$ :**

For all  $i$  and for any  $0 < t < \log \frac{1+\mu}{\mu}$ , we have,

$$\mathbb{E}(e^{tX_i}) = \frac{1}{1 - \mu(e^t - 1)}.$$

Therefore, from Markov Inequality, we have

$$\begin{aligned} \Pr\left(\frac{1}{n} \sum_{i=1}^n X_i > (1+\delta)\mu\right) &= \Pr\left(e^{t \sum_{i=1}^n X_i} > e^{(1+\delta)n\mu t}\right), \\ &\leq e^{-(1+\delta)n\mu t} \prod_{i=1}^n \mathbb{E}(e^{tX_i}), \\ &= e^{-(1+\delta)n\mu t} \left(\frac{1}{1 - \mu(e^t - 1)}\right)^n. \end{aligned}$$

<sup>1</sup> (originally discussed in Angluin and Valiant (1977))

Therefore, we have,

$$Pr \left( \frac{1}{n} \sum_{i=1}^n X_i > (1+\delta)\mu \right) \leq \min_{0 < t < \log \frac{1+\mu}{\mu}} e^{-(1+\delta)n\mu t} \left( \frac{1}{1-\mu(e^t-1)} \right)^n. \quad (\text{D.1})$$

We have,

$$\operatorname{argmin}_{0 < t < \log \frac{1+\mu}{\mu}} e^{-(1+\delta)n\mu t} \left( \frac{1}{1-\mu(e^t-1)} \right)^n = \operatorname{argmin}_{0 < t < \log \frac{1+\mu}{\mu}} -(1+\delta)n\mu t - n \log(1-\mu(e^t-1)), \quad (\text{D.2})$$

Noting that the right hand side in the above equation is a convex function in  $t$ , we obtain the optimal  $t$  by solving for the zero of the derivative. Specifically, at optimal  $t$ , we have

$$e^t = \frac{(1+\delta)(1+\mu)}{1+\mu(1+\delta)}.$$

Substituting the above expression in (D.1), we obtain the following bound.

$$Pr \left( \frac{1}{n} \sum_{i=1}^n X_i > (1+\delta)\mu \right) \leq \left( 1 - \frac{\delta}{(1+\delta)(1+\mu)} \right)^{n\mu(1+\delta)} \left( 1 + \frac{\delta\mu}{1+\mu} \right)^n. \quad (\text{D.3})$$

First consider the setting where  $\mu \in (0, 1)$ .

**Case 1a: If  $\mu \in (0, 1)$ :** From Taylor series of  $\log(1-x)$ , we have that

$$n\mu(1+\delta) \log \left( 1 - \frac{\delta}{(1+\delta)(1+\mu)} \right) \leq -\frac{n\delta\mu}{1+\mu} - \frac{n\delta^2\mu}{2(1+\delta)(1+\mu)^2},$$

From Taylor series for  $\log(1+x)$ , we have

$$n \log \left( 1 + \frac{\delta\mu}{1+\mu} \right) \leq \frac{n\delta\mu}{(1+\mu)},$$

Note that if  $\delta > 1$ , we can use the fact that  $\log(1+\delta x) \leq \delta \log(1+x)$  to arrive at the preceding result. Substituting the preceding two equations in (D.3), we have

$$Pr \left( \frac{1}{n} \sum_{i=1}^n X_i > (1+\delta)\mu \right) \leq \exp \left( -\frac{n\mu\delta^2}{2(1+\delta)(1+\mu)^2} \right), \quad (\text{D.4})$$

**Case 1b: If  $\mu \geq 1$ :** From Taylor series of  $\log(1-x)$ , we have that

$$n\mu(1+\delta) \log \left( 1 - \frac{\delta}{(1+\delta)(1+\mu)} \right) \leq -\frac{n\delta\mu}{1+\mu},$$

If  $\delta < 1$ , from Taylor series for  $\log(1+x)$ , we have

$$n \log \left( 1 + \frac{\delta\mu}{1+\mu} \right) \leq \frac{n\delta\mu}{(1+\mu)} - \frac{n\delta^2\mu^2}{6(1+\mu)^2} \left( 3 - \frac{2\delta\mu}{1+\mu} \right).$$

If  $\delta \geq 1$ , we have  $\log(1+\delta x) \leq \delta \log(1+x)$  and from Taylor series for  $\log(1+x)$ , it follows that,

$$n \log \left( 1 + \frac{\delta\mu}{1+\mu} \right) \leq \frac{n\delta\mu}{(1+\mu)} - \frac{n\delta\mu^2}{6(1+\mu)^2} \left( 3 - \frac{2\mu}{1+\mu} \right).$$

Therefore, substituting preceding results in (D.3), we have

$$Pr \left( \frac{1}{n} \sum_{i=1}^n X_i > (1+\delta)\mu \right) \leq \begin{cases} \exp \left( -\frac{n\delta^2\mu^2}{6(1+\mu)^2} \left( 3 - \frac{2\delta\mu}{1+\mu} \right) \right) & \text{if } \mu \geq 1 \text{ and } \delta \in (0, 1), \\ \exp \left( -\frac{n\delta\mu^2}{6(1+\mu)^2} \left( 3 - \frac{2\mu}{1+\mu} \right) \right) & \text{if } \mu \geq 1 \text{ and } \delta \geq 1 \end{cases} \quad (\text{D.5})$$



**Bounding**  $Pr\left(\frac{1}{n} \sum_{i=1}^n X_i < (1-\delta)\mu\right)$ :

Now to bound the other one sided inequality, we use the fact that

$$\mathbb{E}(e^{-tX_i}) = \frac{1}{1 - \mu(e^{-t} - 1)},$$

and follow a similar approach. More specifically, from Markov Inequality, for any  $t > 0$  and  $0 < \delta < 1$ , we have

$$\begin{aligned} Pr\left(\frac{1}{n} \sum_{i=1}^n X_i < (1-\delta)\mu\right) &= Pr\left(e^{-t \sum_{i=1}^n X_i} > e^{-(1-\delta)n\mu t}\right) \\ &\leq e^{(1-\delta)n\mu t} \prod_{i=1}^n \mathbb{E}(e^{-tX_i}) \\ &= e^{(1-\delta)n\mu t} \left(\frac{1}{1 - \mu(e^{-t} - 1)}\right)^n \end{aligned}$$

Therefore, we have

$$Pr\left(\frac{1}{n} \sum_{i=1}^n X_i < (1-\delta)\mu\right) \leq \min_{t>0} e^{-(1+\delta)n\mu t} \left(\frac{1}{1 - \mu(e^{-t} - 1)}\right)^n,$$

Following similar approach as in optimizing the previous bound (see (D.1)) to establish the following result.

$$Pr\left(\frac{1}{n} \sum_{i=1}^n X_i < (1-\delta)\mu\right) \leq \left(1 + \frac{\delta}{(1-\delta)(1+\mu)}\right)^{n\mu(1-\delta)} \left(1 - \frac{\delta\mu}{1+\mu}\right)^n.$$

Now we will use Taylor series for  $\log(1+x)$  and  $\log(1-x)$  in a similar manner as described for the other bound to obtain the required result. In particular, since  $1-\delta \leq 1$ , we have for any  $x > 0$  it follows that  $(1 + \frac{x}{1-\delta})^{(1-\delta)} \leq (1+x)$ . Therefore, we have

$$Pr\left(\frac{1}{n} \sum_{i=1}^n X_i < (1-\delta)\mu\right) \leq \left(1 + \frac{\delta}{1+\mu}\right)^{n\mu} \left(1 - \frac{\delta\mu}{1+\mu}\right)^n. \quad (\text{D.6})$$

**Case 2a. If  $\mu \in (0, 1)$ :** Note that since  $X_i \geq 0$  for all  $i$ , we have a zero probability event if  $\delta > 1$ . Therefore, we assume  $\delta < 1$  and from Taylor series for  $\log(1-x)$ , we have

$$n \log\left(1 - \frac{\delta\mu}{1+\mu}\right) \leq -\frac{n\delta\mu}{1+\mu},$$

and from Taylor series for  $\log(1+x)$ , we have

$$n\mu \log\left(1 + \frac{\delta}{1+\mu}\right) \leq \frac{n\delta\mu}{(1+\mu)} - \frac{n\delta^2\mu}{6(1+\mu)^2} \left(3 - \frac{2\delta\mu}{1+\mu}\right).$$

Therefore, substituting the preceding equations in (D.6), we have,

$$Pr\left(\frac{1}{n} \sum_{i=1}^n X_i < (1-\delta)\mu\right) \leq \exp\left(-\frac{n\delta^2\mu}{6(1+\mu)^2} \left(3 - \frac{2\delta\mu}{1+\mu}\right)\right). \quad (\text{D.7})$$

**Case 2b. If  $\mu \geq 1$ :** For similar reasons as discussed above, we assume  $\delta < 1$  and from Taylor series for  $\log(1-x)$ , we have

$$n \log \left( 1 - \frac{\delta\mu}{1+\mu} \right) \leq -\frac{n\delta\mu}{1+\mu} - \frac{n\delta^2\mu^2}{2(1+\mu)^2},$$

and from Taylor series for  $\log(1+x)$ , we have

$$n \log \left( 1 + \frac{\delta\mu}{1+\mu} \right) \leq \frac{n\delta}{(1+\mu)}.$$

Therefore, substituting the preceding equations in (D.6), we have,

$$Pr \left( \frac{1}{n} \sum_{i=1}^n X_i < (1-\delta)\mu \right) \leq \exp \left( -\frac{n\delta^2\mu^2}{2(1+\mu)^2} \right). \quad (\text{D.8})$$

The result follows from (D.4), (D.5), (D.7) and (D.8).  $\square$

Now, we will adapt a non-standard corollary from Babaioff et al. (2015) and Kleinberg et al. (2008) to our estimates to obtain sharper bounds.

**Lemma D.1** *Consider  $n$  i.i.d geometric random variables  $X_1, \dots, X_n$  with parameter  $p$ , i.e. for any  $i$ ,  $P(X_i = m) = (1-p)^m p \ \forall m = \{0, 1, 2, \dots\}$ . Let  $\mu = \mathbb{E}_\pi(X_i) = \frac{1-p}{p}$  and  $\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$ . If  $n > 48 \log(\ell+1)$ , then we have,*

1.  $\mathcal{P} \left( |\bar{X} - \mu| > \max \left\{ \sqrt{\bar{X}}, \bar{X} \right\} \sqrt{\frac{48 \log(\ell+1)}{n} + \frac{48 \log(\ell+1)}{n}} \right) \leq \frac{5}{\ell^2}$  for all  $n = 1, 2, \dots$ .
2.  $\mathcal{P} \left( |\bar{X} - \mu| \geq \max \left\{ \sqrt{\mu}, \mu \right\} \sqrt{\frac{24 \log(\ell+1)}{n} + \frac{48 \log(\ell+1)}{n}} \right) \leq \frac{4}{\ell^2}$  for all  $n = 1, 2, \dots$ .

*Proof.* We will analyze the cases  $\mu < 1$  and  $\mu \geq 1$  separately.

**Case-1:**  $\mu \leq 1$ . Let  $\delta = (\mu+1)\sqrt{\frac{6 \log(\ell+1)}{\mu n}}$ . First assume that  $\delta \leq \frac{1}{2}$ . Substituting the value of  $\delta$  in Theorem 5, we obtain,

$$\begin{aligned} \mathcal{P}(2\bar{X} \geq \mu) &\geq 1 - \frac{1}{\ell^2}, \\ \mathcal{P}\left(\bar{X} \leq \frac{3\mu}{2}\right) &\geq 1 - \frac{1}{\ell^2}, \\ \mathcal{P}\left(|\bar{X} - \mu| < (\mu+1)\sqrt{\frac{6\mu \log(\ell+1)}{n}}\right) &\geq 1 - \frac{2}{\ell^2}. \end{aligned}$$

From the above three results, we have,

$$\mathcal{P}\left(|\bar{X} - \mu| < \sqrt{\frac{48\bar{X} \log(\ell+1)}{n}}\right) \geq \mathcal{P}\left(|\bar{X} - \mu| < \sqrt{\frac{24\mu \log(\ell+1)}{n}}\right) \geq 1 - \frac{3}{\ell^2}. \quad (\text{D.9})$$

By assumption,  $\mu \leq 1$ . Therefore, we have  $\mathcal{P}\left(\bar{X} \leq \frac{3}{2}\right) \geq 1 - \frac{1}{\ell^2}$  and,

$$\mathcal{P}\left(\bar{X} \leq \sqrt{\frac{3\bar{X}}{2}}\right) \geq 1 - \frac{1}{\ell^2}.$$

Therefore, substituting above result in (D.9), we have

$$\mathcal{P} \left( |\bar{X} - \mu| > \max \left\{ \sqrt{\bar{X}}, \sqrt{\frac{2}{3}} \bar{X} \right\} \sqrt{\frac{48 \log(\ell+1)}{n}} \right) \leq \frac{4}{\ell^2}. \quad (\text{D.10})$$

Now consider the scenario, when  $(\mu+1)\sqrt{\frac{6 \log(\ell+1)}{\mu n}} > \frac{1}{2}$ . Then, we have,

$$\delta_1 \triangleq \frac{12(\mu+1)^2 \log(\ell+1)}{\mu n} \geq \frac{1}{2},$$

which implies,

$$\begin{aligned} \exp \left( -\frac{n\mu\delta_1^2}{2(1+\delta_1)(1+\mu)^2} \right) &\leq \exp \left( -\frac{n\mu\delta_1}{6(1+\mu)^2} \right), \\ \exp \left( -\frac{n\delta_1^2\mu}{6(1+\mu)^2} \left( 3 - \frac{2\delta_1\mu}{1+\mu} \right) \right) &\leq \exp \left( -\frac{n\mu\delta_1}{6(1+\mu)^2} \right). \end{aligned}$$

Therefore, substituting the value of  $\delta_1$  in Theorem 5, we have

$$\mathcal{P} \left( |\bar{X} - \mu| > \frac{48 \log(\ell+1)}{n} \right) \leq \frac{2}{\ell^2}.$$

Hence, from the above result and (D.10), it follows that,

$$\mathcal{P} \left( |\bar{X} - \mu| > \max \left\{ \sqrt{\bar{X}}, \sqrt{\frac{2}{3}} \bar{X} \right\} \sqrt{\frac{48 \log(\ell+1)}{n}} + \frac{48 \log(\ell+1)}{n} \right) \leq \frac{6}{\ell^2}. \quad (\text{D.11})$$

**Case 2:  $\mu \geq 1$**

Let  $\delta = \sqrt{\frac{12 \log(\ell+1)}{n}}$ , then by our assumption, we have  $\delta \leq \frac{1}{2}$ . From Theorem 5, it follows that,

$$\begin{aligned} \mathcal{P} \left( |\bar{X} - \mu| < \mu \sqrt{\frac{12 \log(\ell+1)}{n}} \right) &\geq 1 - \frac{2}{\ell^2} \\ \mathcal{P} (2\bar{X} \geq \mu) &\geq 1 - \frac{1}{\ell^2} \end{aligned}$$

Hence we have,

$$\mathcal{P} \left( |\bar{X} - \mu| < \bar{X} \sqrt{\frac{48 \log(\ell+1)}{n}} \right) \geq \mathcal{P} \left( |\bar{X} - \mu| < \mu \sqrt{\frac{12 \log(\ell+1)}{n}} \right) \geq 1 - \frac{3}{\ell^2}. \quad (\text{D.12})$$

By assumption  $\mu \geq 1$ . Therefore, we have  $\mathcal{P} (\bar{X} \geq \frac{1}{2}) \geq 1 - \frac{1}{\ell^2}$  and,

$$\mathcal{P} \left( \bar{X} \geq \sqrt{\frac{\bar{X}}{2}} \right) \geq 1 - \frac{1}{\ell^2}. \quad (\text{D.13})$$

Therefore, from (D.12) and (D.13), we have

$$\mathcal{P} \left( |\bar{X} - \mu| > \max \left\{ \bar{X}, \sqrt{\frac{\bar{X}}{2}} \right\} \sqrt{\frac{48 \log(\ell+1)}{n}} \right) \leq \frac{4}{\ell^2}. \quad (\text{D.14})$$

The result follows from (D.10) and (D.14).  $\square$

From the proof of Lemma D.1, the following result follows.

**Corollary D.1** Consider  $n$  i.i.d geometric random variables  $X_1, \dots, X_n$  with parameter  $p$ , i.e. for any  $i$ ,  $P(X_i = m) = (1-p)^m p \ \forall m = \{0, 1, 2, \dots\}$ . Let  $\mu = \mathbb{E}_\pi(X_i) = \frac{1-p}{p}$  and  $\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$ . If  $\mu \leq 1$ , then we have,

1.  $\mathcal{P} \left( |\bar{X} - \mu| > \sqrt{\frac{48\bar{X} \log(\ell+1)}{n}} + \frac{48 \log(\ell+1)}{n} \right) \leq \frac{5}{\ell^2}$  for all  $n = 1, 2, \dots$ .
2.  $\mathcal{P} \left( |\bar{X} - \mu| \geq \sqrt{\frac{24\mu \log(\ell+1)}{n}} + \frac{48 \log(\ell+1)}{n} \right) \leq \frac{4}{\ell^2}$  for all  $n = 1, 2, \dots$ .

**Proof of Lemma A.2** Fix  $i$  and  $\ell$ , define the events,

$$\mathcal{A}_{i,\ell} = \left\{ |\bar{v}_{i,\ell} - v_i| > \sqrt{48\bar{v}_{i,\ell} \frac{\log(\ell+1)}{|\mathcal{T}_i(\ell)|}} + \frac{48 \log(\ell+1)}{|\mathcal{T}_i(\ell)|} \right\}.$$

Let  $\bar{v}_{i,m} = \frac{\sum_{\tau=1}^m \hat{v}_{i,\tau}}{m}$ . Then, we have,

$$\begin{aligned} \mathbb{P}_\pi(\mathcal{A}_{i,\ell}) &\leq \mathbb{P}_\pi \left\{ \max_{m \leq \ell} \left( |\bar{v}_{i,m} - v_i| - \sqrt{48\bar{v}_{i,m} \frac{\log(\ell+1)}{m}} - \frac{48 \log(\ell+1)}{m} \right) > 0 \right\} \\ &= \mathbb{P}_\pi \left( \bigcup_{m=1}^{\ell} \left\{ |\bar{v}_{i,m} - v_i| - \sqrt{48\bar{v}_{i,m} \frac{\log(\ell+1)}{m}} - \frac{48 \log(\ell+1)}{m} > 0 \right\} \right) \\ &\leq \sum_{m=1}^{\ell} \mathbb{P}_\pi \left( |\bar{v}_{i,m} - v_i| > \sqrt{48\bar{v}_{i,m} \frac{\log(\ell+1)}{m}} + \frac{48 \log(\ell+1)}{m} \right) \\ &\stackrel{(a)}{\leq} \sum_{m=1}^{\ell} \frac{5}{\ell^2} \leq \frac{5}{\ell} \end{aligned} \tag{D.15}$$

where inequality (a) in (D.15) follows from Corollary D.1. The inequalities in Lemma A.2 follows from definition of  $v_{i,\ell}^{\text{UCB}}$ , Corollary D.1 and (D.15).  $\square$

Proof of Corollary B.1 is similar to the proof of Lemma A.2.

## E. Proof of Lemma 7.1

We follow the proof of  $\Omega(\sqrt{NT})$  lower bound for the Bernoulli instance with parameters  $\frac{1}{2}$ . We first establish a bound on KL divergence, which will be useful for us later.

**Lemma E.1** Let  $p$  and  $q$  denote two Bernoulli distributions with parameters  $\alpha + \epsilon$  and  $\alpha$  respectively. Then, the KL divergence between the distributions  $p$  and  $q$  is bounded by  $4K\epsilon^2$ ,

$$KL(p||q) \leq \frac{4}{\alpha} \epsilon^2.$$

$$\begin{aligned}
 KL(p\|q) &= \alpha \cdot \log \frac{\alpha}{\alpha + \epsilon} + (1 - \alpha) \log \frac{1 - \alpha}{1 - \alpha - \epsilon} \\
 &= \alpha \left[ \log \frac{1 - \frac{\epsilon}{1 - \alpha}}{1 + \frac{\epsilon}{\alpha}} \right] - \log \left( 1 - \frac{\epsilon}{1 - \alpha} \right) \\
 &= \alpha \log \left( 1 - \frac{\epsilon}{(1 - \alpha)(\alpha + \epsilon)} \right) - \log \left( 1 - \frac{\epsilon}{1 - \alpha} \right)
 \end{aligned}$$

using  $1 - x \leq e^{-x}$  and bounding the Taylor series for  $-\log 1 - x$  by  $x + 2x^2$  for  $x = \frac{\epsilon}{1 - \alpha}$ , we have

$$\begin{aligned}
 KL(p\|q) &\leq \frac{-\alpha\epsilon}{(1 - \alpha)(\alpha + \epsilon)} + \frac{\epsilon}{1 - \alpha} + 4\epsilon^2 \\
 &= \left( \frac{2}{\alpha} + 4 \right) \epsilon^2 \leq \frac{4}{\alpha} \epsilon^2
 \end{aligned}$$

Q.E.D.

Fix a guessing algorithm, which at time  $t$  sees the output of a coin  $a_t$ . Let  $P_1, \dots, P_n$  denote the distributions for the view of the algorithm from time 1 to  $T$ , when the biased coin is hidden in the  $i^{th}$  position. The following result establishes for any guessing algorithm, there are at least  $\frac{N}{3}$  positions that a biased coin could be and will not be played by the guessing algorithm with probability at least  $\frac{1}{2}$ . Specifically,

**Lemma E.2** *Let  $\mathcal{A}$  be any guessing algorithm operating as specified above and let  $t \leq \frac{N\alpha}{60\epsilon^2}$ , for  $\epsilon \leq \frac{1}{4}$  and  $N \geq 12$ . Then, there exists  $J \subset \{1, \dots, N\}$  with  $|J| \geq \frac{N}{3}$  such that*

$$\forall j \in J, \mathcal{P}_j(a_t = j) \leq \frac{1}{2}$$

Let  $N_i$  to be the number of times the algorithm plays coin  $i$  up to time  $t$ . Let  $P_0$  be the hypothetical distribution for the view of the algorithm when none of the  $N$  coins are biased. We shall define the set  $J$  by considering the behavior of the algorithm if tosses it saw were according to the distribution  $P_0$ . We define,

$$J_1 = \left\{ i \mid E_{P_0}(N_i) \leq \frac{3t}{N} \right\}, J_2 = \left\{ i \mid \mathcal{P}_0(a_t = i) \leq \frac{3}{N} \right\} \text{ and } J = J_1 \cap J_2. \quad (\text{E.1})$$

Since  $\sum_i E_{P_0}(N_i) = t$  and  $\sum_i \mathcal{P}_0(a_t = i) = 1$ , a counting argument would give us  $|J_1| \geq \frac{2N}{3}$  and  $|J_2| \geq \frac{2n}{3}$  and hence  $|J| \geq \frac{N}{3}$ . Consider any  $j \in J$ , we will now prove that if the biased coin is at position  $j$ , then the probability of algorithm guessing the biased coin will not be significantly different from the  $P_0$  scenario. By Pinsker's inequality, we have

$$|\mathcal{P}_j(a_t = j) - \mathcal{P}_0(a_t = j)| \leq \frac{1}{2} \sqrt{2 \log 2 \cdot KL(P_0\|P_j)}, \quad (\text{E.2})$$

where  $KL(P_0\|P_j)$  is the KL divergence of probability distributions  $P_0$  and  $P_j$  over the algorithm. Using the chain rule for KL-divergence, we have

$$KL(P_0\|P_j) = E_{P_0}(N_j)KL(p\|q),$$

where  $p$  is a Bernoulli distribution with parameter  $\alpha$  and  $q$  is a Bernoulli distribution with parameter  $\alpha + \epsilon$ . From Lemma E.1 and (E.1), we have that Therefore,

$$KL(P_0\|P_j) \leq \frac{4\epsilon^2}{\alpha},$$

Therefore,

$$\begin{aligned} \mathcal{P}_j(a_t = j) &\leq \mathcal{P}_0(a_t = j) + \frac{1}{2}\sqrt{2\log 2 \cdot KL(P_0\|P_j)} \\ &\leq \frac{3}{N} + \frac{1}{2}\sqrt{(2\log 2)\frac{4\epsilon^2}{\alpha}E_{P_0}(N_j)} \\ &\leq \frac{3}{N} + \sqrt{2\log 2}\sqrt{\frac{3t\epsilon^2}{N\alpha}} \leq \frac{1}{2}. \end{aligned} \tag{E.3}$$

The second inequality follows from (E.1), while the last inequality follows from the fact that  $N > 12$  and  $t \leq \frac{N\alpha}{60\epsilon^2}$  Q.E.D..

**Proof of Lemma 7.1** . Let  $\epsilon = \sqrt{\frac{N}{60\alpha T}}$ . Suppose algorithm  $\mathcal{A}$  plays coin  $a_t$  at time  $t$  for each  $t = 1, \dots, T$ . Since  $T \leq \frac{N\alpha}{60\epsilon^2}$ , for all  $t \in \{1, \dots, T-1\}$  there exists a set  $J_t \subset \{1, \dots, N\}$  with  $|J_t| \geq \frac{N}{3}$  such that

$$\forall j \in J_t, P_j(j \in S_t) \leq \frac{1}{2}$$

Let  $i^*$  denote the position of the biased coin. Then,

$$\mathbb{E}_\pi(\mu_{a_t} | i^* \in J_t) \leq \frac{1}{2} \cdot (\alpha + \epsilon) + \frac{1}{2} \cdot \alpha = \alpha + \frac{\epsilon}{2}$$

$$\mathbb{E}_\pi(\mu_{a_t} | i^* \notin J_t) \leq \alpha + \epsilon$$

Since  $|J_t| \geq \frac{N}{3}$  and  $i^*$  is chosen randomly, we have  $P(i^* \in J_t) \geq \frac{1}{3}$ . Therefore, we have

$$\mu_{a_t} \leq \frac{1}{3} \cdot \left(\alpha + \frac{\epsilon}{2}\right) + \frac{2}{3} \cdot (\alpha + \epsilon) = \alpha + \frac{5\epsilon}{6}$$

We have  $\mu^* = \alpha + \epsilon$  and hence the  $Regret \geq \frac{T\epsilon}{6}$ . □

**Lemma E.3** *Let  $L$  be the total number of calls to  $\mathcal{A}_{MNL}$  when  $\mathcal{A}_{MAB}$  is executed for  $T$  time steps. Then,*

$$\mathbb{E}(L) \leq 3T.$$

*Proof.* Let  $\eta_\ell$  be the random variable that denote the duration, assortment  $S_\ell$  has been considered by  $\mathcal{A}_{MAB}$ , i.e.  $\eta_\ell = 0$ , if we no arm is pulled when  $\mathcal{A}_{MNL}$  suggested assortment  $S_\ell$  and  $\eta_\ell \geq 1$ , otherwise. We have

$$\sum_{\ell=1}^{L-1} \eta_\ell \leq T.$$

Therefore, we have  $\mathbb{E}\left(\sum_{\ell=1}^{L-1} \eta_\ell\right) \leq T$ . Note that  $\mathbb{E}(\eta_\ell) \geq \frac{1}{2}$ . Hence, we have  $\mathbb{E}(L) \leq 2T + 1 \leq 3T$ .

□