PASTA: Pessimistic Assortment Optimization

Juncheng Dong *1 Weibin Mo *2 Zhengling Qi 3 Cong Shi 4 Ethan X. Fang 5 Vahid Tarokh 1

Abstract

We consider a class of assortment optimization problems in an offline data-driven setting. A firm does not know the underlying customer choice model but has access to an offline dataset consisting of the historically offered assortment set, customer choice, and revenue. The objective is to use the offline dataset to find an optimal assortment. Due to the combinatorial nature of assortment optimization, the problem of insufficient data coverage is likely to occur in the offline dataset. Therefore, designing a provably efficient offline learning algorithm becomes a significant challenge. To this end, we propose an algorithm referred to as Pessimistic ASsortment opTimizAtion (PASTA for short) designed based on the principle of pessimism, that can correctly identify the optimal assortment by only requiring the offline data to cover the optimal assortment under general settings. In particular, we establish a regret bound for the offline assortment optimization problem under the celebrated multinomial logit model. We also propose an efficient computational procedure to solve our pessimistic assortment optimization problem. Numerical studies demonstrate the superiority of the proposed method over the existing baseline method.

1. Introduction

One of the most critical problems faced by a seller is to select products for presentation to potential buyers. Often faced with limited display spaces and storage costs in both

Proceedings of the 40th International Conference on Machine Learning, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

brick-and-mortar and online retailing, the seller needs to carefully choose a set of products from the vast collection of all available products for displaying to its customers. In this line, the problem of selecting an assortment, i.e., a collection of products from all available products, in order to maximize the seller's revenue is the assortment optimization problem. Obviously, the choice behavior of customers (McFadden, 1981) is of great importance in the problem of assortment optimization. Without loss of generality, we assume the choice of each customer can be described by a preference vector $\boldsymbol{\theta}$. This subsumes the seminal multinomial logit (MNL) model (McFadden, 1973) which is arguably the most well-studied and widespread models in assortment optimization literature (Please see Section 6 for more details) (Talluri & van Ryzin, 2004; Caro & Gallien, 2007; Rusmevichientong et al., 2010; Davis et al., 2013; Chen et al., 2021a; Aouad et al., 2022).

In practice, θ is often unknown and needs to be estimated. Assuming no historical data of customers, *dynamic assortment optimization* adaptively learns θ in a trial-and-error fashion by updating the assortment and observing the subsequent choices of customers sequentially (Caro & Gallien, 2007; Chen et al., 2020; Rusmevichientong et al., 2020; Chen et al., 2021b; Li et al., 2022). Meanwhile, in our era of Big Data, companies often collect abundant customer data. Therefore, it is often in companies' best interest to learn from the existing (potentially massive) offline datasets rather than starting from scratch. Moreover, offline learning is beneficial since online exploration can sometimes be expensive or infeasible. Hence, we take the first stab to formally study the following important question faced by every seller.

Research Question: Given a pre-collected offline dataset of historically offered assortment, customers choices, and revenue, how can we find an efficient and theoretically justified offline algorithm to estimate the optimal assortment set without unrealistic assumptions on the offline dataset?

When the dataset is not adaptively collected, it is not uncommon to encounter the challenge of insufficient coverage of data. For estimators such as the maximum likelihood estimator (MLE) to approximate θ^{\prime} accurately, the offline dataset must include sufficiently many assortments and customer choices. In other words, the data-collecting process needs to sufficiently explore different assortments (by the

^{*}Equal contribution ¹Department of Electrical and Computer Engineering, Duke University, Durham, NC 27705, United States. ²Mitchell E. Daniels, Jr. School of Business, Purdue University, West Lafayette, IN 47907, United States. ³Department of Decision Sciences, George Washington University, Washington, DC 20052, United States. ⁴Herbert Business School, University of Miami, Coral Gables, FL 33146, United States. ⁵Department of Biostatistics and Bioinformatics, Duke University, Durham, NC 27705, United States. Correspondence to: Zhengling Qi <qizhengling@gwu.edu>.

seller) and different choices (by the customers). This is unlikely to happen for offline datasets because the seller would not choose unreasonable assortments whose expected revenues are obviously suboptimal, and the customers would not choose products against their preferences.

Major Contributions. The main contribution of this work is two-fold. First, based on the principle of pessimism, we propose the Pessimistic ASsortment opTimizAtion (PASTA for short) framework, which correctly identifies the optimal assortment. In particular, our framework only requires that the offline dataset covers the optimal assortment set instead of all possible (combinatorially many) assortment sets. Second, we derive the first finite-sample regret bound for offline assortment optimization under the multinomial logit (MNL) model (Please see Section 6), one of the most widely used models for modeling customers' choices. We subsequently propose an algorithm, also with the name *PASTA*, that can efficiently solve the pessimistic assortment optimization problems. Experiments on the simulated datasets (so that θ is known) corroborate the efficacy of pessimistic assortment optimization.

Paper Organization. We briefly review the related work in Section 2 and the preliminaries in Section 3. We propose the pessimistic assortment optimization in Section 4. In Section 5, we present the theoretical results. In Section 6, we study pessimistic assortment optimization under the MNL model as a concrete example. In Section 7, we propose an algorithm that can solve the problem efficiently. We provide experimental results in Section 8, after which we conclude.

2. Related Work

Assortment Optimization. The assortment optimization problem under the MNL model without any constraints was first studied in (Talluri & van Ryzin, 2004). Then more complicated assortment optimization problems under various types of constraints, including space requirement (Rusmevichientong et al., 2009) and cardinality (Rusmevichientong et al., 2010), were considered. (Davis et al., 2013) proposed a linear programming (LP) formulation of the assortment optimization problem that includes several previous works as special cases corresponding to different constraints in the formulation of LP. This line of work assumes that the true parameters of the customer models are known (or at least can be accurately estimated) (Gallego & Topaloglu, 2014; Feldman & Topaloglu, 2015; Flores et al., 2019; Désir et al., 2020; Liu et al., 2020; Aouad et al., 2021). Another closely related line of work is *dynamic assortment* optimization (Caro & Gallien, 2007). In the setting of dynamic assortment optimization, the seller without any prior information about the customers, has finite selling horizons in which it observes the choices of customers and, based on the observed behaviors, optimize their assortments in

an adaptive, trial-and-error fashion (Sauré & Zeevi, 2013; Wang et al., 2018; Chen et al., 2021a; Rusmevichientong et al., 2020; Chen et al., 2020; 2021b). In comparison with the online setting used in dynamic assortment optimization, our work departs from the existing literature by focusing on the offline setting where the seller only has collected datasets but not any control on the data-collecting process.

Pessimism in Offline Learning. The principle of pessimism has been successfully used in reinforcement learning (RL) for finding an optimal policy with pre-collected datasets. On the empirical side, it has helped with improving the performance of both the model-based approach and value-based approach in offline setting (e.g., Yu et al., 2020; Kidambi et al., 2020; Kumar et al., 2020). The importance of pessimism has been analyzed and verified theoretically in the setting of RL (Jin et al., 2021; Fu et al., 2022). Our work main contribution is to take a pessimistic approach to assortment optimization problems and demonstrate its empirical and theoretical values. Moreover, our work differs from the above works by focusing on a decision-making problem with exponentially many choices.

3. Preliminary

Let $rN ext{ s}^u ext{ t} ext{ 1, 2, ..., } ext{ Nu}$ denote the set of N distinct items. For each item i, a feature vector $x_i ext{ P} ext{ R}^d$ is available. Assume that $tx_i u_{PrN ext{ s}}$ are fixed vectors. Denote the collection of all possible assortments under *consideration* by $S ext{ D} ext{ 2}^{rN ext{ s}} ext{ztHu}$. For the offline data, we define a random vector pS, A, Rq from each customer, where $S ext{ D} ext{ r} ext{ N ext{ s}} ext{ denotes an assortment presented to the customer, } A ext{ P S Y t0u}$ denotes the item purchased by the customer for $A ext{ P S} ext{ (A " 0 where no purchase is made), and } R$ denotes the corresponding revenue. The ultimate goal of assortment optimization is to find an optimal set of items s r P S for all customers to maximize the expected revenue. A specific goal of this work is to study how to leverage the offline data, which consists of i.i.d. samples of the random triplet pS, A, Rq in order to learn an optimal assortment.

For the assortment optimization with offline data, a fundamental question is to estimate the expected revenue for an unexplored assortment SPS. This amounts to addressing the causal relationship between assortment and revenue. Under the celebrated potential outcome framework (Rubin, 1974), let the random variable RpSqbe the potential revenue under an intervention that the assortment is set to be SPS. Our goal is to find an optimal assortment

$$s^{\circ}$$
 Parg max Er R p s qs

Note that the expected potential revenue ErRpsqsdefined in the counterfactual world may not be identifiable from the observed data without additional assumptions. Throughout this paper, we make the following standard consistency and

un-confoundedness assumptions in causal inference.

Assumption 3.1. [CONSISTENCY] With probability one, the observed revenue coincides with the potential revenue of the observed assortment. That is, R " RpSqalmost surely.

Assumption 3.2. [UN-CONFOUNDEDNESS] The potential revenues are independent variables of the observed assortment, i.e., tRpsqups KKS.

Assumption 3.1 ensures that the observed revenue is consistent with the potential revenue of purchasing item A p‰0q or no purchase if A " 0, under the observed assortment S. Assumption 3.2 rules out possible unobserved factors that could confound the causal effect of assortment on revenue¹.

Denote $\pi_S p q$ " PpS " Sq as the probability of observing assortment S in the offline data. To non-parametrically identify ErRpSqs for every SPS, we further require the following positivity assumption (Imbens & Rubin, 2015).

Assumption 3.3. [POSITIVITY EVERYWHERE] For all S P S, the probability $\pi_p Sq$ of observing assortment S is positive (i.e. $\pi_p Sq$ **q**0).

Assumption 3.3 requires that every assortment can be observed with a positive chance in the offline data. *This is a strong assumption that will be later relaxed it Assumption 5.1 (I), i.e., requiring positivity only at optimum.* With Assumptions 3.1–3.3, we can identify the effect of an assortment set via inverse propensity score weighting (Rosenbaum & Rubin, 1983): for any SPS,

$$ErRpSqs "E" \frac{IpS " sqR}{\pi_SpSq}^*, \qquad (1)$$

where the expectation in the right-hand-side is taken with respect to the data distribution of pS, A, Rq However, when the number of possible assortments |S| grows exponentially in N, Assumption 3.3 rarely holds for all SPS in practice, given potentially limited offline data particularly when N is large. Moreover, when an assortment S corresponds to an inferior expected revenue ErRpSqS it may not be considered by the seller at all. As a consequence, the probability of observing such an assortment TSpSq is zero. These may prevent us from estimating (1) for every assortment SPS.

We may tempt to use the following identification strategy:

ErRpSqs "EpR|S" sq pby Assumptions 3.1-3.2q
" ErEpR|S" s, Aqs "
$$\pi_A p^i |s; xq'_{s,i},$$
 (2)

where \mathbf{X} "t \mathbf{X}_{j} \mathbf{U}_{PrNS} are the features across items, π_{A} $\mathbf{p}^{i}|\mathbf{S}$; $\mathbf{x}\mathbf{q}$ "PpA" $i|\mathbf{S}$ " Sq is the customer's choice probability (McFadden, 1973) of purchasing the i-th item given an assortment \mathbf{S} , $r_{s,i}$ "EpR|S" \mathbf{S} , A" iq is the conditional expected revenue given the assortment \mathbf{S} with the i-th item being purchased. For ease of notation, we omit the features \mathbf{X} in π_{A} when there is no confusion. Identifying $\mathbf{E}_{r}\mathbf{P}_{s}\mathbf{S}$ as above requires the knowledge of π_{A} \mathbf{p}^{i} |Sq and $r_{s,i}$, which can be learned from data. Although such an identification approach does not explicitly depend on π_{s} \mathbf{p}^{s} q full identification of π_{A} \mathbf{p}^{i} |Sq and $r_{s,i}$ requires the positivity of π_{s} \mathbf{p}^{s} q for every \mathbf{S} \mathbf{P} \mathbf{S} as assumed above.

Despite the aforementioned challenge of insufficient coverage over assortments, we argue that finding an optimal assortment s may not necessarily require $\pi_s p q$ a everywhere but only at the optimal assortment s. In particular, based on (2), when computing

$$\overset{\circ}{S} \operatorname{P} \underset{SPS}{\operatorname{arg max}} \pi_{A} p^{i} |Sq^{i}_{S,i}, \qquad (3)$$

we may not necessarily need to estimate $\pi_A \not p \mid Sq$ and $r_{s,i}$ well for $s \not \sim s^s$, as long as *sub-optimal assortments can* be safely ruled out during the optimization. Our insight is that the estimation of $\pi_A \not p \mid Sq$ and $r_{s,i}$ for the less seen assortment s in the data often incurs large errors. Deploying pessimism by taking the estimation error into consideration can rule out those assortments (Jin et al., 2021), while standard predict-then-optimize (Bertsimas & Kallus, 2020) or empirical maximization approaches (Zhao et al., 2012) may suffer from an overestimation of ErR pSqs Hence, in our proposed pessimistic assortment optimization framework, we only require the positivity at optimum $\pi_s \not p^s \not q \not q \not q$, which is a much weaker assumption than that of Assumption 3.3.

In this paper, we focus on handling the estimation error from $\pi_A p^i | \text{Sq}$ while assuming that $r_{s,i}$ is known. This is a typical assumption in the literature of assortment optimization (Talluri & van Ryzin, 2004; Davis et al., 2013; Flores et al., 2019; Aouad et al., 2021). Our framework can be naturally extended to the scenario where we need to estimate $r_{s,i}$'s. For optimization tractability, we further assume that $r_{s,i}$ " r_i that the expected revenue depends only on the purchased item but not on the underlying assortment. This assumption is reasonable in many applications where the revenue is a deterministic consequence of a purchased item. This can also be easily extended under our pessimism framework but could result in a more complicated assortment optimization problem.

Below, without loss of generality, we assume that $r_i \not\in 0$ for $i \not\in N$ s, while $r_0 \cap 0$ (no purchase incurs zero revenue). For any vector x, let x^J and $|x||_2$ respectively denote the transpose and ℓ_2 -norm of x. For any set A, let |A| denote the cardinal number of A. For any two sequences $t^{\varpi}p^{\eta}qu_{EA}$

¹With the observed features $t^{X_j} U_{PrN s}$, it can be possible to relax Assumption 3.2 to a more plausible condition: the independence holds conditional on the observed features. However, for notation simplicity, without loss of generality, we consider Assumption 3.2.

and $t \gamma p^{n} q u_{h \in 1}$, we write $\varpi p^{n} q$ Á $\gamma p^{n} q$ (resp. $\varpi p^{N} q$ À $\gamma p^{n} q$) whenever there exist constants $c_{1} \neq 0$ (respectively $c_{2} \neq 0$ such that $\varpi p^{n} q \in c_{1} \gamma p^{n} q$ (resp. $\varpi p^{n} q \quad d^{\prime} c_{2} \gamma p^{n} q$). Moreover, we write $\varpi p^{n} q \approx \gamma p^{n} q$ whenever $\varpi p^{n} q \quad A \gamma p^{n} q$ and $\varpi p^{n} q \quad A \gamma p^{n} q$

4. Pessimistic Offline Assortment Optimization

In this section, we introduce our pessimistic offline assortment optimization framework. To this end, based on Eq. (2), we first estimate the choice probability $\pi_A p^i | Sq$ from offline data. Subsequently we calculate optimizing values in optimization problem (3) using a plug-in estimator of $\pi_A p^a | q$ Consider a generic form of model $\pi_A p^a | q$ with the unknown true parameter $q^a | q$. Again, for ease of notation, we omit the features $q^a | q$ when there is no confusion. We remark that $q^a | q$ could be either finite-dimensional or infinite-dimensional. Given an offline dataset $q^a | q | q$ the sample size, one can estimate the model parameter $q^a | q$ via maximum likelihood estimator (MLE). Specifically, define the likelihood-based loss function $q^a | q^a | q$ as

$$p_n p\theta_q \text{ "} \frac{1}{n} \frac{\hat{y}^n}{\prod_{i=1}^n \log n_i pA_i |S_i; \theta_q|}$$

Then the MLE of the unknown parameter θ is $\theta_{ML,n}$ P arg min $_{\theta P\Theta}$ pq where Θ is a pre-specified parameter space. Let

$$V_{\text{ps}}; \boldsymbol{\theta}_{\text{ML},n} \ \mathbf{q}^{\text{u}} \ \overset{\ddot{\mathbf{y}}}{\underset{i \text{ ps}}{\sum}} \pi_{\mathbf{A}} \, \mathbf{p}^{i} \, | \mathbf{s}; \boldsymbol{\theta}_{\text{ML},n} \ \mathbf{q}^{i}.$$

Here, we define V_{ps} ; θ qas the *value function* of s with the customer choice model for π_A depending on the parameter θ . The plug-in estimator of the optimal assortment based on (3) is

$$\beta_{ML,n}$$
 Parg max $V_{DS}; \theta_{ML,n}$ q.

The MLE-based approach first plugs in the MLE of θ , and then directly optimizes the corresponding estimated value function.

As discussed before, a disadvantage of the above estimate-then-optimize approach is that the estimation error of $\theta_{ML,n}$ caused by insufficient data coverage may result in the overestimation of V_{ps} ; θ q, which will propagate to downstream optimization. Alternatively, we can quantify the estimation uncertainty by considering the following likelihood-ratio-test-based confidence region (Owen, 1990):

$$\Omega_n p \alpha_n q$$
 " $t \theta P \Theta : p p p q$ ' $p p p_{ML,n} q d \alpha_n u$

where α_n **a** 0 is pre-specified. Later we analyze the MNL model as a special case (Please see Section 6). With α_n chosen as Opd[nq] we establish in Theorem 6.1 that θ

 $\Omega_n p^{\alpha}_n q$ with high probability. Such a guarantee does not require any data coverage assumption on assortments.

For now on, for simplicity, we drop α_n and write Ω_n for $\Omega_n p \alpha_n q$ when there is no ambiguity.

In order to robustify assortment optimization against plug-in estimation errors, we consider a pessimistic version of (3) by taking the estimation uncertainty from Ω_n into account. Specifically, we propose the **Pessimistic Assortment opTimizAtion (PASTA)** by solving

$$\beta_{PASTA,n} \quad P \underset{SPS}{\operatorname{arg max min}} V_{PS}; \ \theta q \tag{4}$$

Here, for a fixed assortment S PS, the inner layer of minimization computes the worst-case value among all possible model parameters θ within the confidence set Ω_n . In particular, if the estimated value V_{ps} ; $\theta_{ML,n}$ q for S is highly uncertain due to insufficient data coverage, the worst-case value $\min_{\theta \in \Omega_n} V_{ps}$; θ_q is likely much smaller than V_{ps} ; θ_q In that case, the outer layer of (4) may prefer another assortment with a relatively higher worst-case value. In this way, the inner layer of (4) rules out those assortments with less frequency in the offline data. Hence, one essential advantage of such a strategy is that it avoids an overestimation of the value function. In other words, by the plug-in approach, with a non-negligible chance, the estimated value V_{ps} ; $\theta_{ML,n}$ q can be much larger than the truth V_{ps} ; θ q which further leads to a possibly sub-optimal assortment but optimized by the MLE-based approach. In contrast, PASTA is aware of insufficient data coverage, and hence more pessimistic about those highly uncertain value estimates. In the next section, we theoretically analyze the advantage of the PASTA approach.

5. Theoretical Results

In this section, we show that the PASTA method (4) enjoys a generic regret guarantee under a weak assumption of *positivity at optimum* that is $\pi_s p^s$ q q 0. Specifically, given $\beta_{PASTA,n}$ in (4), we adopt the following regret as the performance metric to evaluate the PASTA's performance

$$R_{p_{\mathbf{p}_{\mathsf{ASTA}}n}} q \text{ "} V_{p_{\mathbf{s}}} : \mathring{\theta} q \text{ '} V_{p_{\mathbf{p}_{\mathsf{ASTA}}n}} : \mathring{\theta} q$$

We aim to derive a regret bound for $Rpp_{PASTA,n}$ q under generic conditions. Denote Lppq " $Er' log n_A p_A | S; eqs$ as the population loss function. All detailed proofs can be found in the Appendix 9.

We first show that whenever $\theta' P \Omega_n$ that the confidence region covers the true parameter, the regret of the PASTA method can be calibrated by the worst-case estimation error among $\theta P \Omega_n$ of the value function at the optimal assortment S.

Lemma 5.1. Let **\(\rightarrow{9}\)**PASTA, be the solution by the PASTA

method defined in (4). If $\theta^{\circ} P \Omega_n$, then

$$R_{p_{\text{PASTA},n}} \neq d_{\text{PPO}_{n}}^{\text{max}} \forall p^{s^{\circ}}; \theta^{\circ} \neq Vp^{s^{\circ}}; \theta \neq Vp^{s^{\circ}}$$

Proof of Lemma 5.1.

Next, we impose the following assumptions to obtain the regret guarantee of our algorithm.

Assumption 5.1.

(I) [POSITIVITY AT OPTIMUM] The probability of observing the optimal assortment is positive, that is, $\pi_s p^s$ q q0. (II) [LIKELIHOOD -BASED CONCENTRATION] For any 0 ă δ ă 1, with probability at least 1 δ , we have: (1) θ P Ω_n , and (2)

$$\sup_{\theta \neq \Omega_n} \mathring{L} p\theta q ' L p\theta ' q ' p_n p\theta q ' p_n p\theta ' q \mathring{d} \alpha_n.$$

We emphasize that *PASTA* only requires the positivity at optimum. Compared to the positivity at all assortments in Assumption 3.3, our Assumption 5.1 (I) is much weaker and hence more plausible to be satisfied. Assumption 5.1 (II) is a generic condition for likelihood-based concentration. We later justify that (II) above indeed holds under the general MNL model in Theorem 6.2. In particular, Statement (1) of Part (II) requires the validity of the likelihood-ratio-test-based confidence region Ω_n while Statement (2) of Part (II) requires the concentration of the likelihood-based localized empirical process (van der Vaart & Wellner, 1996).

The positivity at optimum is associated with a finite constant C_{s^*} " $1\{\pi_s ps^*\}$ q related to the learning performance. We also denote r_{s^*} " $\max_{j \in S^*} r_j$ as the largest possible revenue among all items in s^* . Notice that both constants C_{s^*}

and $r_{s^{\circ}}$ depend on the optimal assortment s° only. In the following lemma, we establish the estimation error bound at the optimal assortment s° .

$$V_{ps^{\circ}}$$
; θ q $V_{ps^{\circ}}$; θ q À $r_{s^{\circ}}$ $C_{s^{\circ}}$? $\overline{\alpha_n}$.

Combining Lemmas 5.1 and 5.2, we summarize the regret bound for PASTA in the following theorem.

Theorem 5.3. Under Assumption 5.1, for any 0 \check{a} $\check{\delta}$ \check{a} 1, with probability 1 $\check{\delta}$, we have

$$Rpp_{PASTA,n} q \lambda r_{s^{\circ}} C_{s^{\circ}} ? \overline{\alpha_n}.$$

6. Application: Multinomial Logit Model

In this section, we consider the Multinomial Logit Model (MNL) for customer choices $\pi_A p^2 | s_q$ This is one of the most widely used models in assortment optimization literature (Feng et al., 2022). Under the MNL model, we will verify Assumption 5.1 (II) and establish the regret bound for PASTA in this case.

Given the item-specific features $\mathbf{t}^{X_i} \mathbf{u}_{\mathsf{Pr}^N s}$, MNL assumes that customer's preference for the i-th item is proportional to $\mathsf{expp}^{X_i^J} \theta^c \mathbf{q}$ where $\theta^c \mathsf{P} \Theta$ is the underlying unknown parameter. Here, we assume that the parameter space $\Theta \check{\mathbf{D}} \mathsf{R}^d$ is compact with θ_{max} " $\mathsf{sup}_{\theta\mathsf{P}\Theta} \} \theta \}_2 \check{\mathbf{a}}$ '8 . Given an assortment s, the customer choice probability under MNL is given by

$$\pi_A \not p \mid s; \mathring{\theta} q " \frac{expp^{X_J} \mathring{\theta} q}{1 \mid ps expp^{X_J} \mathring{\theta} q}, \mathring{\theta} Ps.$$
 (5)

Moreover, the probability of no-purchase is normalized to $\pi_A p0|s$; $\theta \neq 1$ [p1 $\int_{j ps} expp^{x_j} \theta \neq qq$ Based on (3) and the MNL model (5), the objective function for assortment optimization can be written as

$$V_{ps}$$
; $\theta_q = \frac{\tilde{r}_{ips} r_i \exp x_i^j \theta_q}{1 + \frac{1}{1} \exp x_i^j \theta_q}$

We first justify Statement (1) of Assumption 5.1 under the MNL model. To this end, given the compactness of Θ , there exists a finite constant C_A **a** 0 such that for all $\theta P\Theta$, SPS and iPS, we have $1|\pi_A p^i|_S$; $\Theta_A d^iC_A$.

Lemma 6.1. Consider the MNL model (5) with a compact set Θ . Assume that θ P Θ . For any 0 \check{a} $\check{\delta}$ \check{a} 1, with probability at least 1' δ , we have

$$\mathbf{p}_n \mathbf{p} \theta^{\circ} \mathbf{q} \cdot \mathbf{p}_n \mathbf{p} \theta_{ML,n} \mathbf{q} \lambda \frac{C_A d}{n} \log \frac{\theta_{max}}{\delta}$$

Lemma 6.1 suggests that, with α_n chosen as $\frac{C_A d}{n} \log \frac{\theta_{\text{max}}}{\delta}$, we can guarantee that θ P Ω_n with high probability, which justifies Statement (1) of Assumption 5.1 (II). In particular, the order of α_n is $\text{Opd}[n_{\mathbf{q}}]$ Notice that Lemma 6.1 does not depend on the distribution of S, which implies that no data coverage assumption on the observed assortments is required. The assumption that θ P Θ for Θ a compact set requires that given any assortment, every product has a chance of being selected by the customer in the data. This is a mild requirement as θ is always finite.

Next, we justify Statement (2) of Assumption 5.1 (II) in the following theorem.

$$\sup_{\theta p\Omega_n} \overset{\check{\downarrow}}{\stackrel{\backprime}{\not}} L p\theta q \overset{\backprime}{} L p\theta \overset{\backprime}{} q \overset{\backprime}{} \overset{\backprime}{\not} D_n p\theta q \overset{\backprime}{} \stackrel{\backprime}{\not} D_n p\theta \overset{\backprime}{} q \overset{\backprime}{\stackrel{\backprime}{\not}} d \overset{\backprime}{} \alpha_n.$$

Finally, with the Assumption of positivity only at optimum (Assumption 5.1 (I)), we can apply Theorem 5.3 to establish the regret bound for PASTA in MNL.

Theorem 6.3. Consider the MNL model (given in Equation (5)). Assume that the conditions in Lemma 6.2 hold and that $\pi_S p^s \neq 0$. Fix $a \delta p \neq 0$, $\pi_S p^s \neq 0$, $\pi_S p^s \neq 0$. Then with $\pi_S p^s \neq 0$ is output of PASTA with $\pi_S p^s \neq 0$ in $\pi_S p^s \neq 0$. Then with probability at least $\pi_S p^s \neq 0$.

$$R p_{\text{PASTA},n} \neq \lambda r_{s^{\circ}} C_{s^{\circ}} \frac{c}{\frac{C_A d}{n} \log \frac{\theta_{\text{max}}}{\delta}}$$

We remark that under the MNL model (given in Equation (5)), the order of regret is $Op^{\alpha} \overline{a(n)}$ This is due to the concentration rate of MNL's empirical likelihood ratio in Lemma 6.1. Such a rate of regret bound matches those in the literature under parametric model assumptions (Qian & Murphy, 2011; Mo & Liu, 2022). However, existing literature requires the positivity $\pi_S p q$ a 0 at every P S. In contrast, Theorem 6.3 only requires positivity $\pi_{\rm S} p {\rm s}^{\circ} {\rm q}$ a 0 at the optimal assortment S° . Furthermore, we can show that $\min_{i \in \mathbb{N}} \pi_A p^i \mid S$ " r N s, θq d' $1[N \text{ for any } \theta P \Theta, \text{ which}]$ implies that $C_A
i N$. Therefore, our regret is of order at least \overline{N} , where N is the total number of available items. It is an interesting problem to establish the minimax lower bound of offline assortment optimization in terms of N, n, d and the cardinal number of S°. This will investigated in a subsequent work.

7. PASTA Algorithm

In this section, we propose an efficient algorithm for solving the max-min problem given in Optimization Problem (4) for the MNL model. Specifically, let

$$V_{ps}$$
; θ_q " $\frac{\ddot{y}}{1} \frac{r_j \exp p^{X_i^J} \theta_q}{1 - r_{ps} \exp p^{X_j^J} \theta_q}$

and given the confidence set Ω_n , we wish to solve

$$\max_{\text{SPS}} \min_{\theta \in \Omega_n} V_{\text{ps}}; \theta_{\text{q}}$$

The proposed iterative algorithm is executed for a maximum of T iterations. At the t-th iteration, given S_t and θ_t from the previous iteration, we consecutively execute the following two steps:

- Step 1: Compute the optimal assortment S_{t-1} given θ_t (see Section 7.1).
- Step 2: Compute the optimal θ_{t-1} using S_{t-1} (see Section 7.2).

The corresponding pseudo-code is presented in Algorithm 1 below.

Algorithm 1 PASTA

Input: offline dataset $\operatorname{tp} S_i$, A_i , $R_i \operatorname{qq}_{u_1}^n$; α_n ; $\operatorname{tr}_i \operatorname{q}_{i_{u_1}}^N$; $\operatorname{tx}_i \operatorname{q}_{i_{u_1}}^N$; maximum number of iterations T

Output: the solution to pessimistic assortment optimization **6**

tion
$$\beta$$

$$\mathbf{p}_{n} p \theta q^{\alpha'} \stackrel{?}{=} \frac{1}{n} \int_{i=1}^{n} \log \pi_{k} p A_{i} | S_{i}; \theta q$$

$$\mathbf{p}_{\mathsf{ML},n} \quad \mathbf{D} \quad \text{arg } \min_{\theta P \Theta} \mathbf{p}_{n} p \theta q$$

$$\Omega_{n} \quad \mathbf{D} \quad t \theta P \Theta : \mathbf{p}_{n} p \theta q \stackrel{?}{=} \mathbf{p}_{\mathsf{ML},n} \quad q \quad d^{i} \alpha_{n} \mathbf{u}$$

$$t \quad \mathbf{D} \quad 0; \theta \quad \mathbf{D} \quad \mathbf{p}_{\mathsf{ML},n} \qquad /* \quad \text{Initialize } \theta_{0} \text{ as } \mathbf{p}_{\mathsf{ML},n} \quad */$$

$$\mathbf{for} \quad t \quad 1 \quad \mathbf{to} \quad T \quad \mathbf{do}$$

$$S_{t} \quad \mathbf{D} \quad \mathbf{SolveLPp} \theta_{t'-1}, t^{r_{i}} \mathbf{u}_{i''-1}^{N}, t^{x_{i}} \mathbf{u}_{i''-1}^{N} \mathbf{q}$$

/* Section 7.1 */
$$\theta_t \ni \mathbf{SolveGDpS}_t, \Omega_n, \mathsf{tf}_i \mathsf{u}_{i=1}^N, \mathsf{t}_i \mathsf{x}_i \mathsf{u}_{i=1}^N \mathsf{qq}$$
/* Section 7.2 */

end for $S \not D S_T$

7.1. Optimal Assortment Computation

Given the MNL model parameter θ_t , computing the assortment S_{t-1} that maximizes the expected revenue can be formulated as a linear programming (LP) problem.

Suppose that an assortment S can be represented by an N-dimensional binary vector Y P t0, tu^N where Y_j " 1 if and only if j P S. Suppose that S P S corresponds to the following feasible set for Y with M linear inequality constraints:

where the matrix of constraint coefficients $ra_{ij} \, g_{PrM \, gj \, PrN \, g}$ is a totally unimodular matrix (Pang, 2017). In other words,

based on the one-to-one correspondence between S and Y. we have SPS if and only if $YP\Gamma$.

Next, we denote V_i " expp $X_i^J \theta_t q$ as the preference score for the i-th item. The customer choice probability under the MNL model (5) becomes $\pi_A p^i | \text{Sq}^{ii} = \sum_{i \neq s} \frac{v_i}{v_i}$. The optimization for S_{t-1} can be formulated as

$$\max_{y \in \Gamma} \frac{\overset{i}{\underset{i \in \Gamma N_{S}}{r_{i}}} r_{i} v_{i} y_{i}}{1 \overset{i}{\underset{i \in \Gamma N_{S}}{r_{i}}} v_{i} y_{i}}, \tag{6}$$

which is equivalent to the following linear programming problem (Davis et al., 2013):

In particular, we can recover the optimal solution to Problem (6), denoted as Y° , using the optimal solution to Problem (7), denoted by W° , via the following formula:

$$y_{j}^{\circ}$$
 " $\frac{w_{j}^{\circ}}{v_{i} w_{0}^{\circ}}$ @ P N s (8)

To conclude, at the t-th iteration, in order to compute an optimal assortment S_{t-1} for a given θ_t , we first solve an LP problem in (7) for W° . Then we recover Y° via (8). Finally, the updated assortment S_{t-1} is obtained by the correspondence $i PS_{t-1}$ if and only if Y_i° " 1.

7.2. Model Parameter Computation

For a given optimized assortment S_{t-1} from Section 7.1, we aim to search for the worst-case MNL parameter θ_{t-1} from the confidence set Ω_n that minimizes the expected revenue. In particular, we employ a gradient descent with line search (GDLS) method to compute θ_{t-1} by solving the following problem

$$\min_{\theta \in \mathbb{R}^n} V_{\mathfrak{p}} s_{t-1}; \theta \mathfrak{q} \tag{9}$$

Here, we remark that VpS_{t-1} ; θq (9) $\frac{r}{r} \frac{r_{i p s_{t-1}} r_{i} \exp px_{i}^{y} \theta q}{r_{i p s_{t-1}} r_{i p s_{t-1}} \exp px_{i}^{y} \theta q} \text{ is}$ a locally Lipschitz function in θ . Given a feasible initial parameter $\hat{\theta}^{p0q} P \Omega_n$, we run at most L gradient descent steps. Suppose β_{ℓ} is the step size for gradient descent in the ℓ -th step. At each step ℓ " 1, 2," ", L, we do a line search to maintain the feasibility. In particular, given $\theta^{p\ell}$ 1q P Ω_n , we first evaluate the gradient as ξ_{ℓ} " $\nabla_{\theta} V_p S_{t-1}$; $\theta^{\ell-1q} q$ Then we initiate β_{ℓ} with a pre-specified step size β_{ℓ} " β , and check whether $\theta^{p\ell q}$ " $\theta^{p\ell'-1q}$ " $\beta_{\ell}\xi_{\ell}$ is feasible, i.e. $\theta^{p\ell q} P \Omega_n$. If not, we set $\beta_{\ell} D C\beta_{\ell}$ for some pre-specified C P β 0, 1q, and recompute $\theta^{p^{\ell}q}$ " $\theta^{p^{\ell}-1q}$ " $\beta_{\ell}\xi_{\ell}$. Such a search is repeated until $\theta^{p^{\ell}q}$ is feasible. We provide the pseudocode in Algorithm 2 for the overall process. Note that L, B, C are all hyper-parameters. In all of our numerical studies, we set L " 2, β " 0.01 and C" $\frac{1}{2}$, which performs well empirically.

Algorithm 2 Gradient Descent with Line Search (GDLS)

Input: assortment S_{t-1} ; feasible set Ω_n ; initial parameter $\theta^{\hat{poq}}$; initial step size β ; step shrinkage constant c; number of descent steps L

```
Output: the updated parameter \theta_{t-1}
ℓĐ0
for \ell " 1 to \ell do
                                  \xi_{\ell} " \nabla_{\theta} V_p s_{t-1}; \theta^{\ell-1q} q /* compute the gradient */
                               \theta^{p\ell} a \theta^{p\ell} 1 a \theta^{p\ell} 1 a \theta^{p\ell} 8 \theta^{p\ell} 1 a \theta^{p\ell} 8 \theta^{p\ell} 8 \theta^{p\ell} 9 \theta^{p\ell} 9 \theta^{p\ell} 9 \theta^{p\ell} 1 a \theta^{p\ell} 9 \theta
                                     while \theta^{plq} R\Omega_n do
                                                                     \beta_{\ell} \oplus \mathcal{C}\beta_{\ell} /* decrease the step size */ \theta^{p\ell} \oplus \mathcal{C}^{p\ell} 1q / \beta_{\ell} \xi_{\ell}
                                     end while
end for
\theta_{t+1} \oplus \theta^{p\ell}q
```

8. Experiments

We compare the PASTA method with assortment optimization without pessimism (referred to as the baseline method in the sequel). Our method and the baseline method are evaluated on synthetic data for which the optimal assortment S° and true parameter θ are known so that the true regrets can be computed. We describe the data generation process and the baseline method in details below.

8.1. Data Generation

We consider the assortment optimization scenarios described by N, K, d, n and p, where N is the total number of available products; K is the cardinality constraint of the assortments, i.e., S" t s : |S| d'K u; d is the dimension of θ' and $t X_i \cup_{i=1}^N n$ is the sample size of the offline dataset; P is the probability for sampling the optimal assortment S^o. Similar to (Chen et al., 2020), we first generate the true preference vector θ° as a uniformly random unit d-dim vector. For $i P t1, \ldots, NL$ we generate r_i (the reward of product i) uniformly from the range r0.5, 0.8 and generate X_i (the feature of product I) as uniformly random unit d-dim vector such that $expp^{X_j} \theta$ q d'expp 0.6q to avoid degenerate cases, where the optimal assortments include too few items. Given such information, the true optimal assortment S° can be computed. Then, we generate an offline dataset D " tp S_i , A_i , R_i $q u_{i-1}^n$ with n samples. For i P $\mathfrak{t}_1, \ldots, \mathfrak{n}_s$ we generate S_i following the distribution π_S

8.2. Baseline

In our experiments, we use the gradient descent method to find $\theta_{ML,n}$ that minimizes the empirical negative log-likelihood function. Then given $\theta_{ML,n}$, the baseline method solves the assortment optimization problem by solving the linear programming problem in (7).

8.3. Performance Comparison

For a given p^N , K, d, n, p_Q we repeat the data generation process in Section 8.1 to randomly generate 50 offline datasets. The solutions of PASTA and the baseline method are recorded in these experiments. For hyper-parameters, we set α_n " $2p_{ML}$ where p_{ML} " $p_n p_{ML,n}$ q and the maximum of iteration p_{ML} " 30. We measure the performance with two metrics: (1) the *average regret* of the solutions which indicates how far the performance of the solutions is to that of the *optimal* performance (i.e., revenue of p_{ML} 0); (2) the *assortment accuracy* of the solutions (with respect to the optimal assortment p_{ML} 0). The assortment accuracy of an assortment p_{ML} 1 is defined as the ratio of the number of correctly chosen products to the number of products in p_{ML} 2. The key results are summarized below.

Effect of Sample Size. We set N " 40, K " 8, d " 16 and p " 0.9. We then gradually increase the number of samples n. The result is presented in Figure 1 indicating that PASTA significantly outperforms the baseline method. While the performance of the baseline method improves with increasing number of samples, the PASTA method maintains a regret that is less than 25% of that of the baseline method. The same experiment repeated with an increased number of products (N " 60, K " 15) demonstrates that the gain of the PASTA method is stable, as presented in Figure 1.

Effect of Probability of Sampling Optimal Assortment in Offline Data. We set N " 40, K " 8, d " 16, n " 150, and let P P 10.1, 0.3, 0.5, 0.7, 0.9We also study the effect of P in scenarios with an increased total number of products (N " 60, K " 15). As can be seen in Figure 2, the gain of pessimistic assortment optimization is consistent and robust for varying values of P.

Effect of Dimension of Features. We set N " 20, K " 5, P " 0.9, P " 150, and let P B, 20, 32, 64, 128In order to characterize the effect of dimension P d, we generate P elements of P independently from Uniformr' 1, 1s. The results are presented in Figure 3. We observe that while both the regret of the baseline method and that of the pes-

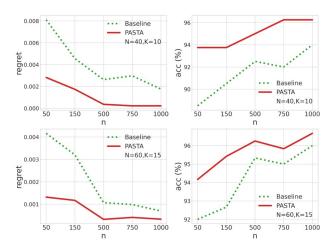


Figure 1. Performance comparison between PASTA and the baseline method with varying number of samples (n). On the left is the average regret (the lower the better) while the assortment accuracy (the higher the better) is on the right.

simistic assortment optimization increase with increasing dimensions of features, the PASTA method maintains its performance gain as the dimension *d* varies.

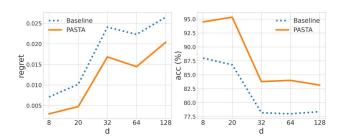


Figure 3. Comparison between PASTA and the baseline method with increasing dimensions of product features (d).

9. Conclusion

This work addresses the issue of insufficient data coverage in offline assortment optimization problems. This becomes more challenging as the number of choices grows quickly as a function of the number of items N. We presented a framework of pessimistic assortment optimization and provided theoretical justifications for our approach. We then performed an in-depth study of the Multinomial Logit Model (MNL), and derived a finite-sample regret bound of pessimistic assortment optimization for this popular model. We presented an efficient algorithm to solve the pessimistic assortment optimization problem for MNL, and demonstrated significant improvements of our approach over the baseline method by extensive numerical studies.

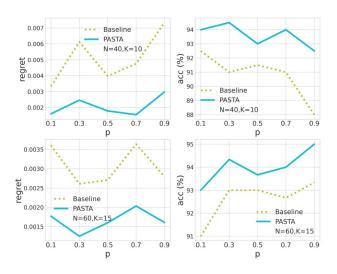


Figure 2. Comparison between PASTA and the baseline method with varying probability of the optimal assortment (P). Top row: N " 40; bottom row: N " 60.

Acknowledgements

Ethan X. Fang is partially supported by NSF DMS-2230795, NSF DMS-2230797. Cong Shi is partially supported by an Amazon Research Award.

References

- Aouad, A., Farias, V., and Levi, R. Assortment optimization under consider-then-choose choice models. *Management Science*, 67(6):3368–3386, 2021.
- Aouad, A., Feldman, J., and Segev, D. The exponomial choice model for assortment optimization: an alternative to the MNL model? *Management Science*, 2022.
- Bertsimas, D. and Kallus, N. From predictive to prescriptive analytics. *Management Science*, 66(3):1025–1044, 2020.
- Caro, F. and Gallien, J. Dynamic assortment with demand learning for seasonal consumer goods. *Management Science*, 53(2):276–292, 2007.
- Chen, X., Wang, Y., and Zhou, Y. Dynamic assortment optimization with changing contextual information. *Journal of Machine Learning Research*, 2020.
- Chen, X., Shi, C., Wang, Y., and Zhou, Y. Dynamic assortment planning under nested logit models. *Production and Operations Management*, 30(1):85–102, 2021a.
- Chen, X., Wang, Y., and Zhou, Y. Optimal policy for dynamic assortment planning under multinomial logit models. *Mathematics of Operations Research*, 46(4): 1639–1657, 2021b.

- Davis, J. M., Gallego, G., and Topaloglu, H. Assortment planning under the multinomial logit model with totally unimodular constraint structures, 2013. URL https://people.orie.cornell.edu/jmd388/publications/MNLConstr.pdf. Working Paper, Cornell University, Ithaca, NY.
- Désir, A., Goyal, V., Segev, D., and Ye, C. Constrained assortment optimization under the Markov chain–based choice model. *Management Science*, 66(2):698–721, 2020.
- Diaconis, P. and Saloff-Coste, L. Logarithmic sobolev inequalities for finite markov chains. *The Annals of Applied Probability*, 6(3):695–750, 1996.
- Feldman, J. and Topaloglu, H. Bounding optimal expected revenues for assortment optimization under mixtures of multinomial logits. *Production and Operations Manage*ment, 24(10):1598–1620, 2015.
- Feng, Q., Shanthikumar, J. G., and Xue, M. Consumer choice models and estimation: A review and extension. *Production and Operations Management*, 31(2):847–867, 2022.
- Flores, A., Berbeglia, G., and Van Hentenryck, P. Assortment optimization under the sequential multinomial logit model. *European Journal of Operational Research*, 273 (3):1052–1064, 2019.
- Fu, Z., Qi, Z., Wang, Z., Yang, Z., Xu, Y., and Kosorok, M. R. Offline reinforcement learning with instrumental variables in confounded markov decision processes, 2022. URL https://arxiv.org/abs/2209.08666. Working Paper, Northwestern University, Evanston, IL.
- Gallego, G. and Topaloglu, H. Constrained assortment optimization for the nested logit model. *Management Science*, 60(10):2583–2601, 2014.
- Harsha, P. Communication complexity, 2011. URL https: //www.tifr.res.in/~prahladh/teaching/ 2011-12/comm/lectures/112.pdf. Lecture Notes, Tata Institute of Fundamental Research (TIFR), Mumbai, India.
- Imbens, G. W. and Rubin, D. B. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, Cambridge, UK, 2015.
- Jin, Y., Yang, Z., and Wang, Z. Is pessimism provably efficient for offline RL? In Meila, M. and Zhang, T. (eds.), Proceedings of the 38th International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pp. 5084–5096. PMLR, 18– 24 Jul 2021.

- Kidambi, R., Rajeswaran, A., Netrapalli, P., and Joachims, T. Morel: Model-based offline reinforcement learning. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 21810–21823. Curran Associates, Inc., 2020.
- Kumar, A., Zhou, A., Tucker, G., and Levine, S. Conservative q-learning for offline reinforcement learning. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS'20, Red Hook, NY, USA, 2020. Curran Associates Inc.
- Li, S., Luo, Q., Huang, Z., and Shi, C. Online learning for constrained assortment optimization under Markov chain choice model, 2022. URL https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4079753. Working Paper, University of Michigan, Ann Arbor, MI.
- Liu, N., Ma, Y., and Topaloglu, H. Assortment optimization under the multinomial logit model with sequential offerings. *INFORMS Journal on Computing*, 32(3):835–853, 2020.
- McFadden, D. Conditional logit analysis of qualitative choice behavior. Institute of Urban and Regional Development, Berkeley, CA, 1973.
- McFadden, D. Econometric models of probabilistic choice. Structural analysis of discrete data with econometric applications, 198272, 1981.
- Mo, W. and Liu, Y. Efficient learning of optimal individualized treatment rules for heteroscedastic or misspecified treatment-free effect models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 84(2): 440–472, 2022.
- Owen, A. Empirical likelihood ratio confidence regions. *The Annals of Statistics*, 18(1):90–120, 1990.
- Pang, R. 1 totally unimodular matrices stanford university, 2017. URL https://theory.stanford.edu/~jvondrak/MATH233B-2017/lec3.pdf.
- Qian, M. and Murphy, S. A. Performance guarantees for individualized treatment rules. *The Annals of Statistics*, 39(2):1180–1210, 2011.
- Rosenbaum, P. R. and Rubin, D. B. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55, 1983.
- Rubin, D. B. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of educational Psychology*, 66(5):688, 1974.
- Rusmevichientong, P., Shen, M., and Shmoys, D. A PTAS

- for capacitated sum-of-ratios optimization. *Operations Research Letters*, 37:230–238, 07 2009.
- Rusmevichientong, P., Shen, Z.-J. M., and Shmoys, D. B. Dynamic assortment optimization with a multinomial logit choice model and capacity constraint. *Operations Research*, 58(6):1666–1680, 2010.
- Rusmevichientong, P., Sumida, M., and Topaloglu, H. Dynamic assortment optimization for reusable products with random usage durations. *Management Science*, 66(7): 2820–2844, 2020.
- Sauré, D. and Zeevi, A. Optimal dynamic assortment planning with demand learning. *Manufacturing & Service Operations Management*, 15(3):387–404, 2013.
- Sen, B. A gentle introduction to empirical process theory and applications, 2018. URL http://www.stat.columbia.edu/~bodhi/Talks/Emp-Proc-Lecture-Notes.pdf. Lecture Notes, Columbia University, New York, NY.
- Talluri, K. and van Ryzin, G. Revenue management under a general discrete choice model of consumer behavior. *Management Science*, 50(1):15–33, 2004.
- van der Vaart, A. W. and Wellner, J. A. Weak Convergence and Empirical Processes: with Applications to Statistics. Springer, New York, NY, 1996.
- Wang, Y., Chen, X., and Zhou, Y. Near-optimal policies for dynamic multinomial logit assortment selection models. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), Advances in Neural Information Processing Systems, volume 31. Curran Associates, Inc., 2018.
- Yu, T., Thomas, G., Yu, L., Ermon, S., Zou, J. Y., Levine, S., Finn, C., and Ma, T. Mopo: Model-based offline policy optimization. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 14129–14142. Curran Associates, Inc., 2020.
- Zhao, Y., Zeng, D., Rush, A. J., and Kosorok, M. R. Estimating individualized treatment rules using outcome weighted learning. *Journal of the American Statistical Association*, 107(499):1106–1118, 2012.

A. Proof of Theoretical Results

Throughout the proofs, we use θ^i to denote the true parameter and n to denote the number of samples. We use Φ_n to denote the empirical negative log-likelihood function, i.e., $\Phi_n \not \mapsto q^i$ " $\frac{1}{n}$ " $\frac{1}{n}$ $\frac{1$

For a general pair of random variables p^{X} , Y q assume that the conditional probability density function of Y given X is parameterically modeled by $Pp^{Y}|x$; θq for parameter θ . For technical reasons, we will consider the following distances.

Definition A.1 (Squared Hellinger Distance).

$$h^2 p^0 p^* k; \theta_1 q p^0 k; \theta_2 q q \frac{1}{2} a \frac{\dot{z}}{p_1 p v | x; \theta_1 q} a \frac{1}{p_2 p v | x; \theta_2 q} dy.$$
 (10)

Definition A.2 (Hellinger Distance).

$$h_p p_p k; \theta_1 q p_p k; \theta_2 q q^{\alpha} \frac{a}{h^2 p_p k; \theta_1 q p_p k; \theta_2 q q}$$
 (11)

Definition A.3 (Generalized Squared Hellinger Distance).

$$H^{2}p\theta_{1}, \theta_{2}q \text{ "E}_{X} h^{2}p^{0}p^{-}|X; \theta_{1}q p^{0}p^{-}|X; \theta_{2}qq.$$

$$\tag{12}$$

Definition A.4 (Generalized Hellinger Distance).

$$H_{p}\theta_{1}, \theta_{2}q \text{ "} E_{X}$$
 "a $h_{2}p^{p}p^{"}X; \theta_{1}q p^{p}"X; \theta_{2}qq$. (13)

In our theoretical results, we particularly consider PpV|x; θq " $\pi_A p^2|s$; θq as the conditional density of A given S (hereafter denoted as A|S).

A.1. Proof of Lemma 5.2

Under Assumption 5.1, for any 0 \check{a} $\check{\delta}$ \check{a} 1, with probability at least 1 $\check{\delta}$, we have for any $\theta P \Omega_n$,

$$V_{ps}^{\circ}$$
; θ_{q}° , V_{ps}° ; θ_{q} λ_{s}° , C_{s}° , $\overline{\alpha_{n}}$, (14)

where r_{s^*} a max_{$j \in S^*$} r_j is the largest possible revenue among all items in the optimal assortment.

Proof of Lemma 5.2. For any θ such that $\mathbf{p}_n \mathbf{p} \theta_{\mathbf{q}} \cdot \mathbf{p}_n \mathbf{p} \theta_{\mathbf{ML},n} \mathbf{q} d' \alpha_n$, i.e., $\theta \mathbf{P} \Omega_n$, we have

$$Vps^{\circ}; \theta \neq Vps^{\circ}; \theta \neq d^{\vee}Vps^{\circ}; \theta \neq Vps^{\circ}; \theta \neq d^{\vee}Vps^{\circ}; \theta \neq d^{\vee}Vps^{\vee}; \theta \neq d^{\vee}Vp$$

With Lemma B.2, we have that for any $\theta P \Theta$,

$$\overset{\circ}{\downarrow} V_{ps^{\circ}}; \theta_{q} \land V_{ps^{\circ}}; \theta_{ML,n} \overset{\circ}{q^{\circ}} d' r_{s^{\circ}} C_{s^{\circ}} \stackrel{}{\mathsf{E}}_{S} || \pi_{A} \mathsf{p}^{\circ}|s; \theta_{q} \land \pi_{A} \mathsf{p}^{\circ}|s; \theta_{ML,n} \mathsf{q}| \downarrow$$

where $|| "_1|$ is the ℓ_1 -norm, r_{s^*} " $\max_{j \in S^*} r_j$ is the largest possible revenue among all items and C_{s^*} " $1\{\pi_S p^{S^*} \neq 1\}$ In Lemma B.3, we establish that

"
$$E_{S} \mid [\pi_{A} p^{\circ}|S; \theta q \uparrow \pi_{A} p^{\circ}|S; \theta_{ML,n} q] \mid d^{2} \stackrel{?}{=} \frac{b}{2} \frac{1}{H^{2} p^{\theta}, \theta_{ML,n} q}$$

where H^2 is the generalized squared Hellinger distance defined in (12) with PpV|x; θq " $\pi_A p^2|s$; θq as the conditional density of A|S.

Combining the above two inequalities, we have that for any $\theta P \Theta$,

$$\stackrel{\circ}{\downarrow} V_{ps^{\circ}}; \theta_{q} \stackrel{\circ}{} V_{ps^{\circ}}; \theta_{ML,n} \stackrel{\circ}{q} \stackrel{\circ}{A} r_{s^{\circ}} C_{s^{\circ}} \stackrel{b}{\longrightarrow} \frac{H^{2}p\theta, \theta_{ML,n}}{q} \qquad (15)$$

In the following, we use the fact that $\log x$ d' $2p^{?} \overline{X}$ 1q for any $X \in 0$ to show that for any $S \cap S$ and any θ :

which implies that

$$Lp\theta q Lp\theta q e 2H^2p\theta; \theta q$$
 (16)

By Lemma B.4, we have that for any $\theta P \Omega_n$,

$$H^2p\theta$$
, $\theta_{ML,n}$ q d' $2H^2p\theta$, θ q ` $2H^2p\theta$, $\theta_{ML,n}$ q d' $Lp\theta_{ML,n}$ q ' $Lp\theta$ ' q ` $tLp\theta$ q ' $Lp\theta$ ' qu (17)

From Assumption 5.1, we have that with probability at least 1 $\dot{\delta}$, for any $\theta P \Omega_0$,

$$\stackrel{\circ}{\downarrow}$$
Lp θ q 'Lp θ 'q ' \triangleright_n p θ q ' \triangleright_n p θ 'q 'd' α_n .

In other words, under Assumption 5.1, with probability at least 1 $^{\prime}$ $^{\delta}$ for any $\theta P \Omega_n$,

Plugging Eq. (18) into Eq. (17), we have that with probability at least 1' δ , for any $\theta P \Omega_0$,

$$H^2 p\theta, \theta_{ML,n} \neq d^4 \alpha_n.$$
 (19)

Combining the above inequality and Eq. (15), we have that, with probability at least 1 δ , Vps° ; $\theta_{ML} q\dot{q} \dot{A}$ $r_{s^{\circ}} C_{s^{\circ}} \dot{\alpha}_{n}$ for all $\theta P \Omega_{n}$. This concludes the proof.

A.2. Proof of Lemma 6.1

Consider the MNL model (5) with a compact set Θ . Assume that θ° P Θ . For 0 \check{a} $\check{\delta}$ \check{a} 1, with probability at least 1 $\check{\delta}$, we have

$$\mathbf{p}_{n} \mathbf{p}^{\theta} \mathbf{q} \wedge \mathbf{p}_{n} \mathbf{p}^{\theta}_{\mathsf{ML},n} \mathbf{q} \dot{\mathbf{A}} \frac{C_{\mathsf{A}} d}{n} \log \frac{\theta_{\mathsf{max}}}{\delta}. \tag{20}$$

Proof of Lemma 6.1. Fix 0 ă δ ă 1. Suppose $\alpha_n = \frac{C_A d}{\rho} \log \frac{2\theta_{\text{max}}}{\delta}$. Define an oracle confidence set as

$$\Omega_n$$
 " $t \theta P \Theta : Lp\theta q ' Lp\theta ' q d'\alpha_n u$.

In particular, $\theta' P \Omega_0$. By Lemmas B.1 and A.2, we also have with probability at least 1' δ [2,

$$L p \theta_{ML,n} q L p \theta q d' 2 C_A H^2 p \theta_{ML,n}$$
, $\theta q \lambda \alpha_n$,

 $\sup_{f \in PF_n} |pP_n|' = \Pr_{f \in PF_n} |pP_n|'$

$$pP_{n} \stackrel{?}{\cdot} Pqp^{f} q \stackrel{\grave{A}E}{} P_{n} \stackrel{?}{\cdot} P]_{F_{n}} \stackrel{\flat}{\cdot} \frac{\#}{n} plog C_{A} \not qE P_{n} \stackrel{?}{\cdot} P]_{F_{n}} \stackrel{\iota}{\cdot} \sup_{f pF_{n}} Ep^{f} \stackrel{?}{\cdot} Ef \not q^{2} log \frac{2}{\delta} \stackrel{\iota}{\cdot} \frac{log C_{A}}{n} log \frac{2}{\delta}.$$
 (21)

For the variance term, we have

$$\sigma_{E_n}^2$$
 " $\sup_{f \in PE_n} \operatorname{Epf}' = \operatorname{Epf}'$

For the expected envelope, our goal below is to apply Sen (2018, Theorem 7.13) (stated in Theorem B.6). Consider the covering number $N \not\in F_n$, $L^2 \not\in Q$ and finitely supported probability measure Q. By Lemma A.1, based on the MNL model (5), for some L $\check{\mathbf{a}}$ '8, F_n is a class of L-Lipschitz functions with respect to the index space $p\Theta$, $2\mathbf{q}$. Then in terms of the bracketing number N_{rs} and covering number N, for any $\epsilon \mathbf{q}$ 0 and probability measure \mathbf{Q} , we have

$$N \text{ peL}, F_n, L^2 \text{pQqq d'N}_{rs} \text{p2} \epsilon \text{L}, F_n, L^2 \text{pQqq d'N pe}, \Theta_i$$

By $\Theta \check{D} R^d$ and Θ is compact, we further have $N \not = 0$, Θ , $\Omega \not = 0$. Therefore,

$$N \not\in F_n$$
, $L^2 \not\in Qqq \ A \stackrel{\frown}{\epsilon} \stackrel{d}{\epsilon}$

By Theorem B.6, we further have

$$\mathsf{E}[\mathsf{P}_n \ ' \ \mathsf{P}]_{F_n} \ \grave{\mathsf{A}} \frac{\mathsf{d}}{n} \frac{\mathsf{d}}{\sigma_{F_n}} \log \frac{\mathsf{L}}{\sigma_{F_n}} - \frac{\mathsf{d}}{n} \ \hat{\mathsf{d}} 2 \log \mathcal{C}_A \operatorname{\mathsf{q}} \log \frac{\mathsf{L}}{\sigma_{F_n}} \ \grave{\mathsf{A}} \frac{\mathsf{c}}{n} \frac{\overline{\mathcal{C}_A d}}{n} \alpha_n = \frac{C_A d}{n} \frac{\mathsf{c}}{\log \frac{2\theta_{\mathsf{max}}}{\delta}} \ \grave{\mathsf{A}} \alpha_n.$$

The Talagrand's inequality (21) becomes

$$pP_n$$
 ' $Pqpfq \lambda \alpha_n$.

In particular, in the case of $\theta_{ML,n}$ $P\Omega_n$ corresponding to $f p\theta_{ML,n}$ $q PF_n$, we have

This complete the proof.

Lemma A.1. Consider the MNL model:

$$\pi_{A} p^{j} | s; \theta q$$
 " $\frac{\exp p^{X_{j}^{J}} \theta q}{1 - i + i + i + i + i}; \quad \pi_{A} p 0 | s; \theta q$ " $\frac{1}{1 - i + i + i + i + i + i}; \quad \text{@ Ps. sps. } \theta P \Theta.$

Let θ_{\max} " $\max_{\theta \in \Theta} \|\theta\|_2$, X_{\max} " $\max_{j \in PrN_S} X_j \|_2$. If Θ is compact, that is, θ_{\max} ă `8 , then the log-likelihood ratio $\log \frac{\pi_A p^A \|S; \theta^*\|_2}{\pi_A p^A \|S; \theta^*\|_2}$ is a uniformly Lipschitz function in $\theta \in \Theta$.

Proof of Lemma A.1.

$$\begin{cases}
\frac{B}{B\theta} \log \frac{\pi_A pA | S; \theta q}{\pi_A pA | S; \theta q}, & \frac{1}{\pi_A pA | S; \theta q}, & \frac{B}{B\theta}, & \frac{B}{B\theta},$$

That is, $\theta \triangleright \tilde{N}\log \frac{\pi_A p^A | S; \theta^* q}{\pi_A p^A | S; \theta q}$ is L-Lipschitz.

Lemma A.2 (Concentration of Parametric MLE in Hellinger Distance). Consider the MNL model (5). For $0 \ \mbox{\'a} \ \mbox{\'a} \ \mbox{1}$, with probability at least $1' \ \mbox{\'o}$, we have

$$H^2 p \theta_{ML,n}$$
, $\theta \neq \lambda \frac{d}{n} \log \frac{\theta_{max}}{\delta}$.

Proof of Lemma A.2. We follow from Fu et al. (2022, Corollary 2) as a special case, where our data are generated i.i.d. instead of being a general Markov chain.

A.3. Proof of Lemma 6.2

Consider the MNL model (5). Suppose conditions in Lemma 6.1 hold, and $L_p\theta_q$ and $L_np\theta_q$ are uniformly and strongly convex. Let $\alpha_n = \frac{C_A d}{n} \log \frac{\theta_{max}}{\delta}$. For $0 \ \Bar{a} \ \Bar{a} \ \Bar{a}$, with probability at least $1' \ \Bar{\delta}$, we have

$$\sup_{\theta \in \Omega_n} \stackrel{\checkmark}{\stackrel{L}{\not}} p\theta q \stackrel{\prime}{/} L p\theta \stackrel{\circ}{q} \stackrel{\prime}{/} p \stackrel{\prime}{q} \stackrel{\prime}{/} p \stackrel{\prime}{q} \stackrel{\prime}{/} q \stackrel{\prime}{/} \alpha_n.$$

Proof of Lemma 6.2. Fix $0 \, \check{a} \, \check{b} \, \check{a} \, 1$. By the strong convexity assumption on $L \, p\theta q$ and $L \, n \, p\theta q$ there exists a constant $\mu \, \check{a} \, 0$ such that for any $\theta \, P \, \Theta$,

$$\{\mu\}$$
 $\{\theta' \in \theta'\}$ $\{\theta'\}$ $\{\theta'\}$ $\{\theta'\}$ $\{\theta'\}$ $\{\phi'\}$ $\{\phi'\}$

By Lemma 6.1, with probability at least 1 δ [2, we have $\theta_{ML,n}$ $P \Omega_n$. Then for any $\theta P \Omega_n$, we have

$$\begin{vmatrix} \theta & \theta \\ 1 & b \end{vmatrix}^2 \frac{d' \mid \theta \wedge \Phi_{ML,n} \mid 2 \mid }{|\Phi_n| p \theta q \wedge |\Phi_n| p \Phi_{ML,n} \mid q} \begin{vmatrix} \Phi_{ML,n} \mid b \\ \frac{1}{\mu} \end{vmatrix} \frac{d' \mid 2}{|\Phi_{ML,n}| q \wedge |\Phi| p \Phi_{ML,n} \mid q}$$
 pby triangular inequality q pby strong convexity q pby $\theta \mid \Omega_n$ and $\Phi_{ML,n} \mid P \mid \Omega_n$ respectively q

The above implies that with probability at least 1' $\delta[2]$, we have $\Omega_n \stackrel{\bullet}{D} \Omega_n$, where Ω_n is a ball centered around θ' with radius $\overline{\alpha}_n$:

$$\mathfrak{Q}_n$$
 " $\neg \theta P\Theta : \{\theta', \theta'\}_2 d'^{?} \overline{\alpha}_n^{()}$.

Define \mathcal{B}_n " $\log \frac{\pi_A \, p^A \, |S;\theta^* \, q}{\pi_A \, p^A \, |S;\theta\, q} : \theta \, P \, \mathfrak{D}_n$. Let $\{P_n \, ' \, P\}_{\mathcal{B}_n}$ " $\sup_{f \in \mathcal{B}_n} |pP_n|' \, Pqpf \, q$ be the envelop. By Talagrand's inequality, with probability at least 1 ' $\delta[2]$, we have for any $f \, P \, \mathcal{B}_n$,

For the variance term, we have

$$\sigma_{R_n}^2$$
 " $\sup_{f \in R_n} E f$ ' $E f q^2 d$ ' $\sup_{\theta \in \Omega_n} E \left[\log \frac{\pi_A p^A | S; \theta q}{\pi_A p^A | S; \theta q} \right]^2$

A $\sup_{\theta \in \Omega_n} \left\{ \theta \right\}_2^2$ pby Lipschitzness in Lemma A.1q d' α_n pby definition of $\Omega_n q$

For the expected envelope, by Theorem B.6, we further have

$$\mathsf{E} \mathsf{P}_n \, \dot{} \, \, \mathsf{P} \mathsf{P}_{\mathsf{R}_n} \, \, \dot{\mathsf{A}} \, \, \frac{\mathsf{d}}{n} \frac{\mathsf{d}}{n} \sigma_{\mathsf{R}_n}^2 \log \frac{\mathsf{L}}{\sigma_{\mathsf{R}_n}} - \frac{\mathsf{d}}{n} \, \, \, 2 \log \mathcal{C}_{\mathsf{A}} \operatorname{qlog} \frac{\mathsf{L}}{\sigma_{\mathsf{R}_n}} \, \, \, \, \dot{\mathsf{A}} \, \, \, \frac{\mathsf{d}}{n} \alpha_n \, \, \dot{\mathsf{A}} \, \, \alpha_n.$$

Therefore, the Talagrand's inequality (23) gives that for any $f P \not \vdash B_n$

$$|pP_n|'$$
 Pqp f q $|$ $\lambda \alpha_n$.

In other words, with probability at least 1 ' δ , $\Omega_n \stackrel{\bullet}{D}$ Ω_n corresponding to $t^f p \theta_{Q} u_{P\Omega_n} \stackrel{\bullet}{D}$ Ω_n , we have

$$\sup_{\theta \in \Omega_n} \hat{L} p\theta q L p p p q L p p q L p p p q L p p p q L p p q$$

B. Technical Lemmas

$$|Lp\theta q \cdot Lp\theta^{\circ} q| d^{2}C_{A}H^{2}p\theta, \theta q$$
 (24)

Proof of Lemma B.1. By definition,

$$|L p\theta'| q |L p\theta q|$$
 "E $|E_A| \log \frac{\pi_A p^A |S; \theta'| q}{\pi_A p^A |S; \theta q}$ "S

In particular, for a fixed SP Ns, we have

E"
$$\log \frac{\pi_A p^A | S; \theta \neq S}{\pi_A p^A | S; \theta \neq S}$$
" S" KL $\pi_A p^B | S; \theta \neq T_A p^B | S; \theta \neq$

pby log-Sobolev inequality (Diaconis & Saloff-Coste, 1996, Theorem A.1) q " $\frac{C_A \log p C_A}{C_A}$ 2 $\frac{1}{2}$ $\frac{1}$

Therefore,

$$|Lp\theta^{\circ} q \wedge Lp\theta q| dE^{-2}C_A h^2 p \pi_A p^{\circ} |S; \theta^{\circ} q \pi_A p^{\circ} |S; \theta qq^{\circ} = 2C_A H^2 p \theta^{\circ}, \theta q$$

Lemma B.2. Let C_{s^*} " $\frac{1}{\pi_s ps^* q}$ and r_{s^*} " $\max_{j ps^*} r_j$, then the following inequality holds for any θ_1 , $\theta_2 p \Theta$: V_{ps} ; $\theta_{1}q \cdot V_{ps}$; $\theta_{2}q \cdot r_{s} \cdot C_{s} \cdot E_{s} \cdot |\pi_{A}p \cdot S; \theta_{1}q \cdot \pi_{A}p \cdot S; \theta_{2}q|$

where || " || denotes the L 1 norm.

Proof of Lemma B.2.

$$\stackrel{\circ}{\downarrow}Vps^{\circ}; \theta_{1}q^{\prime}Vps^{\circ}; \theta_{2}q^{\circ} \stackrel{\circ}{\stackrel{\circ}{\downarrow}}Es \stackrel{\circ}{\stackrel{\circ}{\downarrow}Es} \frac{\mathsf{IpS} \stackrel{\circ}{\stackrel{\circ}{\downarrow}}s^{\circ}q^{\vee}}{\pi_{S}pSq} \stackrel{\circ}{\iota_{PS}}r_{i,S} \stackrel{\circ}{\stackrel{\circ}{\downarrow}}\pi_{A}p^{i}|S; \theta_{1}q^{\prime}\pi_{A}p^{i}|S; \theta_{2}q \stackrel{\circ}{\stackrel{\circ}{\downarrow}}$$
(25)

$$d' E_{S} \stackrel{\text{(i)}}{\stackrel{\smile}{\downarrow}} \frac{r_{i,S}}{\pi_{S} p S_{q}} \stackrel{\text{(i)}}{\stackrel{\smile}{\downarrow}} r_{i,S} \stackrel{\text{(i)}}{\pi_{A}} p^{i} |S; \theta_{1} q \wedge \pi_{A} p^{i} |S; \theta_{2} q \stackrel{\text{(i)}}{\stackrel{\smile}{\downarrow}}$$

$$(26)$$

$$d' E_{S} \stackrel{\text{(Y)}}{\stackrel{\text{(Y)}}{\longrightarrow}} \frac{g}{\pi_{S} p S q} \stackrel{\text{(Y)}}{\stackrel{\text{(Y)}}{\longrightarrow}} r_{i,S} \stackrel{\text{(Y)}}{\longrightarrow} r_{i,S} \stackrel{\text{(Y)}$$

$$d' C_{s^{\circ}} E_{S} IpS "s^{\circ} q \overset{\ddot{y}}{\overset{\checkmark}{\downarrow}}_{i,S} \pi_{A} p^{i} |S; \theta_{1} q ' \pi_{A} p^{i} |S; \theta_{2} q \overset{\checkmark}{\overset{\checkmark}{\downarrow}}$$

$$(28)$$

"
$$r_s$$
 C_s E_s $| \Pi_A p$ S ; $\theta_1 q$ $\Pi_A p$ S ; $\theta_2 q$ $| \cdot \rangle$, (30)

where Eq. (25) comes from the sample-based estimation of ErRpsqs(Eq. (1)), Eq. (27) comes from the Hölder's inequality, Eq. (28) comes from the fact that $\overset{\cdots}{\underset{\pi_{S} pSq}{\vdash}} \overset{\cdots}{\underset{\pi_{S}}{\vdash}} \overset{\cdots}{\underset{\pi_{S} pSq}{\vdash}} \overset{\cdots}{\underset{\pi_{S}}{\vdash}} C_{S^{\circ}}$ because $\frac{IpS^{\circ} \circ q}{\pi_{S} pSq}$ has the value zero everywhere except at $S^{\circ} \circ S^{\circ}$. The last equality follows from the definition of L^1 norm.

Lemma B.3. For any θ_1 , $\theta_2 P \Theta$,

"
$$E_{S} \mid \Pi_{A} p^{"}|S; \theta_{I} q ' \Pi_{A} p^{"}|S; \theta_{2} q \mid \downarrow d' 2 \stackrel{?}{=} \frac{1}{2} \frac{1}{H^{2} p \theta_{1}, \theta_{2} q}$$

Proof of Lemma B.3. We first use the facts that (1) L^{1} " $\frac{1}{2}$ TV where TV is the total variation distance and (2) TV d' $\overline{2}h$ (Harsha, 2011) where h is the Hellinger distance to have that for any SPS,

$$| \mathcal{T}_{A} p^{"} \sharp; \theta_{i} q^{\prime} \mathcal{T}_{A} p^{"} \sharp; \theta_{i} q | \downarrow d^{\prime} 2^{?} \overline{2} h^{\prime} \mathcal{T}_{A} p^{"} \sharp; \theta_{i} q^{\prime} \mathcal{T}_{A} p^{"} \sharp; \theta_{i} q^{\prime}$$

$$(31)$$

From (31), we have

Taking expectation with respect to S on both sides, we have

"
$$Es || \mathcal{T}_A p^* |S; \theta_1 q \wedge \mathcal{T}_A p^* |S; \theta_2 q|_1^2 \quad d' \quad 8Es \quad h^2 p \mathcal{T}_A p^* |S; \theta_1 q \mathcal{T}_A p^* |S; \theta_2 q q,$$
"
$$Es || \mathcal{T}_A p^* |S; \theta_1 q \wedge \mathcal{T}_A p^* |S; \theta_2 q|_1^2 \quad d' \quad 8H^2 p \theta_1, \theta_2 q$$

By the Jensen's inequality, we have

This implies that

"
$$E_{S} \mid | \pi_{A} p^{"} | S; \theta_{I} q \mid \pi_{A} p^{"} | S; \theta_{2} q | \downarrow d' 2 \quad \frac{2}{2} \frac{a}{H^{2} p \theta_{1}, \theta_{2} q}$$

Lemma B.4 (Properties of H and H^2). For any θ_1 , θ_2 , $\theta_3 P\Theta$, the following inequalities hold:

$$H_{p}\theta_{1}, \theta_{2}q \quad d' \quad H_{p}\theta_{1}, \theta_{3}q \quad H_{p}\theta_{2}, \rho_{3}q$$
 $H_{p}\theta_{1}, \theta_{2}q^{2} \quad d' \quad H^{2}p\theta_{1}, \theta_{2}q \quad d'H_{p}\theta_{1}, \theta_{2}q$
 $H^{2}p\theta_{1}, \theta_{2}q \quad d' \quad 2H^{2}p\theta_{1}, \theta_{3}q \quad 2H^{2}p\theta_{2}, \theta_{3}q$

Proof of Lemma B.4. For ease of notation, for i " 1, 2, 3 we use p_i to denote π_A parametrized by θ_i , i.e., $p_i p_i p_i | sq$ " $\pi_A p_i | sq$ "

(1) Notice that for any S P S, a $h^{2}p^{0}_{1}p^{-}$ ^{B}q $^{B}p^{0}_{2}p^{-}$ ^{B}q is just the regular Hellinger distance that satisfies the triangular inequality. Hence we have

$$a \frac{1}{h^2 p \rho_1 p^* sq \rho_2 p^* sqq} a \frac{1}{h^2 p \rho_1 p^* sq \rho_3 p^* sqq} a \frac{1}{h^2 p \rho_2 p^* sq \rho_3 p^* sqq}$$
(32)

Take expectation of both side of Eq. (32) with respect to S, we have

in the off Eq. (32) with respect to 5, we have
$$\stackrel{\text{"a}}{\text{Es}} \frac{1}{h^2 p_1 p^{\text{"}} \text{Sq} p_2 p^{\text{"}} \text{Sqq}} \stackrel{\text{I}}{\text{qf}} \stackrel{\text{"a}}{\text{Es}} \frac{1}{h^2 p_2 p^{\text{"}} \text{Sq}} \stackrel{\text{I}}{\text{qf}} \stackrel{\text{"a}}{\text{pf}} \stackrel{\text{I}}{\text{Sqq}} \stackrel{\text{I}}{\text{sq}} \stackrel{\text{I}}{\text$$

By the definition of H, this means that

$$H_p\theta_1$$
, θ_2q $dH_p\theta_1$, θ_3q $H_p\theta_2$, θ_3q

(2) For the first inequality, by applying the Jensen's inequality, we have

Then the inequality follows.

For the second inequality, we have

$$H p\theta_1, \theta_2 q 'H^2 p\theta_1, \theta_2 q "Es hp\theta_1 p" Sq p_2 p" Sqq 'h^2 p\theta_1 p" Sq p_2 p" Sqq$$

$$(34)$$

" E_S 1'
$$h_p p_1 p^p p_2 p^p p_3 q h_p p_1 p^p p_4 p^p p_4 p^2 p^p p_4 q.$$
 (35)

Note that for any S , 1' $^{h}p^{0}_{1}p^{-}$ ^{S}q $^{0}p^{-}$ ^{S}q is a non-negative function because the Hellinger distance is no larger than 1, and $^{h}p^{0}_{1}p^{-}$ ^{S}q $^{0}p^{-}$ $^{0}p^{-}$

$$E_S$$
 1 ' $h_p p_1 p^- Sq p_2 p^- Sqq h_p p_1 p^- Sq p_2 p^- Sqq e 0.$

Then the inequality follows.

(3) Notice that for any a, b, Çwe have pa' bq² d' 2pa' cq² \ 2pb' cq². With this fact, we have that for any S,

$$d' \frac{1}{2} \frac{\dot{z}}{2} = a \frac{1}{p_1 p a |sq} + a \frac{1}{p_3 p a |sq$$

"
$$2^{-1}\frac{1}{2}$$
 ' $a \frac{1}{p_1p^2|sq}$ ' $a \frac{1}{p_3p^3|sq}$ ' $a \frac{1}{p_3p^3|sq}$ ' $a \frac{1}{p_2p^3|sq}$ ' $a \frac{1}{p_3p^3|sq}$ '

This implies that

$$H^2p\theta_1$$
, θ_2q $d^2H^2p\theta_1$, θ_3q $2H^2p\theta_2$, θ_3q (40)

Lemma B.5 (Log-Density Ratio Variance Bound). Suppose X, P is an R-valued random variable with probability density function P, and P_1 , P_2 are two other probability density functions for X such that P_1 and P_2 are uniformly bounded from below by C^{-1} on the support of P. Then we have

$$E_{X_{,,p}} \log \frac{p_1 p^{X} q}{p_2 p^{X} q}^{2} d' 2Ch^2 p^{0}_{1}, p_2 q$$

where h^2 is the squared Hellinger distance in (10).

Proof of Lemma B.5. By $logp^X q$ d' $2p^{?} X'$ 1q for any $X \in Q$, we have

Theorem B.6 (Sen (2018, Theorem 7.13)). Let F be a measurable function class, such that $\sup_{f \in F} f_{8}$ d' f_{\max} for some constant f_{\max} g_{∞} a '8 . Assume that for f_{∞} e f_{\max} , f_{∞} d' f_{∞} , and every finitely supported probability measure f_{∞} , we have the covering number (Sen, 2018) as:

$$N \not \in F, L^2 \not \in Qqq \ A \stackrel{\circ}{=} \stackrel{d}{=} .$$
 (41)

Let σ_F^2 " $\sup_{f \in F} \mathsf{E} p^f$ ' $\mathsf{E} f q^2$. Then we have

$$\mathsf{E}]\mathsf{P}_n \ ' \ \mathsf{P}]_F \ \grave{\mathsf{A}} \ ^{\mathbf{C}} \ \frac{\overline{d}}{\overline{n}} \sigma_F^2 \ \log \frac{A}{\sigma_F} \ _ \ ^{\mathbf{''}} \ \frac{d}{\overline{n}} f_{\max} \ \log \frac{A}{\sigma_F} \ .$$

Proof of Theorem B.6. In this proof, we denote X as the underlying random variable, $tX_iU_{i=1}^n$ are n i.i.d. copies of X, and for any $f \ PF$, $P_n \not f \ q^u \frac{1}{n} \int_{i=1}^n f pX_i q$, $P_p f \ q^u E_r f pX$ qs Without loss of generality, assume that $0 \ PF$, and for any $f \ PF$, $P_p f \ q^u 0$. Let $t \in U_{i=1}^n$ be i.i.d. Rademacher random variables that are independent of $tX_iU_{i=1}^n$. By symmetrization, we have

$$E]P_{n} \stackrel{\wedge}{\cdot} P]_{F} \stackrel{d'}{d'} 2E \sup_{f \in F} \stackrel{\circ}{\underset{j=1}{\overset{\circ}{\vee}}} \frac{1}{n} \stackrel{y^{n}}{\underset{j=1}{\overset{\circ}{\vee}}} \epsilon_{i} f p X_{i} \stackrel{\circ}{\mathbf{q}}.$$

$$(42)$$

Conditional on $t^{X_i} u^n_{i=1}$, by Dudley's entropy bound, we have

$$E_{\epsilon} \sup_{f \in F} \bigvee_{j=1}^{s} \epsilon_{i} \frac{f \operatorname{pX}_{i} \operatorname{q}}{n} \bigvee_{0}^{s} \operatorname{d}^{i} \frac{\dot{z} \sigma_{F,n}}{n} \quad \text{a} \quad \frac{1}{\log N \operatorname{pU}_{i} F_{i} L^{2} \operatorname{pP}_{n} \operatorname{q}} \operatorname{qd} u, \tag{43}$$

where we consider the L^2pP_nq as the metric on F, that is, for any f PF, $f F_{L^2pP_nq} = \frac{1}{n} \int_{-i^n-1}^{i^n} f p^{i^n} d^n$. We also denote $\sigma_{F,n}^2 = \sup_{f \in F} f F_{L^2pP_nq}$ and E_{ϵ} to emphasize that the expectation is taken with respect to the Rademacher random variables $f \in U_{n-1}^{i^n}$ but holding $f \times_i U_{n-1}^{i^n}$ as fixed. By (41) with Q chosen as P_n , we have

(43)
$$\grave{A}$$
 $\stackrel{\grave{z}}{a} \stackrel{\sigma_{F,n}}{a} \stackrel{c}{\log} \frac{A}{\delta} d\delta$

$$\stackrel{{}^{\diamond}}{d} \stackrel{\cdots}{2} \frac{A}{\partial \sigma_{F,n}} \log \frac{A}{\sigma_{F,n}} \qquad \text{pby Lemma B.7q}$$
(44)

In particular, $logp^A \{ \sigma_{F,n} \neq elogp^A \} \{ f_{max} \neq elogp^A$

where E takes expectation with respect to $t^{X_i}u^n_{i^*}$ ₁. Notice that

where we define F^2 " t^2 : f PF u. We aim to apply (42), (43), (44), (45) to F^2 . Notice that $\sup_{f^2 \in P_{n,q}} \{f^2\}_{L^2 pP_{n,q}} \{g^2\}_{L^2 pP_{n,q}} \{g^2\}_{L^2$

$$N p2f_{\text{max}} \epsilon, F^2, L^2 pP_n qq dN p\epsilon, F, L^2 pP_n qq \grave{A} \frac{\hat{A}^{d}}{\epsilon}$$

Therefore, applying (42), (43), (44) and (45) to F^2 , we have

E]
$$P_n ' P$$
 $F_2 \grave{A} = \frac{c}{n} + \frac{d}{n} + \frac{1}{2} \frac{f_{max}^2 Ep\sigma_{F,n}^2}{f_{max}^2 Ep\sigma_{F,n}^2} = \frac{4A^2}{Ep\sigma_{F,n}^2}$

Define B " b $\frac{1}{2} Ep\sigma_{F,n}^2 \text{ qlog } \frac{A^2}{Ep\sigma_{F,n}^2}$. Then we have

$$\mathsf{Ep} \sigma_{F,n}^2 \ \mathsf{q} \ ' \ \sigma_F^2 \ \mathsf{A} \ \frac{\mathsf{d}}{n} f_{\mathsf{max}} B.$$

By $u \, \tilde{\mathbf{N}} \, u \, \log \frac{A^2}{u}$ is non-decreasing on $u \, \mathbf{P} \, \hat{\mathbf{p}}$, $A^2 \{ e \, \text{s} \, \text{and non-increasing on} \, u \, \mathbf{P} \, A^2 \{ e, \, \mathbf{e} \, \text{s} \, \text{o} \, \text{o$

$$B^2 \text{ "} \frac{1}{2} \mathsf{E} p \sigma_{F,n}^2 \operatorname{qlog} \frac{A^2}{\mathsf{E} p \sigma_{F,n}^2 \operatorname{q}} \grave{\mathsf{A}} \frac{1}{2} \overset{\#^*}{\sigma_F^2} \overset{\mathsf{C}}{\overset{\mathsf{C}}{\sigma_F^2}} \overset{\mathsf{C}}{\overset{\mathsf{d}}{n}} f_{\max} B \overset{\mathsf{A}}{\overset{\mathsf{A}}{\circ}} \overset{\mathsf{A}^2}{\overset{\mathsf{e}}{\bullet}} \operatorname{log} \overset{\mathsf{A}^2}{\overset{\mathsf{D}}{\sigma_F^2}} \overset{\mathsf{A}^2}{\overset{\mathsf{d}}{\sigma_F}} f_{\max} B \overset{\mathsf{A}^2}{\overset{\mathsf{A}^2}{\overset{\mathsf{e}}{\bullet}}}.$$

In particular, $B \stackrel{b}{d} = \frac{1}{2e}$, $\sigma_F^2 \stackrel{d}{d} f_{\text{max}}^2 \stackrel{d}{d} A^2 \{ e^2 \times A^2 \{ e \text{ Then the cap } A^2 \{ e \text{ is inactive as } d [n \ \tilde{N} \ 0 \text{ asymptotically.} \}$ Therefore,

$$B^2 \grave{A} \quad \sigma_F^2 \stackrel{c}{\stackrel{\cdot}{\cdot}} \frac{\overline{d}}{n} f_{\text{max}} B \stackrel{s}{\stackrel{\cdot}{\cdot}} \log \frac{A}{\sigma_E}$$

In particular, B is bounded by both roots of the corresponding quadratic equation:

$$B \grave{A} \frac{1}{2} \frac{\overset{\$}{\&} c}{\overset{\overline{d}}{n}} f_{\text{max}} \log \frac{A}{\sigma_{F}} \stackrel{\text{d}}{\overset{\text{d}}{n}} f_{\text{max}} \log \frac{A}{\sigma_{F}} \stackrel{\text{d}}{\overset{\text{d}}{n}} f_{\text{max}} \log \frac{A}{\sigma_{F}} \stackrel{\text{d}}{\overset{\text{d}}{\sigma_{F}}} \log \frac{A}{\sigma_{F}} \log \frac{A}{$$

Combined with (45), we further have

$$\mathsf{E}]\mathsf{P}_n \ ' \ \mathsf{P}_{F} \ \grave{\mathsf{A}} \ \overset{\mathsf{C}}{\overset{\mathsf{C}}{\underline{d}}} \overset{\mathsf{C}}{\underline{d}} B \ \grave{\mathsf{A}} \ \overset{\mathsf{C}}{\overset{\mathsf{d}}{\underline{d}}} \sigma_F^2 \log \frac{A}{\sigma_F} - \ \overset{\mathsf{d}}{\underline{d}} f_{\max} \log \frac{A}{\sigma_F} \ .$$

Lemma B.7. Suppose a, A a 0 such that logpA aq e 1. Then we have

$$\dot{z}_a C \underline{\qquad} C \underline{\qquad} C \underline{\qquad} \log \frac{A}{u} du d' 2a \log \frac{A}{a}.$$

Proof of Lemma B.7. Define

$$f_{paq}" = \begin{cases} # & b & \overline{\log \frac{A}{a}}, & \S_a \\ 0 & \log \frac{A}{u} \end{cases}$$

$$\begin{cases} \S_a \\ \log \frac{A}{u} \end{cases} du, \quad a \neq 0; \\ a \text{ } 0. \end{cases}$$

Then f is continuous at 0. Moreover, for $a \neq 0$, we have

$$f \operatorname{paq} = \operatorname{log} \frac{A}{a} + \operatorname{log} \frac{A}{a}$$

which is nonnegative if $logpA[aq \ \check{e}\ 1.\ As\ a\ \tilde{N}\ 0]$, we further have

Therefore, for any $\overset{a}{\bullet} \overset{e}{\bullet} 0$, we have $\overset{f}{\circ} p\overset{a}{\bullet} 0$, which concludes the proof.