

# Self-fulfilling Bandits: Dynamic Selection in Algorithmic Decision-making

Jin Li, Ye Luo, Xiaowei Zhang\*

Faculty of Business and Economics, The University of Hong Kong, Pokfulam Road, Hong Kong SAR, jli1@hku.hk, kurtluo@hku.hk, xiaoweiz@hku.hk

This paper identifies and addresses dynamic selection problems that arise in online learning algorithms with endogenous data. In a contextual multi-armed bandit model, we show that a novel bias (*self-fulfilling bias*) arises because the endogeneity of the data influences the choices of decisions, affecting the distribution of future data to be collected and analyzed. We propose a class of algorithms to correct for the bias by incorporating instrumental variables into leading online learning algorithms. These algorithms lead to the true parameter values and meanwhile attain low (logarithmic-like) regret levels. We further prove a central limit theorem for statistical inference of the parameters of interest. To establish the theoretical properties, we develop a general technique that untangles the interdependence between data and actions.

*Key words:* Self-fulfilling bias, dynamic selection, endogeneity spillover, contextual multi-armed bandit model

---

\* We are grateful for helpful conversations with Chunrong Ai, Xi Chen, Pingyang Gao, Bob Gibbons, John Klopfer, Danielle Li, Wing Suen, Chang Sun, Brian Viard, and seminar participants at University of Hong Kong and Chinese University of Hong Kong (Shenzhen). All remaining errors are ours.

## 1. Introduction

Due to better data availability and improved computing power in recent years, organizations are increasingly using algorithms to make decisions that depend on the characteristics of individual customers. For example, Netflix recommends movies based on the past browsing history of its users, Amazon suggests products based on individual shopping history, and Facebook displays advertisements that are customer-dependent. Unlike the widely used A/B testing<sup>1</sup>, these algorithms do not wait for the final outcome of the experiments before choosing their actions. Instead, they engage in *online* learning so that their choices of actions are constantly adjusted in response to the new information and feedback from the behaviors of the targeted users.

In addition to e-commerce recommendations, there are many business applications for these online learning algorithms (Blattner, Nelson, and Spiess 2021). Banks can use past data to decide whether to approve a loan for applicants of particular characteristics, and if approved, the information on how the applicant repays the loan can be used to update their decision rules in the future. Schools and companies can use existing data to decide whether to interview a candidate of a certain background, and the outcome of the interview can be used to update their future interview criterion. As more organizations are going through digital transformation, these online learning algorithms will become more important. The competitiveness of the organizations will also increasingly depend on the effectiveness of these algorithms because they influence key business policies such as hiring and customer recommendation.

The existing studies on online learning algorithms, carried out mostly by computer scientists (Chu et al. 2011, Agrawal and Goyal 2013, Slivkins 2019), have typically assumed that the model for the data generating process accurately reflects the problem to be studied. In particular, they have not paid attention to the endogeneity problem in data, such as the omitted variables, sample selection, and model mis-specification and etc. The endogeneity problems are, of course, essential for analyzing human behaviors and are prevalent in the applications of these online learning algorithms. The endogeneity issue also introduces a new type of *dynamic selection* problem in online learning algorithms because the data analysis and decision-making interact. Essentially, the endogeneity in the data will influence the decisions—whether to approve a loan application or to interview a candidate—and since the counterfactual outcomes of those who are turned away are never observed, the decisions lead to a selected sample to be analyzed, causing further bias in the data analysis.

<sup>1</sup> The use of A/B testing or randomized controlled experiments in social sciences has received substantial interests in recent years (Blake, Nosko, and Tadelis 2015, Athey and Imbens 2017, Koning, Hasan, and Chatterji 2019, Kohavi, Tang, and Xu 2020, Azevedo et al. 2020.)

The purpose of this paper is to identify and address the dynamic selection problems that arise in online learning algorithms with data endogeneity. Our setup is a contextual multi-armed bandit problem. An agent (e.g., a website) in each period first observes a vector of covariates (e.g., visitor characteristics such as age, gender, browsing history, etc.), and then decides on which arm to select (e.g., banners to display). The decision gives the agent a random reward in each period. The agent does not know the expected rewards of pulling the arms. In deciding which arm to pull, he balances exploitation (pulling the arm with a higher expected reward) with exploration (pulling the arm that generates information that allows him to make better decisions in the future). This bandit problem is contextual in nature because the agent’s expected reward from pulling an arm depends on the vector of covariates: the website’s choice of which banner to display depends on the characteristics of the visitor.

Unlike typical contextual bandit problems, we assume that there exists endogeneity in the data, which may arise because of a number of reasons, such as measurement error, model mis-specification, sample selection, omitted variables, and so on. For example, suppose that when the consumer sentiment is high, the visitors to a website are more likely to click on the banners. Suppose also that the website is more likely to attract older customers when the consumer sentiment is high. In this case, the age of the customers (covariate) becomes correlated with the clicks (reward) through an omitted variable (consumer sentiment), thereby leading to an endogeneity problem.

Our first contribution is to identify a novel type of bias that arises in this situation. We refer to this bias as *self-fulfilling bias* because actions (choices of arms) and beliefs regarding the rewards of the arms (that arise from the data analysis) are entangled. Wrong beliefs regarding the rewards can lead to wrong actions, and these actions can then generate data that actually vindicate the beliefs. In online learning environments, the actions are not purely random as the decision policy adapts to past data. When the noise is correlated with the covariates, the noise also becomes correlated with the actions, causing them to become endogeneous. We refer to this as the *endogeneity spillover*, which is the key reason underlying the self-fulfilling bias.

Our second contribution is to propose algorithms that not only correct for the self-fulfilling bias but also generate actions that obtain low levels of regret (loss in payoff compared to optimal actions). Our method incorporates the instrumental variable (IV) approach to leading online learning algorithms for contextual bandit problems—that is, algorithms of Greedy or upper-confidence-bound (UCB) types. Our proposed algorithms consist of three phases. In Phase 1, arms are randomly selected. At the end of Phase 1, the two-stage least squares (2SLS) method is implemented on the data associated with each arm to compute the arm-specific coefficient estimates. In Phase 2, the

arm with the highest expected reward (using the coefficient estimates in Phase 1) is selected in each period. The covariance matrix of the estimates are calculated at the end of Phase 2. In Phase 3, arms are selected according to the particular algorithm used—Greedy or UCB—and the coefficients and covariance matrices are updated using “joint-2SLS” in each period. Under these algorithms, the coefficients converge to the true value in the long run. Depending on whether the Greedy or UCB algorithm is used, the regret of our algorithm is of order  $\mathcal{O}(\log(T))$  and  $\mathcal{O}(\log^2(T))$ , respectively.

Our third contribution is to develop a technique that facilitates the theoretical analysis of online learning algorithms with dynamic selection problems. As mentioned above, a key feature of online learning algorithms is that data and actions are intertwined. This complicates the analysis because the effects of past actions can be long-lasting. In particular, when a poor estimate of the coefficients of the covariates causes the algorithm to select a wrong arm, it contaminates the generated data by changing the distribution of actions that are endogenous in the future (because the data associated with the optimal arm is no longer observed). The contaminated data could then worsen the future coefficient estimates that are calculated based on it. The algorithm is potentially trapped in a vicious cycle—the coefficient estimates are inconsistent and the policy is sub-optimal even in the long run.

Mathematically, in our algorithms and in online learning algorithms in general, this complication manifests in a two-way dependence between (i) the estimates of the coefficients and (ii) the estimates of the covariance matrix between the IVs and an augmented vector—that is induced by the actions—of the covariates. On the one hand, in 2SLS, the coefficient estimates depend on the covariance estimates. On the other hand, each covariance estimate depends on *all* the past coefficient estimates because they affect the past actions and the past realization of data. This latter dependence is new and central to the analysis of online learning algorithms with selection problems. Moreover, the dependence is complicated both because the coefficient estimates are serially correlated and because they affect the estimates of the covariance matrix in a nonlinear manner.

Our technique deals with this two-way dependence by performing an induction in a zig-zag manner. In each induction step, we first derive the properties of the covariance estimates (zig), and then use this property to establish the induction hypothesis on coefficient estimates (zag). The heart of the induction is the zig-part, which employs a special technique to eliminate the dependence of the covariate estimates on all past coefficient estimates. The theoretical analysis of this dependence also indicates that, for the coefficient estimates to converge to the true value, the covariance estimates at the beginning of Phase 3 must be sufficiently close to the covariance matrix induced by the optimal policy. This motivates Phase 2 in our algorithm, which stabilizes the covariance estimates.

Finally, this technique establishes the rate of convergence of the estimates. These rates are used to analyze the regret of the algorithm. They are also used to establish the asymptotic distribution of the parameter estimators at the terminal period of the algorithm and, therefore, allow for the construction of confidence intervals and hypothesis testing.

Our paper contributes to the literature on dynamic learning. Within this literature, a key emphasis has been the trade-off between exploration and exploitation (March 1991). The workhorse for analyzing this trade-off is the multi-armed bandit models. Economists have extensively studied their theoretical and empirical implications (see Bergemann and Välimäki 2008 for a review). Recently, there is a growing sub-literature that studies algorithmic decision-making based on multi-armed bandit models (Perchet et al. 2016, Caria et al. 2020, Currie and MacLeod 2020, Kasy and Sautmann 2021, Kawaguchi 2021, Kawaguchi, Uetake, and Watanabe 2021)<sup>2</sup>. The closest one to our study is that of Li, Raymond, and Bergman (2020), which also examines the performances of UCB-type algorithms on the contextual bandit problem. We add to this sub-literature by studying the biases that arise in algorithmic decision-making when there is an endogeneity problem, and we also show how to correct for it.

Our paper also contributes to the vast literature of casual inference (see, for example, Imbens and Rubin 2015 for a review). Most of this literature does not examine the environment of online learning, where the actions influence and are influenced by data. Two notable exceptions are Nambiar, Simchi-Levi, and Wang (2019) and Li, Luo, and Zhang (2021). Nambiar, Simchi-Levi, and Wang (2019) consider a dynamic pricing problem with learning, where the endogeneity problem arises from mis-specification of demand functions. Li, Luo, and Zhang (2021) consider a Markov decision process setup with learning and show that online learning environments may exacerbate the estimation bias. Different from these papers, our focus is on the dynamic selection problem where current data, through affecting the choice of actions, affects how future data is generated. This dynamic selection problem also differs from selection problems in Heckman (1979, 1990), Rosenbaum and Rubin (1983), Altonji, Elder, and Taber (2005), and Oster (2019), which do not feature how past selection problems affect future data selection.

The remainder of the paper is organized in the following manner. Section 2 introduces the background of the contextual bandit problem. Section 3 discusses the case where there exists at least one endogenous variable in the contextual bandit problem, which leads to the dynamic

<sup>2</sup> While we focus on online learning environments, there is a rich line of work that studies offline learning of optimal policies *after* the experimental or observational data is collected (Kitagawa and Tetenov 2018, Narita, Yasui, and Yata 2019, Athey and Wager 2021, Kallus and Zhou 2021, Narita, Yasui, and Yata 2021, Shi et al. 2021, Zhan et al. 2021).

selection problem and self-fulfilling bias. Section 4 presents a set of regularity conditions for IVs and the contextual bandit setup as well as proposes a class of IV-based bandit algorithms. Section 5 establishes the main results including regret bounds and asymptotic distribution for the proposed algorithms. Section 6 demonstrates the performance of these algorithms in a series of simulations. Section 7 presents the conclusions. Technical proofs are included the online Supplemental Material.

## 2. Background

Consider an agent who makes a sequence of decisions over a time horizon of  $T$  periods. At each time  $t = 1, \dots, T$ , the agent observes a vector of covariates  $\mathbf{v}_t \in \mathbb{R}^p$ , takes an action (i.e., pulls an arm)  $a_t \in \mathcal{A} = \{1, \dots, M\}$  with  $M \geq 2$  being a fixed positive integer, and receives a *random* reward  $R_t$ :

$$R_t = \sum_{i=1}^M \mathbb{I}(a_t = i) \mu_i(\mathbf{v}_t) + \epsilon_t, \quad (1)$$

where  $\mu_i: \mathbb{R}^p \mapsto \mathbb{R}$  is an unknown function, representing the expected reward associated with arm  $i$  for a given covariate,  $i = 1, \dots, M$ , and  $\epsilon_t$  is the noise having mean zero and standard deviation (SD)  $\sigma$ . We use  $\mathcal{H}_t := (\mathbf{v}_1, a_1, R_1, \dots, \mathbf{v}_t, a_t, R_t)$  to denote the history up to the end of period  $t$ . The agent seeks a *nonanticipating* policy  $\pi$  that maps  $(\mathcal{H}_{t-1}, \mathbf{v}_t)$  to  $a_t$  to maximize the expected cumulative reward over time. We write  $R_t^\pi = R_t$  to emphasize its dependence on  $\pi$ .

This sequential decision-making problem is commonly known as the contextual bandit problem (Slivkins 2019, chapter 8). For this class of problems, a standard metric for measuring the performance of  $\pi$  is its *expected cumulative regret*, which is referred to just as “regret” hereafter. Specifically, we compare  $\pi$  with an oracle policy  $\pi^*$ , which knows the functions  $\{\mu_i\}_{i=1}^M$  *a priori* and pulls arm  $\pi^*(\mathbf{v}_t) := \arg \max_{i=1, \dots, M} \mu_i(\mathbf{v}_t)$  at each time  $t$ . Following policy  $\pi$ , the agent incurs regret

$$\text{Regret}^\pi(T) := \sum_{t=1}^T \mathbb{E}[R_t^{\pi^*} - R_t^\pi], \quad (2)$$

which measures the difference in expected cumulative reward between  $\pi^*$  and  $\pi$ . Thus, the agent aims to approach the oracle’s performance by learning the reward functions gradually. In the sequel, we assume, for simplicity, that for each arm  $i$ , the reward function  $\mu_i$  is linear in the covariate vector—that is,  $\mu_i(\mathbf{v}) = \mathbf{v}^\top \boldsymbol{\alpha}_i$ , where  $\boldsymbol{\alpha}_i \in \mathbb{R}^p$  is a vector of unknown parameters. Then, the problem is reduced to learning these linear coefficients over time.

When the covariates are exogenous, a number of algorithms of Greedy or UCB types have been proposed to learn the coefficients. These algorithms differ in their finite-time and asymptotic properties, but they learn the true value of coefficients in the long run and, therefore, have a *sub-linear* regret—that is,  $T^{-1} \text{Regret}^\pi(T) \rightarrow 0$  as  $T \rightarrow \infty$ . This basically implies that at least the agent is taking the optimal action in most of the time periods as the horizon increases.

### 3. Self-fulfilling Bias: An Example

We now study a case in which the covariates are endogenous—that is,  $\mathbb{E}[\epsilon_t | \mathbf{v}_t] \neq 0$ . We show that the endogeneity of covariates leads to distortion in actions, which generates another type of endogeneity; that is, it creates “endogeneity spillover”. The additional endogeneity generates a dynamic selection problem and leads to a new type of bias, which we refer to as self-fulfilling bias.

Consider the following special case of our setup. Assume that there exist two arms, and between the two arms  $\mathcal{A} = \{1, 2\}$ , one is the safe arm, having a *known* expected reward that is independent of the covariates, while the other is the risky arm, whose expected reward has an unknown linear dependence on the covariates. Further, assume that the covariates are one-dimensional. In this simple example, the reward functions are

$$\mu_1(v) \equiv c \quad \text{and} \quad \mu_2(v) = \alpha v, \quad \forall v \in \mathbb{R}.$$

It follows that, if  $\alpha$  is known and positive, then the optimal policy for the agent is to pull arm 2 if  $v > c/\alpha$ , and pull arm 1 otherwise.

However, because  $\alpha$  is unknown, the agent needs to estimate its value over time and makes his decisions accordingly. In contrast to the standard regression analysis in which data is taken as given, in the dynamic environment considered here, the data used for estimating  $\alpha$  is only available when arm 2 is pulled (pulling arm 1 generates noisy observations of  $c$ , which are irrelevant to  $\alpha$ ). Now, suppose the agent uses ordinary least squares (OLS) to estimate  $\alpha$  and the estimates converge to a limit, say  $\hat{\alpha}$ . Then, in the long run,

$$\hat{\alpha} = \frac{\text{Cov}[v_t, R_t | R_t \text{ is generated from arm 2}]}{\text{Var}[v_t | R_t \text{ is generated from arm 2}]}.$$
 (3)

Note that given  $\hat{\alpha}$ , the agent’s policy in the long run is to pull arm 2 when  $v > c/\hat{\alpha}$ . Thus, the conditioning event in equation (3) is identical to  $\{v_t > c/\hat{\alpha}\}$ . Using the fact that  $R_t = \alpha v_t + \epsilon_t$  for arm 2, we have

$$\hat{\alpha} = \alpha + \frac{\text{Cov}[v_t, \epsilon_t | v_t > c/\hat{\alpha}]}{\text{Var}[v_t | v_t > c/\hat{\alpha}]}.$$
 (4)

It is evident from equation (4) that  $\hat{\alpha}$  is a fixed point as it appears on both sides of the equation. It also shows that the bias of the estimate,  $\hat{\alpha} - \alpha$ , has an expression that is similar to, but different from, that of the usual OLS bias, which is given by  $\frac{\text{Cov}[v_t, \epsilon_t]}{\text{Var}[v_t]}$ . We refer to the difference between them,  $(\hat{\alpha} - \alpha) - \frac{\text{Cov}[v_t, \epsilon_t]}{\text{Var}[v_t]}$ , as the *self-fulfilling bias* to reflect that the limit policy of the agent is induced by his limit belief (i.e., the limit estimate  $\hat{\alpha}$ ), and the limit belief is confirmed by the data generated from the limit policy.

The self-fulfilling bias arises because there are two—rather than one—types of endogeneity problems. First, when correlation exists between the noise and the covariates, it directly creates a bias in the estimate in a manner that is similar to the usual OLS bias. Second, this bias affects the agent’s actions, creating another source of endogeneity, which is essentially a type of bias caused by sample selection. Moreover, such sample selection is dynamic because when the policy that determines agent’s actions, the estimate of  $\alpha$  changes over time.

To illustrate the problem formally, we rewrite equation (1) as

$$R_t = c + \alpha v_t \mathbb{I}(a_t = 2) - c \mathbb{I}(a_t = 2) + \epsilon_t. \quad (5)$$

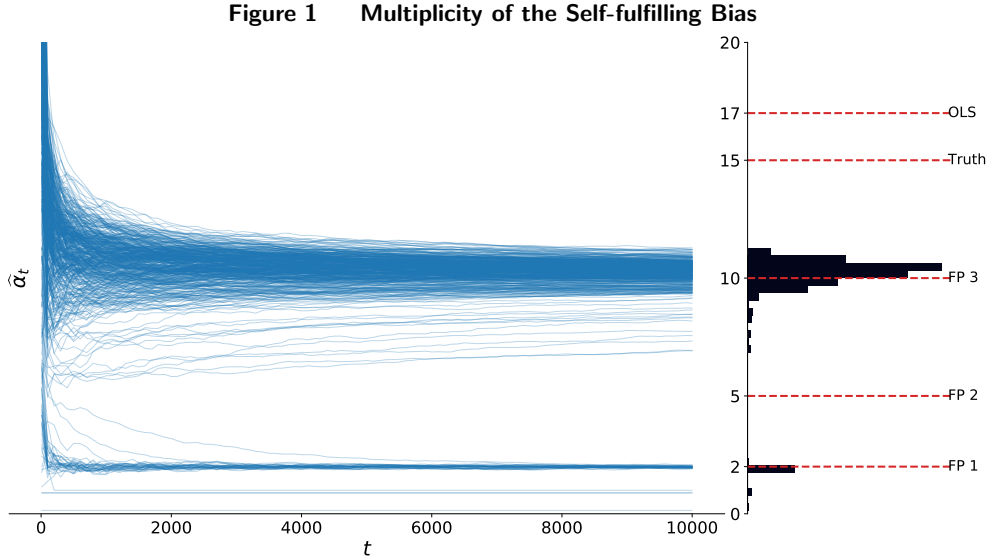
When  $\epsilon_t$  and  $v_t$  are correlated, the only endogenous variable in this regression is  $v_t \mathbb{I}(a_t = 2)$  if the arm  $a_t$  is selected independent of  $v_t$ . But because  $a_t$  depends on  $v_t$ —through the relation  $\mathbb{I}(a_t = 2) = \mathbb{I}(v_t > c/\hat{\alpha})$  in the long run—the correlation between  $\epsilon_t$  and  $v_t$  leads to a correlation between  $\mathbb{I}(a_t = 2)$  and  $\epsilon_t$ . Therefore, the regression equation (5) contains two endogenous variables— $v_t \mathbb{I}(a_t = 2)$  and  $\mathbb{I}(a_t = 2)$ , reflecting the endogeneity spillover problem.

The additional endogenous variable complicates the dependence of the value of the limit coefficient estimate ( $\hat{\alpha}$ ) on the joint distribution of  $v_t$  and  $\epsilon_t$ . In particular, if the term  $\frac{\text{Cov}[v_t, \epsilon_t | v_t > s]}{\text{Var}[v_t | v_t > s]}$  is the same for all cutoff  $s$ —as is the case when  $v_t$  and  $\epsilon_t$  have a joint normal distribution<sup>3</sup>—then  $\hat{\alpha}$  is identical to the OLS estimate. In general, however, the coefficient estimate can be either higher or lower than the OLS estimate, depending on the joint distribution of  $v_t$  and  $\epsilon_t$ . When  $\mathbb{E}[\epsilon_t | v_t]$  is a nonlinear function of  $v_t$ , there will be additional bias in general.

More importantly, the value of  $\hat{\alpha}$  (and hence the value of the self-fulfilling bias) may not be unique, which happens when equation (4) contains multiple solutions. Figure 1 illustrates one such case and shows that the agents’ beliefs converge to two different values in the long run. With multiple self-fulfilling biases, this example shows that agents can adopt different policies in the long run, even if the underlying production environment is the same. It therefore provides one explanation for why seemingly similar enterprises may choose different practices and result in persistent differences in performance, which is a central question in organizational economics and strategy (see, for example, Gibbons and Henderson 2012 for a review).

COMMENT 1. There appear to be little restrictions on the relationship between the true value  $\alpha$  and the belief that an agent has in the long run. Effectively, for a given  $\alpha$ , the beliefs regarding  $\alpha$  may converge in the long run to *any* number of *arbitrarily chosen* values in the long run. In other words, for a group of agents with the same initial beliefs, they can hold many very different beliefs in the long run. We discuss why this is the case in the Appendix.  $\square$

<sup>3</sup> If  $\mathbb{E}[\epsilon_t | v_t]$  is a linear function of  $v_t$ , then  $\frac{\text{Cov}[v_t, \epsilon_t | v_t > s]}{\text{Var}[v_t | v_t > s]}$  is a constant with respect to  $s$ . That  $(v_t, \epsilon_t)$  has a joint normal distribution is a special case of this condition.



*Note.* Suppose that  $c = 1$ ,  $\alpha = 15$ ,  $v_t$  is uniformly distributed on  $[0, 1]$ , and  $\epsilon_t = \sum_{k=0}^3 \beta_k v_t^k + \eta_t$ , where  $\beta_k$ 's are some constants and  $\eta_t$  is an independent standard normal random variable. Then, (i) the limit of the OLS estimate is 17 and (ii) equation (4) is reduced to a cubic equation in terms of  $\hat{\alpha}$ , thereby yielding three distinct real fixed points (FPs)—2, 5, and 10—if  $\beta_k$ 's are appropriately selected (see the Appendix for details). The left portion of this figure shows 500 sample paths of estimates of  $\alpha$  following a greedy policy, while the right portion shows the histogram of the terminal values of these sample paths. The simulation results indicate that the greedy estimates mostly converge to either 2 or 10.

We illustrate the self-fulfilling bias using the simplest possible example. The underlying problem is more general because in any dynamic environment where the agent's actions depend on the covariates, endogeneity in the covariates will likely generate additional endogeneity with the actions. Note that when the coefficient estimate differs from the true parameter in the long run because of the self-fulfilling bias, the fraction of time a wrong decision is taken is positive in the limit. In other words, not only is the policy suboptimal in the long run, the number of periods in which a wrong action is taken is proportional to the time horizon  $T$ . Thus, the regret of the policy with self-fulfilling bias is in the order of  $\mathcal{O}(T)$ .

To correct for the self-fulfilling bias, it is tempting to have the agent to carry out the exploration in the conventional sense, i.e., to choose actions randomly sometimes. However, as long as the data from non-random actions is used, the agent's parameter estimate is biased because the endogeneity spillover problem renders  $\mathbf{v}_t \mathbb{I}(a_t = 2)$  an endogenous term in the regression equation (5). To correct for the endogeneity in  $\mathbf{v}_t$ , the solution is not more exploration in the conventional sense. Rather, it is *causal exploration*—that is, to use the IVs to “perturb” the data-generating process. We discuss how to correct for the self-fulfilling bias in the next section.

## 4. IV-Bandit Algorithms

In this section, we propose a class of algorithms that incorporates IVs into the online bandit algorithm. These algorithms update a policy function  $\pi(\cdot|v)$  that approaches the optimal one gradually, while retaining a regret of the order of  $\mathcal{O}(\log(T))$  or  $\mathcal{O}(\log^2(T))$  in the long run.

Throughout the rest of the paper, we use  $\|\mathbf{v}\|$  to denote the Euclidean norm if  $\mathbf{v}$  is a vector in a finite dimensional Euclidean space. The following definitions are also needed.

**DEFINITION 1.** Let  $\mathbf{M} \in \mathbb{R}^{k \times \ell}$  be a real matrix. If  $k \geq \ell$ , then a non-negative real number  $s$  is said to be a *singular value* of  $\mathbf{M}$  if  $s^2$  is an eigenvalue of  $\mathbf{M}^\top \mathbf{M} \in \mathbb{R}^{\ell \times \ell}$ . If  $k < \ell$ , then a non-negative real number  $s$  is said to be a singular value of  $\mathbf{M}$  if  $s^2$  is an eigenvalue of  $\mathbf{M} \mathbf{M}^\top \in \mathbb{R}^{k \times k}$ . Hence,  $\mathbf{M} \in \mathbb{R}^{k \times \ell}$  has  $\min\{k, \ell\}$  singular values in total. Let  $\phi_{\max}(\mathbf{M})$  and  $\phi_{\min}(\mathbf{M})$  denote the largest and smallest singular values of  $\mathbf{M}$ , respectively.

**DEFINITION 2.** The *spectral norm* of a matrix  $\mathbf{M} \in \mathbb{R}^{k \times \ell}$  is defined as  $\|\mathbf{M}\| := \sup_{\mathbf{u} \in \mathbb{R}^\ell: \|\mathbf{u}\|=1} \|\mathbf{M}\mathbf{u}\|$ , and its *max norm* is defined as  $\|\mathbf{M}\|_{\max} := \max_{1 \leq i \leq k, 1 \leq j \leq \ell} |M_{ij}|$ , where  $M_{ij}$  denotes the  $(i, j)$ -th entry of  $\mathbf{M}$ .

**DEFINITION 3.** A random variable  $X$  is said to be *sub-Gaussian* with variance proxy  $\varsigma^2$  if  $\mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] \leq e^{\lambda^2 \varsigma^2 / 2}$  for all  $\lambda \in \mathbb{R}$ .

### 4.1. General Setup

Similar to (5), because  $\sum_{i=1}^M \mathbb{I}(a_t = i) = 1$ , we rewrite the reward function (1) as

$$R_t = \mathbf{v}_t^\top \boldsymbol{\alpha}_1 + \mathbb{I}(a_t = 2) \mathbf{v}_t^\top (\boldsymbol{\alpha}_2 - \boldsymbol{\alpha}_1) + \cdots + \mathbb{I}(a_t = M) \mathbf{v}_t^\top (\boldsymbol{\alpha}_M - \boldsymbol{\alpha}_1) + \epsilon_t. \quad (6)$$

With this formulation, we note that there are at most  $Mp$  endogenous variables in (6), where  $p$  is the dimensionality of  $\mathbf{v}_t$ . To see this, note that, first, there are at most  $p$  endogenous variables within  $\mathbf{v}_t$ . In addition, as our discussion following (5) illustrates, the optimal policy is

$$a_t^* := \arg \max_{i=1, \dots, M} \{\mathbf{v}_t^\top \boldsymbol{\alpha}_i\},$$

which is dependent on  $\mathbf{v}_t$ . As the action  $a_t$  converges to  $a_t^*$ ,  $a_t$  will also be endogenous. This brings another  $p$  endogenous variables apart from  $\mathbf{v}_t$  itself. To cope with these endogenous variables, we assume that there are at least  $Mp$  IVs.

**ASSUMPTION 1.** Let  $\{(\mathbf{v}_t, \mathbf{z}_t, \epsilon_t) : t = 1, \dots, T\}$  be an *i.i.d.* sequence, where  $\mathbf{z}_t \in \mathbb{R}^q$  is a vector of valid IVs with  $q \geq Mp$ .

(i)  $\epsilon_t$  is sub-Gaussian with variance proxy  $\varsigma^2$

(ii)  $\mathbb{E}[\mathbf{z}_t \epsilon_t] = 0$  and  $\text{Var}(\epsilon_t) = \sigma^2 < \infty$ .

(iii) There exist positive constants  $\bar{v}$  and  $\bar{z}$ , such that  $\|\mathbf{z}_t\|_\infty \leq \bar{z}$  and  $\|\mathbf{v}_t\|_\infty \leq \bar{v}$  almost surely, where  $\|\cdot\|_\infty$  denotes the  $L^\infty$  norm of a vector.

(iv) Let  $\boldsymbol{\delta}_{i,j} := \boldsymbol{\alpha}_i - \boldsymbol{\alpha}_j$ ,  $\boldsymbol{\Sigma}_{vz}^* := \mathbb{E}[\mathbf{v}_t^* \mathbf{z}_t^\top]$ , and  $\boldsymbol{\Sigma}_{zz} := \mathbb{E}[\mathbf{z}_t \mathbf{z}_t^\top]$ , where

$$\mathbf{v}_t^* := \begin{pmatrix} \mathbb{I}(a_t^* = 1) \mathbf{v}_t \\ \vdots \\ \mathbb{I}(a_t^* = M) \mathbf{v}_t \end{pmatrix} = \begin{pmatrix} \mathbb{I}(\min_{j \neq 1} \mathbf{v}_t^\top \boldsymbol{\delta}_{1j} > 0) \mathbf{v}_t \\ \vdots \\ \mathbb{I}(\min_{j \neq M} \mathbf{v}_t^\top \boldsymbol{\delta}_{Mj} > 0) \mathbf{v}_t \end{pmatrix} \in \mathbb{R}^{Mp}.$$

Then,  $\boldsymbol{\Sigma}_{vz}^*$  is a full-rank matrix and  $\boldsymbol{\Sigma}_{zz}$  is a positive definite matrix.

(v) There exists a constant  $L > 0$ , such that  $\Pr(|\mathbf{v}_t^\top \boldsymbol{\delta}_{i,j}| \leq c) \leq Lc$  for all  $c > 0$  and  $1 \leq i \neq j \leq M$ .

Condition (i) is common in the bandit literature and is key for establishing a  $\log(T)$ -like regret bound. For example, Goldenshluger and Zeevi (2013) assume normality on  $\epsilon_t$ , and Zhong, Hong, and Liu (2021) assume boundedness, both of which are special cases of Condition (i). Bastani, Bayati, and Khosravi (2021) impose the same sub-Gaussian condition as we do. If Condition (i) is replaced with moment restrictions on  $\epsilon_t$ , one might still obtain a sublinear regret, but not a  $\log(T)$  regret. Conditions (ii) and (iii) are again common assumptions in the IV and contextual bandit literature, respectively.

Condition (iv) implies that  $\boldsymbol{\Sigma}_{vz}^* \boldsymbol{\Sigma}_{zz}^{-1} (\boldsymbol{\Sigma}_{vz}^*)^\top \in \mathbb{R}^{(Mp) \times (Mp)}$  is a positive definite matrix. This is a key condition for parameter identification. It requires that the optimal policy must be dependent on  $\mathbf{v}_t$ ; more specifically,  $\mathbb{I}(a_t^* = 1) \mathbf{v}_t, \dots, \mathbb{I}(a_t^* = M) \mathbf{v}_t$  should not be co-linear when projected to the  $\mathbf{z}_t$  space. Condition (iv) holds in general if  $\boldsymbol{\delta}_{i,j} \neq \mathbf{0}$  for all  $1 \leq i \neq j \leq M$ .

Further, Condition (v) states that the density of  $\mathbf{v}_t$  is Lipschitz continuous in the vicinity of the hyperplane  $\{\mathbf{v} : \mathbf{v}^\top \boldsymbol{\delta}_{i,j} = 0\}$ . This condition is often referred to as the ‘‘margin condition’’ in statistical learning literature (Tsybakov 2004). It is also a standard condition in the literature on contextual bandits. Moreover, Conditions (iv) and (v) together imply that  $\Pr(\mathbf{v}_t^\top \boldsymbol{\delta}_{i,j} < 0) > 0$  and  $\Pr(\mathbf{v}_t^\top \boldsymbol{\delta}_{i,j} > 0) > 0$ , which is referred to as the ‘‘diversity condition’’ in Goldenshluger and Zeevi (2013).

COMMENT 2. The covariate vector  $\mathbf{v}_t \in \mathbb{R}^p$  may be decomposed into two parts:  $\mathbf{x}_t$  and  $\mathbf{d}_t$ , where  $\mathbf{x}_t \in \mathbb{R}^\ell$  is exogenous,  $\mathbf{d}_t \in \mathbb{R}^m$  is potentially endogenous, and  $p = \ell + m$ . Then, following the discussion below (6), there are  $((M-1)\ell + Mm)$  number of endogenous variables in (6), with  $m$  of them from the term  $\mathbf{v}_t^\top \boldsymbol{\alpha}_1$  and the other  $(M-1)(\ell+m)$  of them from the term  $\mathbb{I}(a_t = i) \mathbf{v}_t^\top (\boldsymbol{\alpha}_i - \boldsymbol{\alpha}_1)$ ,  $i = 2, \dots, M$ . Thus, the minimal number of IVs that we need is  $((M-1)\ell + Mm)$ . However, for the convenience of exposition, we stipulate in Assumption 1 that there are a total of  $q \geq Mp = M(\ell+m)$  IVs.  $\square$

COMMENT 3. To construct a sufficient number of IVs, following the decomposition of  $\mathbf{v}_t$  in Comment 2, practitioners may begin with a set of IVs for  $\mathbf{d}_t$ , denoted as  $\check{\mathbf{z}}_t \in \mathbb{R}^{\check{q}}$  with  $\check{q} \geq m$ . The exogenous variables  $\mathbf{x}_t$  can serve as IVs for themselves. A natural means to construct additional IVs for  $\mathbb{I}(a_t = i)\mathbf{v}_t$ ,  $i = 2, \dots, M$ , in (6) is to use such variables as  $\mathbb{I}(\mathbf{x}_t \geq \hat{x})$ ,  $\mathbb{I}(\check{\mathbf{z}}_t \geq \hat{z})$ , and their interaction terms—for example,  $\mathbb{I}(\mathbf{x}_t \geq \hat{x})\mathbf{x}_t$  and  $\mathbb{I}(\check{\mathbf{z}}_t \geq \hat{z})\check{\mathbf{z}}_t$  for some threshold values  $\hat{x}$  and  $\hat{z}$ . See the simulation example in Section 6 for an illustration. Other nonlinear transformations of  $\mathbf{x}_t$  and  $\check{\mathbf{z}}_t$  may also work as IVs, if

$$\mathbb{E}[\epsilon_t | \mathbf{x}_t, \check{\mathbf{z}}_t] = 0. \quad (7)$$

This condition and the use of splines as IVs are widely employed in the conditional moment restriction literature (Chen and Pouzo 2012). This implies Condition (iii) in Assumption 1 and, thus, it is slightly stronger. For optimal selection of IVs under the conditional moment restriction (7), we refer to Donald, Imbens, and Newey (2009) and Belloni et al. (2012).  $\square$

## 4.2. Structure

In this subsection, we propose a class of IV-based algorithms for linear contextual bandits and name it *IV-Bandit*. This algorithm consists of three phases. In the first phase—which we refer to as *Coefficient Stabilization*—we seek a stable estimate of the coefficient vector  $\alpha$ , thereby implying that the estimate should stay close to the true value with a high probability. For this purpose, our *behavior policy* is to select each arm randomly—irrespective of the observed covariates—with equal probability. At the end of Phase 1, we run *arm-specific-2SLS* to obtain an estimate of the coefficient vector  $\alpha$ . Specifically, depending on the arm pulled in each period—whether  $a_t = 1, 2$ , etc.—we divide the time periods into  $M$  groups and then conduct 2SLS separately on the data (i.e., rewards, covariates, and IVs) in each group. Note that the estimate is unbiased because the actions are generated randomly, and the IVs are used to correct for the endogeneity in the covariates. But the absence of bias is achieved at the cost of a regret that is linear in the number of periods  $T_1$ . To balance the regret and the purpose of coefficient stabilization, we choose the duration of Phase 1 to be  $T_1 = \mathcal{O}(\log(T))$ .

In the second phase—which we refer to as *Covariance Stabilization*—we seek a stable estimate of the covariance matrix  $\mathbf{\Omega}^* := (\mathbf{\Sigma}_{vz}^* \mathbf{\Sigma}_{zz}^{-1} (\mathbf{\Sigma}_{vz}^*)^\top)^{-1}$ . In this phase, the behavioral policy takes the form of a greedy policy—that is, it uses the estimate of  $\alpha$  from Phase 1 to evaluate the expected payoff of each arm for a given value of the covariates, and then selects an arm that maximizes the expected payoff. During this phase, which lasts from period  $T_1 + 1$  to  $T_2$ , no updates on the coefficients are performed. As explained below, fixing the coefficients helps the algorithm to obtain a

good estimate of  $\mathbf{\Omega}^*$  at the beginning of Phase 3, which is crucial for the estimate of  $\boldsymbol{\alpha}$  to converge to the true value. Similar to Phase 1, we set the duration of this phase to be  $T_2 - T_1 = \mathcal{O}(\log(T))$  so that sufficient data is generated without causing excessive regret.

In the third phase—which we refer to as *Policy Improvement*—our algorithm allows for a class of possible behavioral policies, depending on the choice of a parameter  $\theta$ . In the simplest form, which corresponds to the case of  $\theta = 0$ , the behavioral policy is again greedy: we select the arm that maximizes the estimated expected reward in each period. However, in contrast to Phase 2, the criterion for evaluating the expected reward is constantly updated. In each period, the expected reward for each arm is calculated from the most updated estimate of  $\boldsymbol{\alpha}$ . When  $\theta > 0$ , the behavioral policy in Phase 3 of our algorithm becomes a UCB type. UCB is a common strategy for solving contextual bandit problems (Chu et al. 2011, Dimakopoulou et al. 2019, Zhong, Hong, and Liu 2021). In contrast to the greedy policy, UCB takes into account the uncertainty regarding its expected reward when deciding an arm, thereby balancing between exploitation and exploration. Specifically, UCB selects the arm that maximizes a weighted sum of the estimated expected reward and the estimated SD. In our algorithm, the value of  $\theta$  represents the relative significance of the SD, thereby reflecting the level of exploration. Again, both the expected reward and SD in each period are estimated via the joint-2SLS using all the available data from the beginning of Phase 2.

Two features of our algorithms, the use of joint-2SLS in Phase 3 and the existence of Phase 1 and 2, are crucial for generating unbiased estimation by dealing with the endogeneity spillover issue. In particular, one might consider applying the arm-specific-2SLS in Phase 3 (as done in Phase 1) to obtain estimates of  $\boldsymbol{\alpha}$  and  $\mathbf{\Omega}^*$ . This approach leads to biased estimates because it copes with only the endogeneity in the covariate  $\mathbf{v}_t$ , but not the endogeneity in the action  $a_t$ , which—as explained in detail in Section 5—is propagated from that in  $\mathbf{v}_t$  through the behavior policy that again depends on  $\mathbf{v}_t$ . To account for this problem, joint-2SLS conducts 2SLS with  $M$  regressors  $\{\mathbb{I}(a_t = i)\mathbf{v}_t : i = 1, \dots, M\}$  in addition to the IVs  $\mathbf{z}_t$ . This formulation clarifies that the minimal number of IVs should be  $Mp$ , rather than  $p$ , as the number of regressors is  $M$  times as many.

In addition, one might consider beginning the algorithm immediately with Phase 3—without Phases 1 and 2—using the joint-2SLS starting from the first time period to estimate both the coefficients  $\boldsymbol{\alpha}$  and the covariance matrix  $\mathbf{\Omega}^*$  and select the arms accordingly. However, doing so again leads to biased estimates in the long run. This problem also arises because of the interactions between the estimates and the chosen actions. Poor estimates of the coefficients and the covariance matrix in the past, through their effects on the chosen actions, have persistent effects in the future and lead to biased estimates in the long run. Phase 1 ensures that the initial coefficient estimate is

sufficiently accurate. Then, Phase 2 uses the greedy policy based on the initial coefficient estimate to generate “proxies” of  $\mathbf{v}_t^*$ ’s, which, by definition, are induced by the unknown optimal policy. This allows us to obtain a stable estimate in the first-stage regression of the joint-2SLS, which is critical for producing a sample path of  $\hat{\boldsymbol{\alpha}}_t$  that converges to the true value  $\boldsymbol{\alpha}$ . In Section 5, we discuss how our algorithm deals with this problem.

COMMENT 4. Because the reward associated with an arm is sampled only if that arm is selected, and the decision to select an arm depends on the estimate of  $\boldsymbol{\alpha}$ , each arm effectively faces a sample selection problem. Apart from our algorithm, one may apply the bias correction technique (Heckman 1979, 1990) to address the sample selection problem. This method often requires a normality assumption or some other parametric assumption on the distribution of the error terms. A less restrictive method is the semi-parametric maximum likelihood estimation (Cosslett 2004). Both of these methods usually require the error terms to be identically distributed, while the IV-Bandit algorithm may work in the presence of heteroskedasticity. Moreover, they are developed in a static environment. Extending them to an online learning context is beyond the scope of the present paper and we defer the investigation to future research.  $\square$

### 4.3. Details

To facilitate the subsequent presentation, we introduce the following notations. For each  $i = 1, \dots, M$  and  $t$ , let  $\mathcal{J}_{i,t} := \{s \leq t | a_s = i\}$  be the set of time periods when arm  $i$  is pulled within the first  $t$  periods. Let  $\mathbf{V}_{i,t} \in \mathbb{R}^{|\mathcal{J}_{i,t}| \times p}$  denote the matrix for data collected in time periods in  $\mathcal{J}_{i,t}$ , composed of rows  $\mathbf{v}_s^\top$  for all  $s \in \mathcal{J}_{i,t}$ . Likewise, we denote by  $\mathbf{R}_{i,t} \in \mathbb{R}^{|\mathcal{J}_{i,t}|}$  the vector composed of elements  $R_{i,s}$  and by  $\mathbf{Z}_{i,t} \in \mathbb{R}^{|\mathcal{J}_{i,t}| \times q}$  the vector composed of rows  $\mathbf{z}_{i,s}^\top$ , for all  $s \in \mathcal{J}_{i,t}$ . Let

$$\boldsymbol{\alpha} := \begin{pmatrix} \boldsymbol{\alpha}_1 \\ \vdots \\ \boldsymbol{\alpha}_M \end{pmatrix} \in \mathbb{R}^{Mp}, \quad \hat{\boldsymbol{\alpha}}_t := \begin{pmatrix} \hat{\boldsymbol{\alpha}}_{1,t} \\ \vdots \\ \hat{\boldsymbol{\alpha}}_{M,t} \end{pmatrix} \in \mathbb{R}^{Mp}, \quad \tilde{\mathbf{v}}_t := \begin{pmatrix} \mathbb{I}(a_t = 1)\mathbf{v}_t \\ \vdots \\ \mathbb{I}(a_t = M)\mathbf{v}_t \end{pmatrix} \in \mathbb{R}^{Mp}.$$

For any  $t_1 < t_2$ , let

$$\tilde{\mathbf{V}}_{t_1:t_2} := \begin{pmatrix} \tilde{\mathbf{v}}_{t_1+1}^\top \\ \vdots \\ \tilde{\mathbf{v}}_{t_2}^\top \end{pmatrix} \in \mathbb{R}^{(t_2-t_1) \times (Mp)}, \quad \mathbf{Z}_{t_1:t_2} := \begin{pmatrix} \mathbf{z}_{t_1+1}^\top \\ \vdots \\ \mathbf{z}_{t_2}^\top \end{pmatrix} \in \mathbb{R}^{(t_2-t_1) \times q}, \quad \mathbf{R}_{t_1:t_2} := \begin{pmatrix} R_{t_1+1} \\ \vdots \\ R_{t_2} \end{pmatrix} \in \mathbb{R}^{t_2-t_1}.$$

We define a projection operator  $\mathcal{P}$ , which takes effect on a matrix  $\mathbf{Z}$ , as  $\mathcal{P}[\mathbf{Z}] = \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top$ . The IV-Bandit algorithm is formally presented in Algorithm 1. For the ease of reference, we call the algorithm *IV-Greedy* when  $\theta = 0$  and *IV-UCB* when  $\theta > 0$ .

**Algorithm 1: IV-Bandit**


---

**Initialization** : Set  $\theta$ ,  $T_1$ ,  $T_2$ , and  $T$

**Coefficient Stabilization** : **for**  $t = 1, \dots, T_1$  **do**  
 | Observe  $(\mathbf{v}_t, \mathbf{z}_t)$ , pull arm  $a_t = 1, \dots, M$  purely randomly, and collect  $R_t$ .  
**end**

**for**  $i = 1, \dots, M$  **do**  
 | Run arm-specific-2SLS on  $(\mathbf{V}_{i,T_1}, \mathbf{Z}_{i,T_1}, \mathbf{R}_{i,T_1})$  to estimate  $\boldsymbol{\alpha}$ :  

$$\hat{\boldsymbol{\alpha}}_{i,T_1} := (\mathbf{V}_{i,T_1}^\top \mathcal{P}[\mathbf{Z}_{i,T_1}] \mathbf{V}_{i,T_1})^{-1} \mathbf{V}_{i,T_1}^\top \mathcal{P}[\mathbf{Z}_{i,T_1}] \mathbf{R}_{i,T_1}.$$
  
**end**

**Covariance Stabilization** : **for**  $t = T_1 + 1, \dots, T_2$  **do**  
 | Observe  $(\mathbf{v}_t, \mathbf{z}_t)$ , pull arm  $a_t = \arg \max_{i=1, \dots, M} \{\mathbf{v}_t^\top \hat{\boldsymbol{\alpha}}_{i,T_1}\}$ , and collect  $R_t$ .  
**end**

**Policy Improvement** : **for**  $t = T_2 + 1, \dots, T$  **do**  
 | **if**  $t = T_2 + 1$  **then**  
 | | **for**  $i = 1, \dots, M$  **do**  
 | | | Set  $\hat{\boldsymbol{\alpha}}_{i,T_2} := \hat{\boldsymbol{\alpha}}_{i,T_1}$  and  $\hat{\boldsymbol{\Omega}}_{i,T_2} := (\mathbf{V}_{i,T_1}^\top \mathcal{P}[\mathbf{Z}_{i,T_1}] \mathbf{V}_{i,T_1})^{-1}$ .  
 | | **end**  
 | | Set  $\hat{\sigma}_{T_2} := (T_1^{-1} \sum_{i=1}^M \|\mathbf{R}_{i,T_1} - \mathbf{V}_{i,T_1} \hat{\boldsymbol{\alpha}}_{i,T_1}\|^2)^{\frac{1}{2}}$ .  
 | **end**  
 | Observe  $(\mathbf{v}_t, \mathbf{z}_t)$ , pull arm  $a_t = \arg \max_{i=1, \dots, M} \left\{ \mathbf{v}_t^\top \hat{\boldsymbol{\alpha}}_{i,t-1} + \hat{\sigma}_{t-1} \sqrt{2\theta \log(t - T_1) \cdot \mathbf{v}_t^\top \hat{\boldsymbol{\Omega}}_{i,t-1} \mathbf{v}_t} \right\}$ , and collect  $R_t$ .  
 | Run joint-2SLS on  $(\tilde{\mathbf{V}}_{T_1:t}, \mathbf{Z}_{T_1:t}, \mathbf{R}_{T_1:t})$  to update the estimates of  $\boldsymbol{\Omega}^*$  and  $\boldsymbol{\alpha}$ :  

$$\hat{\boldsymbol{\Omega}}_t := (\tilde{\mathbf{V}}_{T_1:t}^\top \mathcal{P}[\mathbf{Z}_{T_1:t}] \tilde{\mathbf{V}}_{T_1:t})^{-1} \quad \text{and} \quad \hat{\boldsymbol{\alpha}}_t := \hat{\boldsymbol{\Omega}}_t \tilde{\mathbf{V}}_{T_1:t}^\top \mathcal{P}[\mathbf{Z}_{T_1:t}] \mathbf{R}_{T_1:t}.$$
  
 | Set  $\hat{\boldsymbol{\Omega}}_{i,t}$  as the  $i$ -th  $p \times p$  diagonal block of  $\hat{\boldsymbol{\Omega}}_t \in \mathbb{R}^{(Mp) \times (Mp)}$ , respectively.  
 | **if**  $\theta > 0$  **then**  
 | | Update the estimate of  $\sigma$ :  $\hat{\sigma}_t = (t - T_1)^{-\frac{1}{2}} \|\mathbf{R}_{T_1:t} - \tilde{\mathbf{V}}_{T_1:t} \hat{\boldsymbol{\alpha}}_t\|$ .  
 | **end**  
**end**

---

COMMENT 5. Algorithm 1 may be enhanced by using  $\check{\boldsymbol{\alpha}}_t := \frac{T_1}{t} \hat{\boldsymbol{\alpha}}_{T_1} + (1 - \frac{T_1}{t}) \hat{\boldsymbol{\alpha}}_t$ , a weighted average between  $\hat{\boldsymbol{\alpha}}_{T_1}$  and  $\hat{\boldsymbol{\alpha}}_t$ , as the estimate of  $\boldsymbol{\alpha}$  in Phase 3. This may somewhat reduce the necessary duration of Phase 2, which nevertheless would still be of order  $\mathcal{O}(\log(T))$ . In other words, this variant would not affect the order of the regret bound. Due to the additional technical complexity this variant would introduce, we focus on the current formulation of Algorithm 1.  $\square$

## 5. Theoretical Results

In this section, we analyze the theoretical properties of the IV-Bandit algorithms. We show that IV-Greedy produces a regret of order  $\mathcal{O}(\log(T))$ , while IV-UCB produces a regret of order  $\mathcal{O}(\log^2(T))$ . We also show that these algorithms deliver consistent and asymptotically normal estimates of  $\boldsymbol{\alpha}$ .

### 5.1. IV-Greedy

Below, we first overview the regret analysis of the IV-Greedy algorithm, and then discuss a new technique—zig-zag induction—that we develop for analyzing algorithms associated with online learning environments

#### 5.1.1. Overview of the Regret Analysis

**THEOREM 1.** *Let  $T_1 = C_{T_1} \log(T)$  and  $T_2 = (C_{T_1} + C_{T_2}) \log(T)$  for some sufficiently large constants  $C_{T_1}$  and  $C_{T_2}$ . Then, under Assumption 1, the regret of the IV-Greedy Algorithm (i.e., Algorithm 1 with  $\theta = 0$ ) is  $\text{Regret}(T) = \mathcal{O}(\log(T))$  for all  $T > T_2$ .*

Note that in the presence of endogeneity, an algorithm that does not produce consistent estimates of the coefficient vector  $\boldsymbol{\alpha}$  will generally suffer a linear regret (Nambiar, Simchi-Levi, and Wang 2019). Therefore, it is necessary to establish consistency of the coefficient estimate. This requires the analysis of  $\Pr(\|\hat{\boldsymbol{\alpha}}_t - \boldsymbol{\alpha}\| < \eta)$  for small  $\eta$  over Phase 3 of the algorithm. In Phase 3, the 2SLS formula yields

$$\begin{aligned} \hat{\boldsymbol{\alpha}}_t &= \boldsymbol{\alpha} + \hat{\boldsymbol{\Omega}}_t \tilde{\mathbf{V}}_{T_1:t}^\top \mathcal{P}[\mathbf{Z}_{T_1:t}] \boldsymbol{\epsilon}_{T_1:t} \\ &= \boldsymbol{\alpha} + (\hat{\boldsymbol{\Sigma}}_{vz, T_1:t} \hat{\boldsymbol{\Sigma}}_{zz, T_1:t}^{-1} \hat{\boldsymbol{\Sigma}}_{vz, T_1:t}^\top)^{-1} \hat{\boldsymbol{\Sigma}}_{vz, T_1:t} \hat{\boldsymbol{\Sigma}}_{zz, T_1:t}^{-1} \hat{\mathbf{G}}_{T_1:t}, \end{aligned} \quad (8)$$

where for any  $t_1 < t_2$ ,  $\boldsymbol{\epsilon}_{t_1:t_2} := (\epsilon_{t_1+1}, \dots, \epsilon_{t_2})^\top \in \mathbb{R}^{t_2-t_1}$ , and

$$\begin{aligned} \hat{\boldsymbol{\Sigma}}_{vz, t_1:t_2} &:= \frac{1}{t_2 - t_1} \sum_{s=t_1+1}^{t_2} \tilde{\mathbf{v}}_s \mathbf{z}_s^\top = \frac{1}{t_2 - t_1} \tilde{\mathbf{V}}_{t_1:t_2}^\top \mathbf{Z}_{t_1:t_2} \in \mathbb{R}^{(Mp) \times q}, \\ \hat{\boldsymbol{\Sigma}}_{zz, t_1:t_2} &:= \frac{1}{t_2 - t_1} \sum_{s=t_1+1}^{t_2} \mathbf{z}_s \mathbf{z}_s^\top = \frac{1}{t_2 - t_1} \mathbf{Z}_{t_1:t_2}^\top \mathbf{Z}_{t_1:t_2} \in \mathbb{R}^{q \times q}, \\ \hat{\mathbf{G}}_{t_1:t_2} &:= \frac{1}{t_2 - t_1} \sum_{s=t_1+1}^{t_2} \mathbf{z}_s \epsilon_s = \frac{1}{t_2 - t_1} \mathbf{Z}_{t_1:t_2}^\top \boldsymbol{\epsilon}_{t_1:t_2} \in \mathbb{R}^q. \end{aligned}$$

Expression (8) shows that  $\hat{\boldsymbol{\alpha}}_t$  depend on  $\hat{\boldsymbol{\Sigma}}_{zz, T_1:t}$ ,  $\hat{\boldsymbol{\Sigma}}_{vz, T_1:t}$ , and  $\hat{\mathbf{G}}_{T_1:t}$ . The calculations of  $\hat{\boldsymbol{\Sigma}}_{zz, T_1:t}$  and  $\hat{\mathbf{G}}_{T_1:t}$  are standard. They are averages of samples of i.i.d. random variables and their statistical properties are standard. For expositional ease, we do not discuss these two terms below and relegate their analysis to the proof in the Supplemental Material. In the following discussion, we refer to  $\hat{\boldsymbol{\alpha}}_t$  as the coefficient estimate and to  $\hat{\boldsymbol{\Sigma}}_{vz, T_1:t}$  as the covariance estimate.

However, the calculation of  $\hat{\boldsymbol{\Sigma}}_{vz, T_1:t}$  is not standard because it depends on the past estimates  $\{\hat{\boldsymbol{\alpha}}_s : T_1 \leq s \leq t-1\}$  through  $\tilde{\mathbf{V}}_{T_1:t}$ —that is,

$$\hat{\boldsymbol{\Sigma}}_{vz, T_1:t} = \frac{1}{t - T_1} \sum_{s=T_1+1}^t \begin{pmatrix} \mathbb{I}(a_s = 1) \mathbf{v}_s \\ \vdots \\ \mathbb{I}(a_s = M) \mathbf{v}_s \end{pmatrix} \mathbf{z}_s^\top, \quad (9)$$

where  $a_s$  is selected by the greedy algorithm and depends on  $\widehat{\boldsymbol{\alpha}}_s$ .

Note that if the optimal arm  $a_s^*$  were known and selected in each period, we would simply calculate

$$\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}^* := \frac{1}{t - T_1} \sum_{s=T_1+1}^t \begin{pmatrix} \mathbb{I}(a_s^* = 1) \mathbf{v}_s \\ \vdots \\ \mathbb{I}(a_s^* = M) \mathbf{v}_s \end{pmatrix} \mathbf{z}_s^\top, \quad (10)$$

and use it in place of  $\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}$  in (8). The calculation of  $\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}^*$  is again standard because it is the average of samples from i.i.d. random variables.

However, because the optimal arms are unknown, the analysis of  $\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}$  poses a challenge. First, since  $\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}$  depends on the choice of arm  $a_s$  for all  $s$  from  $T_1 + 1$  to  $t$ , any past coefficient estimate  $\widehat{\boldsymbol{\alpha}}_s$ , through its effect on the selected arm, has an effect on  $\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}$ . Second, these coefficient estimates are serially correlated because they are all calculated from the same data source. Lastly, not only does  $\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}$  depend on  $\{\widehat{\boldsymbol{\alpha}}_s : T_1 \leq s \leq t - 1\}$ ,  $\widehat{\boldsymbol{\alpha}}_t$  also depends on  $\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}$ . This two-way dependence implies that the coefficient estimates and the covariance estimates are entangled. The past coefficient estimate  $\widehat{\boldsymbol{\alpha}}_s$  affects the choice of arm  $a_s$ , which affects  $\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}$  for all  $t > s$ . These covariance estimates, through expression (8), again affects  $\widehat{\boldsymbol{\alpha}}_t$ .

The two-way dependence, a common feature of online learning environments, makes it difficult for the coefficient estimates to converge to the true value in the long run. A poor coefficient estimate can lead to a wrong action, which contaminates the generated data (since the data associated with the right arm is no longer observed). The wrong data then worsens the future coefficient estimates and, thus, the future data-generating process. In summary, errors made in the past do not disappear; instead, they make future errors more likely.

To deal with this problem, we now describe an approach that untangles the dependence via a “zig-zag” induction. As an overview of its structure, assume that, in period  $t - 1$  in Phase 3, the coefficient estimate  $\widehat{\boldsymbol{\alpha}}_s$  is sufficiently close to the true value for *all*  $s$  from  $T_2 + 1$  to  $t - 1$  with a sufficiently high probability. The goal is to show that this remains true for period  $t$ , and this takes two steps. First, we show that the covariance estimate  $\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}$  is sufficiently accurate in period  $t$ . This is the zig-step. Second, we use this result from the zig-step to show that the induction hypothesis holds at time  $t$ . The difficult part of the induction is the zig-step. As mentioned earlier, the covariance estimate in the zig-step depends on all the coefficient estimates in the past, which are serially correlated. Our technique finds a means to replace the dependence (of the covariance estimates) on the past coefficient estimates with the dependence on the realizations of past covariates, which are i.i.d. over time.

**5.1.2. Details of the Zig-Zag Induction** We now provide a detailed description of the zig-zag induction. Our description comprises three parts. First, we define the relevant objects for the induction. Second, we describe the zig-zag decomposition within the induction step. Third, we describe the observation that allows us to deal with the difficulties in the zig-step of the induction.

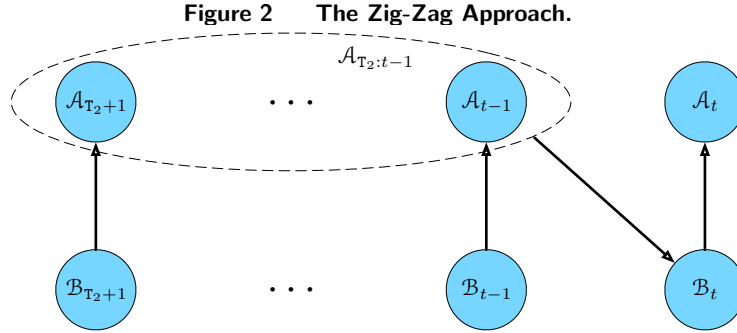
**Part I.** Let

$$\mathcal{A}_t := \left\{ \|\hat{\alpha}_t - \alpha\| \leq C_A \sqrt{\frac{\log(T)}{t - T_1}} \right\} \quad \text{and} \quad \mathcal{B}_t := \{ \|\hat{\Sigma}_{vz, T_1:t} - \Sigma_{vz}^*\| \leq \eta \},$$

for some constant  $C_A > 0$ . Further, let  $\mathcal{A}_{T_2:t} := \bigcap_{s=T_2+1}^t \mathcal{A}_s$ , where we note that the intersection begins at the beginning of Phase 2. This is the event in which all the past coefficient estimates from Phase 2 are sufficiently accurate. The goal of the induction is to show that there exists a constant  $K > 0$ , such that

$$\Pr(\mathcal{A}_{T_2:t}) \geq 1 - K(t - T_2)T^{-2}, \quad t = T_2 + 1, \dots, T.$$

This is achieved via a two-step induction using  $\mathcal{B}_t$  as a bridge—first from  $\mathcal{A}_{T_2:t-1}$  to  $\mathcal{B}_t$  and then from  $\mathcal{B}_t$  to  $\mathcal{A}_t$  (see Figure 2).



Note that we perform the induction on  $\mathcal{A}_{T_2:t}$  instead of  $\mathcal{A}_t$ . This is because the zig-zag induction step requires the analysis of the covariance estimate. As we emphasize above, the covariance estimate depends on the past coefficient estimates. Therefore, showing that the covariance estimate is sufficiently accurate requires event  $\mathcal{A}_{T_2:t}$  to occur—that is, all the past coefficient estimates from Phase 2 on to be sufficiently accurate.

**Part II.** We decompose the induction step into the zig-zag steps. In particular, suppose the induction claim on  $\mathcal{A}_{T_2:t-1}$  holds:

$$\Pr(\mathcal{A}_{T_2:t-1}) \geq 1 - K(t - 1 - T_2)T^{-2}.$$

Note that

$$\Pr(\mathcal{A}_{T_2:t}) = \Pr(\mathcal{A}_t \cap \mathcal{A}_{T_2:t-1}) = \Pr(\mathcal{A}_{T_2:t-1}) - \Pr(\mathcal{A}_t^c \cap \mathcal{A}_{T_2:t-1}).$$

Therefore, the induction step requires establishing an upper bound on  $\Pr(\mathcal{A}_t^c \cap \mathcal{A}_{T_2:t-1})$ . The zig-zag approach decomposes  $\mathcal{A}_t^c \cap \mathcal{A}_{T_2:t-1}$  into two events depending on whether the covariance estimate in period  $t$  is sufficiently accurate (i.e.,  $\mathcal{B}_t$  or  $\mathcal{B}_t^c$ ):

$$\begin{aligned} \Pr(\mathcal{A}_t^c \cap \mathcal{A}_{T_2:t-1}) &= \Pr(\mathcal{A}_t^c \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_t) + \Pr(\mathcal{A}_t^c \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_t^c) \\ &\leq \Pr(\mathcal{A}_t^c \cap \mathcal{B}_t) + \Pr(\mathcal{A}_{T_2:t-1} \cap \mathcal{B}_t^c). \end{aligned} \quad (11)$$

The analysis of  $\Pr(\mathcal{A}_t^c \cap \mathcal{B}_t)$ , which corresponds to the second step of the induction (from  $\mathcal{B}_t$  to  $\mathcal{A}_t$  in Figure 2), is standard. When  $\widehat{\Sigma}_{vz, T_1:t}$  is sufficiently accurate—that is, in the event of  $\mathcal{B}_t$ —the formula for  $\widehat{\alpha}_t$  in (8) can be used to show that  $\mathcal{A}_t^c$  occurs with a sufficiently small probability.

However, the analysis of  $\Pr(\mathcal{A}_{T_2:t-1} \cap \mathcal{B}_t^c)$ , which corresponds to the first step of the induction (from  $\mathcal{A}_{T_2:t-1}$  to  $\mathcal{B}_t$  in Figure 2) is nonstandard. Recall that if the optimal arms were known and selected in each period, then it would suffice to analyze  $\widehat{\Sigma}_{vz, T_1:t}^*$  defined in (10) instead of analyzing  $\widehat{\Sigma}_{vz, T_1:t}$ . The analysis of the former would be standard, since it would only involve realizations of the covariates and the IVs  $\{(\mathbf{v}_s, \mathbf{z}_s) : s = T_1 + 1, \dots, t\}$ , which are i.i.d. over time. In contrast, because the optimal arms are not known and the selected arm depends on the past coefficient estimates, the calculation of  $\widehat{\Sigma}_{vz, T_1:t}$  depends on all the past coefficient estimates, which are not i.i.d.

**Part III.** Finally, we make the following key observation. When the induction assumption holds in period  $t - 1$  of Phase 3, we can show (with details in the proof) that there exists a constant  $K' > 0$ , such that

$$\{a_s \neq a_s^*\} \subseteq \cup_{1 \leq j \neq a_s^* \leq M} \{|\mathbf{v}_s^\top \boldsymbol{\alpha}_j - \mathbf{v}_s^\top \boldsymbol{\alpha}_{a_s^*}| \leq K' \eta\}, \quad (12)$$

for all  $s = T_2 + 1, \dots, t - 1$ . Therefore, this relation states that a wrong arm is selected only when the difference in the actual expected reward between this arm and some other arm is small. (We can also show, with a similar argument, that (12) holds all  $s = T_1 + 1, \dots, T_2$ , i.e., Phase 2 of the algorithm, provided that the coefficient estimate at the end of Phase 1 is sufficiently accurate.)

The intuition for this observation is that, when the coefficient estimate is sufficiently accurate, the estimated expected reward of choosing each arm (calculated from the coefficient estimate) is also close to the actual expected reward; and thus the arm selected by the algorithm is optimal if the actual expected reward from this arm is sufficiently different than that from any other arm. In other words, when a wrong arm is selected, it must be that the difference in the actual expected reward between this arm and some other arm is small.

A consequence of this observation is that, instead of estimating the probability that a wrong arm is selected, which depends on the past coefficient estimates, it suffices to estimate the probability that the covariate vector  $\mathbf{v}_s$  falls within a bound. Specifically, we can show that

$$\|\widehat{\Sigma}_{vz, T_1:t} - \widehat{\Sigma}_{vz, T_1:t}^*\| \leq \tilde{K} \frac{1}{t - T_1} \sum_{s=T_1+1}^t \mathbb{I}(a_s \neq a_s^*)$$

$$\leq \tilde{K} \frac{1}{t - \mathbf{T}_1} \sum_{s=\mathbf{T}_1+1}^t \sum_{1 \leq j \neq a_s^* \leq M} \mathbb{I}(|\mathbf{v}_s^\top \boldsymbol{\alpha}_j - \mathbf{v}_s^\top \boldsymbol{\alpha}_{a_s^*}| \leq K'\eta), \quad (13)$$

for some constant  $\tilde{K} > 0$ . Because the covariate vectors are i.i.d. over time periods  $s$ , so are  $\mathbb{I}(|\mathbf{v}_s^\top \boldsymbol{\alpha}_j - \mathbf{v}_s^\top \boldsymbol{\alpha}_i| \leq K'\eta)$ . In other words, the right-hand-side of the inequality (13) is a sample average of i.i.d. random variables; and therefore, bounds on its tail probability can be established using standard concentration inequalities. This allows us to show that, with a sufficiently high probability,  $\hat{\boldsymbol{\Sigma}}_{vz, \mathbf{T}_1:t}$  is close to  $\hat{\boldsymbol{\Sigma}}_{vz, \mathbf{T}_1:t}^*$ . Moreover, since it is standard to show that  $\hat{\boldsymbol{\Sigma}}_{vz, \mathbf{T}_1:t}^*$  is close to the true value with a sufficiently high probability, we can therefore show that the probability of  $\mathcal{B}_t^c$  is small. This completes the analysis of  $\Pr(\mathcal{A}_{\mathbf{T}_2:t-1} \cap \mathcal{B}_t^c)$  in the decomposition (11), and together with the analysis of  $\Pr(\mathcal{A}_t^c \cap \mathcal{B}_t)$  discussed earlier, we complete the induction step.

The induction shows that, with a sufficiently high probability, the coefficient estimates remains sufficiently accurate over the entire Phase 3 of the algorithm. Standard calculation (Bastani, Bayati, and Khosravi 2021) can then be applied to show that the regret in Phase 3 is  $\mathcal{O}(\log(T))$ . Since the durations of Phases 1 and 2 are both  $\mathcal{O}(\log(T))$ , this shows that the total regret is  $\mathcal{O}(\log(T))$ . Goldenshluger and Zeevi (2013) prove that in a standard linear contextual bandit problem—with no endogeneity involved—the best possible lower bound for any algorithm is  $C \log(T)$  for some constant  $C$ . Since our setup allows endogeneity in the covariates and includes standard linear contextual bandits as special cases,  $C \log(T)$  is also a lower bound on regret for our setup. Hence, Theorem 1 shows that IV-Greedy achieves the asymptotically minimal regret.

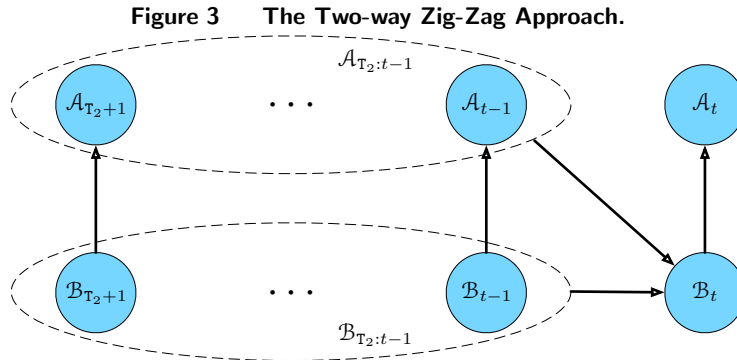
## 5.2. IV-UCB

**THEOREM 2.** *Let  $\mathbf{T}_1 = C_{\mathbf{T}_1} \log(T)$  and  $\mathbf{T}_2 = (C_{\mathbf{T}_1} + C_{\mathbf{T}_2}) \log(T)$  for some sufficiently large constants  $C_{\mathbf{T}_1}$  and  $C_{\mathbf{T}_2}$ . Then, under Assumption 1, the regret of the IV-UCB algorithm (i.e., Algorithm 1 with  $\theta > 0$ ) is  $\text{Regret}(T) = \mathcal{O}(\log^2(T))$  for all  $T > \mathbf{T}_2$ .*

For the IV-UCB algorithm, in each period  $t$  of Phase 3, the selected arm  $a_t$  depends not only on the coefficient estimate  $\hat{\boldsymbol{\alpha}}_{t-1}$ , as is in the case of IV-Greedy, but also on the covariance estimate  $\hat{\boldsymbol{\Sigma}}_{vz, \mathbf{T}_1:t-1}$  (again, through the calculation of  $\hat{\boldsymbol{\Omega}}_{t-1}$ ). This arm selection criterion deepens the two-way dependence between the coefficient estimates and the covariance estimates. In particular, whereas the covariance estimates only affect the arm choices indirectly (through their effects on the coefficient estimates) in IV-Greedy, they now affect the arm choices directly. This implies that, to establish the convergence of the coefficient estimates, both the past coefficient estimates and the past covariance estimates must be well controlled.

Despite the additional complexity, the zig-zag induction approach can be adapted to untangle the two-way dependence in Phase 3 of the algorithm. Different from the approach for IV-Greedy,

the induction step from period  $t - 1$  to period  $t$  will now also require the occurrence of event  $\mathcal{B}_{T_2:t-1} := \cap_{T_2+1}^{t-1} \mathcal{B}_s$  (i.e., all the past covariance estimates) to be sufficiently accurate, as opposed to only the occurrence of  $\mathcal{B}_{t-1}$ . Figure 3 illustrates the flow of the induction and indicates that in order to show the covariance is well estimated in period  $t$ , there is a requirement on both the coefficient and covariance estimates.



In contrast to IV-Greedy, the regret of IV-UCB is  $\mathcal{O}(\log^2(T))$ . The additional factor of  $\log(T)$  is incurred because—in addition to the dependence on the coefficient estimate—the direct dependence of the UCB criterion on the covariance estimates introduces an additional source of error in selecting the optimal arm. It takes more time to eliminate the additional error, thereby leading to a higher regret. Note that Zhong, Hong, and Liu (2021) also establish a  $\mathcal{O}(\log^2(T))$  regret of a UCB-type algorithm for standard linear contextual bandits. However, the result that IV-UCB has a higher regret than IV-Greedy must be interpreted in an asymptotic sense as the time horizon increases to infinity. The former may still have a lower regret for a given finite time horizon.

### 5.3. Statistical Inference

The contextual bandit literature mostly focuses on regret analysis, whereas results on statistical inference are relatively scarce. We provide such results for the coefficient estimates.

**THEOREM 3.** *Let  $T_1 = C_{T_1} \log(T)$  and  $T_2 = (C_{T_1} + C_{T_2}) \log(T)$  for some sufficiently large constants  $C_{T_1}$  and  $C_{T_2}$ . Then, under Assumption 1, the estimator  $\hat{\alpha}_T$  of Algorithm 1 satisfies*

$$\sqrt{T - T_1}(\hat{\alpha}_T - \alpha) \rightsquigarrow \mathcal{N}\left(\mathbf{0}, \sigma^2 (\Sigma_{vz}^* \Sigma_{zz}^{-1} (\Sigma_{vz}^*)^\top)^{-1}\right), \quad (14)$$

as  $T \rightarrow \infty$ , where  $\Sigma_{vz}^* = \mathbb{E}[\mathbf{v}_t^* \mathbf{z}_t^\top]$ ,  $\Sigma_{zz} = \mathbb{E}[\mathbf{z}_t \mathbf{z}_t^\top]$ , and  $\mathcal{N}(\mathbf{0}, \Sigma)$  denotes the multivariate normal distribution with zero mean vector and covariance matrix  $\Sigma$ .

To construct a confidence interval for  $\boldsymbol{\alpha}$ , we need a consistent estimator of the covariance matrix  $\sigma^2(\boldsymbol{\Sigma}_{vz}^* \boldsymbol{\Sigma}_{zz}^{-1} (\boldsymbol{\Sigma}_{vz}^*)^\top)^{-1}$  in (14). To this end, note that  $\hat{\sigma}_T^2$ ,  $\hat{\boldsymbol{\Sigma}}_{vz, T_1:T}$ , and  $\hat{\boldsymbol{\Sigma}}_{zz, T_1:T}$  are consistent estimators of  $\sigma^2$ ,  $\boldsymbol{\Sigma}_{vz}^*$ , and  $\boldsymbol{\Sigma}_{zz}$ , respectively. Moreover, by the definition of  $\hat{\boldsymbol{\Omega}}_T$  in Algorithm 1,

$$(T - T_1) \hat{\boldsymbol{\Omega}}_T = (\hat{\boldsymbol{\Sigma}}_{vz, T_1:T} \hat{\boldsymbol{\Sigma}}_{zz, T_1:T}^{-1} (\hat{\boldsymbol{\Sigma}}_{vz, T_1:T})^\top)^{-1}.$$

Thus, we use  $(T - T_1) \hat{\sigma}_T^2 \hat{\boldsymbol{\Omega}}_T$  to consistently estimate  $\sigma^2(\boldsymbol{\Sigma}_{vz}^* \boldsymbol{\Sigma}_{zz}^{-1} (\boldsymbol{\Sigma}_{vz}^*)^\top)^{-1}$ . Then, standard calculation leads to the following result, which facilitates inference on  $\boldsymbol{\alpha}$ .

**COROLLARY 1.** *Let  $T_1 = C_{T_1} \log(T)$  and  $T_2 = (C_{T_1} + C_{T_2}) \log(T)$  for some sufficiently large constants  $C_{T_1}$  and  $C_{T_2}$ . Then, under Assumption 1, the estimator  $\hat{\boldsymbol{\alpha}}_T$  of Algorithm 1 satisfies*

$$\begin{aligned} (T - T_1)^{\frac{1}{2}} (\hat{\sigma}_T^2 \hat{\boldsymbol{\Omega}}_T)^{-\frac{1}{2}} (\hat{\boldsymbol{\alpha}}_T - \boldsymbol{\alpha}) &\rightsquigarrow \mathcal{N}(0, \mathbf{I}_{Mp}), \\ (T - T_1) (\hat{\boldsymbol{\alpha}}_T - \boldsymbol{\alpha})^\top (\hat{\sigma}_T^2 \hat{\boldsymbol{\Omega}}_T)^{-1} (\hat{\boldsymbol{\alpha}}_T - \boldsymbol{\alpha}) &\rightsquigarrow \chi_{Mp}^2. \end{aligned}$$

as  $T \rightarrow \infty$ , where  $\hat{\sigma}_T$  and  $\hat{\boldsymbol{\Omega}}_T$  are defined in Algorithm 1,  $\mathbf{I}_{Mp}$  is the identity matrix of size  $(Mp) \times (Mp)$ , and  $\chi_{Mp}^2$  denotes the chi-square distribution with  $Mp$  degrees of freedom.

## 6. Simulation Experiments

In this section, we present a simulation study to demonstrate the performances of IV-Greedy and IV-UCB and compare them with the following alternatives. These alternatives—similar to both IV-Greedy and IV-UCB—all have a randomization phase of duration  $T_1$ , in which each arm is selected uniformly at random.

- (i) Naïve-IV-UCB. In each period  $t = T_1 + 1, \dots, T$ , this algorithm performs arm-specific 2SLS to obtain an estimate of  $\boldsymbol{\alpha}$  and selects an arm following the same UCB criterion as the IV-UCB algorithm, except that  $\hat{\sigma}_t$  and  $\hat{\boldsymbol{\Omega}}_{i,t}$  are defined and computed from the arm-specific 2SLS.
- (ii) OLS-UCB. In each period  $t = T_1 + 1, \dots, T$ , this algorithm performs arm-specific OLS to obtain an estimate of  $\boldsymbol{\alpha}$  and selects an arm following the same UCB criterion as the IV-UCB algorithm, except that  $\hat{\sigma}_t$  and  $\hat{\boldsymbol{\Omega}}_{i,t}$  are defined and computed from the arm-specific OLS.
- (iii) Randomize-then-commit (RTC). After the first  $T_1$  periods, this algorithm performs arm-specific 2SLS to obtain an estimate of  $\boldsymbol{\alpha}$  and commit to it. In each period  $t = T_1 + 1, \dots, T$ , the algorithm continues to use this estimate without updates to select an arm in a greedy manner. Note that IV-UCB, Naïve-IV-UCB, and OLS-UCB all involve a parameter  $\theta$  in their respective UCB criterion, and we set  $\theta = 0.5$  for them all. For all the five algorithms, we set  $T_1 = 50$  and the total run length  $T = 20000$ . In addition, we set  $T_2 = 100$  for both IV-Greedy and IV-UCB.

We evaluate these algorithms via regret and number of wrong arms pulled, both of which are common metrics for online learning algorithms. We also perform statistical inference on the coefficients in the reward function upon termination of the algorithms.

To this end, we consider a two-armed linear contextual bandit problem with three-dimensional covariates, i.e.,  $M = 2$ . Let the covariate vector  $\mathbf{v}_t = (1, x_t, d_t)^\top$ , where  $x_t \in \mathbb{R}$  is an exogenous variable and  $d_t \in \mathbb{R}$  is an endogenous variable. The endogeneity is introduced in the following manner. Suppose the random reward at time  $t$  is

$$R_t = \sum_{i=1}^2 \mu_i(\mathbf{v}_t) \mathbb{I}(a_t = i) + \epsilon_t,$$

where  $\mu_i(\mathbf{v}_t) = \beta_{i,0} + \beta_{i,1}x_t + \gamma_i d_t$ , for some parameters  $\beta_{i,0}$ ,  $\beta_{i,1}$ , and  $\gamma_i$ . Suppose also

$$d_t = \sqrt{x_t} + \rho_{\tilde{z}} \tilde{z}_t + \rho_{\eta} \eta_t \quad \text{and} \quad \epsilon_t = \tilde{\epsilon}_t + 2\eta_t,$$

where  $\rho_{\tilde{z}}$  and  $\rho_{\eta}$  are some positive constants, and  $x_t$ ,  $\tilde{z}_t$ ,  $\eta_t$ , and  $\tilde{\epsilon}_t$  are independent random variables with the following distributions:

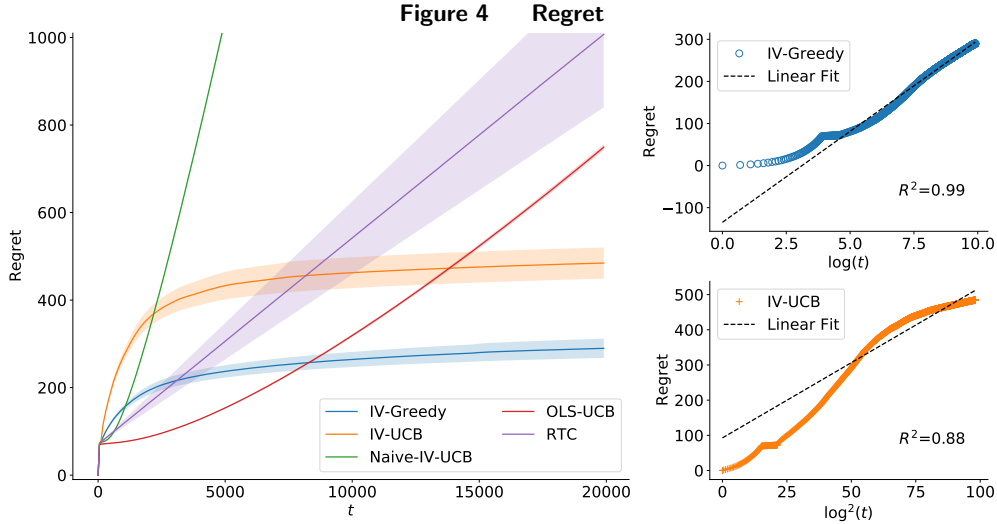
$$\begin{aligned} x_t &\sim \text{TruncNorm}(\mu_x, \sigma_x^2, l_x, u_x), & \tilde{z}_t &\sim \text{TruncNorm}(\mu_{\tilde{z}}, \sigma_{\tilde{z}}^2, l_{\tilde{z}}, u_{\tilde{z}}), \\ \eta_t &\sim \text{TruncNorm}(\mu_{\eta}, \sigma_{\eta}^2, l_{\eta}, u_{\eta}), & \tilde{\epsilon}_t &\sim \mathcal{N}(0, \sigma_{\tilde{\epsilon}}^2). \end{aligned}$$

Here, we use  $\text{TruncNorm}(\mu, \sigma^2, l, u)$  to denote the truncated normal distribution that is derived from bounding the  $\mathcal{N}(\mu, \sigma^2)$  distribution on the interval  $(l, u)$ . The parameters are specified as follows:

$$\begin{aligned} (\beta_{1,0}, \beta_{1,1}, \gamma_1) &= (1, 4, 4), & (\beta_{2,0}, \beta_{2,1}, \gamma_2) &= (8, 2, 2), \\ \rho_{\tilde{z}} &= 0.5, & \rho_{\eta} &= 1.5, & \sigma_{\tilde{\epsilon}}^2 &= 0.25, \\ (\mu_x, \sigma_x^2, l_x, u_x) &= (0, 1, 0, 10), & (\mu_{\tilde{z}}, \sigma_{\tilde{z}}^2, l_{\tilde{z}}, u_{\tilde{z}}) &= (0, 4, 0, 10), \\ (\mu_{\eta}, \sigma_{\eta}^2, l_{\eta}, u_{\eta}) &= (0, 0.25, -5, 5). \end{aligned}$$

Note that  $\epsilon_t$  and  $d_t$  are correlated through  $\eta_t$  and, thus,  $d_t$  becomes endogenous. Moreover,  $\tilde{z}_t$  is correlated with  $d_t$  but is uncorrelated with  $\epsilon_t$ , so  $\tilde{z}_t$  serves as an IV for  $d_t$ . However, for IV-Greedy and IV-UCB to work, the number of IVs must be at least twice as many as the number of endogenous variables in order to account for the endogeneity spillover issue; see Comment 2. As discussed in Comment 3, we create additional IVs using  $\mathbb{I}(x_t \geq \hat{x})$ ,  $\mathbb{I}(\tilde{z}_t \geq \hat{z})$ , and the interaction terms between them and  $x_t$  or  $\tilde{z}_t$ , for some threshold values  $\hat{x}$  and  $\hat{z}$ . Specifically, we use the following list of IVs:

$$\{1, x_t, \tilde{z}_t, \mathbb{I}(x_t \geq 1), \mathbb{I}(x_t \geq 1)\tilde{z}_t, \mathbb{I}(x_t \geq 1.5), \mathbb{I}(x_t \geq 1.5)\tilde{z}_t, \mathbb{I}(\tilde{z}_t \geq 2), \mathbb{I}(\tilde{z}_t \geq 2)\tilde{z}_t\}.$$

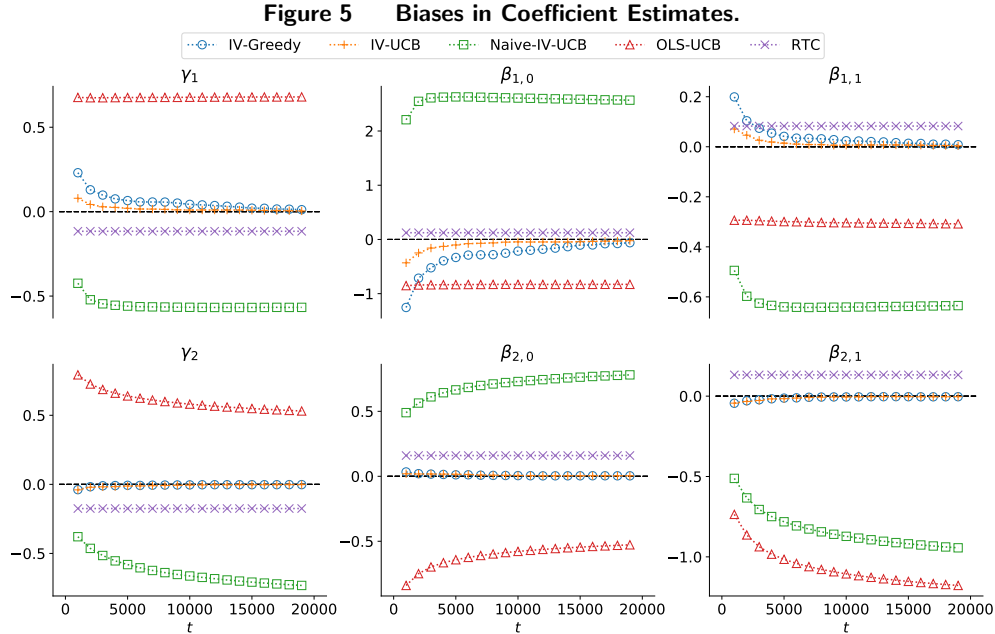


*Note.* The shaded areas represent 95% confidence bands calculated from 1000 replications.

Further, we run the five algorithms (IV-Greedy, IV-UCB, Naïve-IV-UCB, OLS-UCB, and RTC) 1000 times each and present the results in Figures 4 and 5 as well as Table 1.

Figure 4 presents the regret of each algorithm. First, in contrast to the others, both IV-Greedy and IV-UCB achieve a regret with a diminishing growth rate after a sufficient number of iterations (e.g.,  $T = 5000$ ). The confidence bands further show that the performances of these two algorithms are stable in the learning process. Second, to verify our theoretical results in Theorems 1 and 2, we conduct a linear regression analysis of the regrets of IV-Greedy and IV-UCB on  $\log(t)$  and  $\log^2(t)$ , respectively. The value of  $R^2$  is 0.99 for IV-Greedy and 0.88 for IV-UCB, thereby providing strong empirical evidence for our theory. Third, although the theoretical analysis requires that Phases 1 and 2 of IV-Greedy and IV-UCB should be “sufficiently long,” the simulation results indicate that they can be set to be moderate values ( $T_1 = 50$  and  $T_2 = 100$ ) in practice.

Figure 5 presents the biases of the coefficient estimates. For the estimates produced by either IV-Greedy or IV-UCB, the bias vanishes as the number of iterations increases, while the other algorithms are generally biased. In particular, the fact that Naïve-IV-UCB produces inconsistent estimates of the coefficients demonstrates the existence of endogeneity spillover in this online learning example (the endogeneity in the covariates is addressed by the use of IVs in the arm-specific 2SLS, but the endogeneity that spills over to the actions remains). Recall that RTC first performs arm-specific 2SLS to data from randomization and then commits to the coefficient estimates. This results in a dilemma: on the one hand, if the randomization period is long, then the regret is large, because the randomization produces a linear regret; on the other hand, if the randomization period



*Note.* For each parameter  $\alpha$ , the bias  $\hat{\alpha}_t - \alpha$  is calculated for 1000 replications to obtain an average.

is short, then the coefficient estimates are inaccurate,<sup>4</sup> which also leads to a large regret in the long run. Thus, the regret of RTC is, at best, a fractal polynomial of  $T$ , which is significantly worse than  $\log(T)$ .

In Table 1, we report main statistics for the coefficient estimates after 20000 iterations of the five algorithms. For all the coefficients, both IV-Greedy and IV-UCB result in almost zero bias and the corresponding confidence intervals achieve proper coverage of the true values (the coverage is close to 95% for 95% confidence intervals). Note that Naïve-IV-UCB produces both substantially biased estimates and confidence intervals with zero coverage, which, again, demonstrates the endogeneity spillover issue. Note also that the quality of the inference results produced by RTC is fair, but the SD is much larger than both IV-Greedy and IV-UCB.

## 7. Conclusions

We study the dynamic selection problems that arise in algorithmic decision-making when the data has endogeneity problems. In a contextual bandit model, endogeneity in data influences the arms pulled, causing a non-random sampling of data that results in self-fulfilling bias. We correct for the bias by incorporating IVs to existing online learning algorithms. The new algorithms lead to the true parameter values and meanwhile attain low (logarithmic-like) regret levels. We also prove

<sup>4</sup> In general, the 2SLS estimates have a second-order bias if the sample size is small or the IVs are weak (Andrews, Stock, and Sun 2019).

**Table 1** Statistics at Time  $T = 20000$ .

	$\gamma_1$			$\beta_{1,0}$			$\beta_{1,1}$		
	Bias	SD	Coverage	Bias	SD	Coverage	Bias	SD	Coverage
IV-Greedy	0.009	0.112	0.961	-0.049	0.527	0.960	0.007	0.092	0.960
IV-UCB	0.004	0.079	0.958	-0.020	0.352	0.960	0.003	0.080	0.959
Naïve-IV-UCB	-0.566	0.028	0.000	2.565	0.075	0.000	-0.634	0.025	0.000
OLS-UCB	0.680	0.016	0.000	-0.828	0.059	0.000	-0.308	0.021	0.000
RTC	-0.116	1.177	0.913	0.124	1.218	0.921	0.083	1.002	0.937

	$\gamma_2$			$\beta_{2,0}$			$\beta_{2,1}$		
	Bias	SD	Coverage	Bias	SD	Coverage	Bias	SD	Coverage
IV-Greedy	-0.002	0.032	0.955	0.004	0.041	0.955	-0.004	0.037	0.959
IV-UCB	-0.002	0.034	0.954	0.004	0.043	0.956	-0.004	0.039	0.959
Naïve-IV-UCB	-0.736	0.064	0.000	0.784	0.057	0.000	-0.949	0.051	0.000
OLS-UCB	0.527	0.030	0.000	-0.524	0.029	0.000	-1.183	0.050	0.000
RTC	-0.174	1.543	0.928	0.160	1.757	0.933	0.132	1.299	0.942

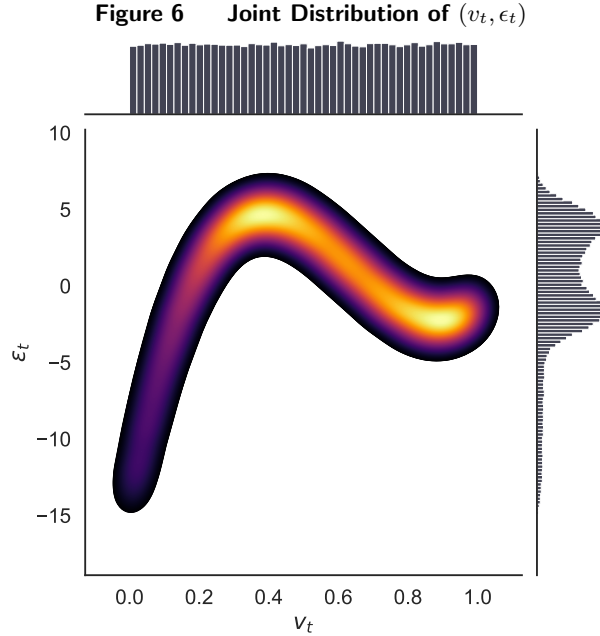
*Note.* The coverage is calculated for the 95% confidence interval based on 1000 replications.

a central limit theorem for statistical inference of the parameters of interest. To establish these properties, we develop a general technique that untangles the interdependence between data and actions.

The use of IVs corrects the bias because they perturb the data-generating process. Therefore, our approach features “ex-ante randomization”, which contrasts with the “ex-post randomization”—that is, randomizing the choices given to the users once they arrive. The ex-post randomization, typically under the name of exploration, is the focus of most research. However, in the presence of endogenous data, our analysis highlights the importance of ex-ante randomization. Understanding how best to carry out ex-ante randomization and how these two types of randomization interact offer fruitful venues for future research.

## Appendix. Example of Multiplicity of the Self-fulfilling Bias

Consider the example in Section 3: the expected reward for the two arms are  $\mu_1(v) \equiv c$  and  $\mu_2(v) = \alpha v$ . Let  $v_t$  be a uniform random variable on  $[0, 1]$ ,  $\eta_t$  be an independent standard normal random variable, and  $\epsilon_t = \sum_{k=0}^3 \beta_k v_t^k + \eta_t$  for some constants  $\beta_k \in \mathbb{R}$ ,  $k = 0, 1, 2, 3$ . See Figure 6 for the joint distribution of  $(v_t, \epsilon_t)$ .



Note.  $c = 1$ ,  $\alpha = 15$ , and  $\{\beta_k : k = 0, 1, 2, 3\}$  are determined by equations (18) and (19) with  $(r_1, r_2, r_3) = (2, 5, 10)$ .

Direct calculations reveals that for any  $b \in [0, 1]$ ,

$$\begin{aligned}\text{Var}[v_t | v_t > b] &= \frac{1}{12}(1-b)^2, \\ \text{Cov}[v_t, v_t^2 | v_t > b] &= \frac{1}{12}(1-b)^2(1+b), \\ \text{Cov}[v_t, v_t^3 | v_t > b] &= \frac{1}{40}(1-b)^2(3b^2 + 4b + 3).\end{aligned}$$

Consequently,

$$\frac{\text{Cov}[v_t, \epsilon_t | v_t > b]}{\text{Var}[v_t | v_t > b]} = \frac{\sum_{k=1}^3 \beta_k \text{Cov}[v_t, v_t^k | v_t > b]}{\text{Var}[v_t | v_t > b]} = \beta_1 + \beta_2(1+b) + \frac{3}{10}\beta_3(3b^2 + 4b + 3). \quad (15)$$

By setting  $b = 0$  in equation (15), we derive the limit of the OLS estimate of  $\alpha$  in the following manner:

$$\hat{\alpha}_{\text{OLS}} := \alpha + \frac{\text{Cov}[v_t, \epsilon_t]}{\text{Var}[v_t]} = \alpha + \frac{\text{Cov}[v_t, \epsilon_t | v_t > 0]}{\text{Var}[v_t | v_t > 0]} = \alpha + \beta_1 + \beta_2 + \frac{9}{10}\beta_3.$$

Moreover, applying equation (15) to equation (4) with  $b = c/\hat{\alpha}$  yields that

$$\hat{\alpha}^3 - \left(\alpha + \beta_1 + \beta_2 + \frac{9}{10}\beta_3\right)\hat{\alpha}^2 - c\left(\beta_2 + \frac{6}{5}\beta_3\right)\hat{\alpha} - \frac{9}{10}c^2\beta_3 = 0, \quad (16)$$

which is a cubic equation in terms of  $\hat{\alpha}$ .

Assume that equation (16) has three real roots:  $r_k$ ,  $k = 1, 2, 3$ . Then,  $\prod_{k=1}^3 (\hat{\alpha} - r_k) = 0$ —that is,

$$\hat{\alpha}^3 - (r_1 + r_2 + r_3)\hat{\alpha}^2 + (r_1r_2 + r_1r_3 + r_2r_3)\hat{\alpha} - r_1r_2r_3 = 0. \quad (17)$$

By matching the coefficients of equation (16) and equation (17), we have

$$\beta_3 = \frac{10r_1r_2r_3}{9c^2}, \quad \beta_2 = \frac{-(r_1r_2 + r_1r_3 + r_2r_3)}{c} - \frac{6}{5}\beta_3, \quad \text{and} \quad \beta_1 = r_1 + r_2 + r_3 - \left(\alpha + \beta_2 + \frac{9}{10}\beta_3\right). \quad (18)$$

Further, we may set  $\beta_0$  such that  $\mathbb{E}[\epsilon_t] = \sum_{k=0}^3 \beta_k v_t^k = 0$ —that is,

$$\beta_0 = - \sum_{k=1}^3 \beta_k \mathbb{E}[v_t^k] = - \sum_{k=1}^3 \frac{1}{k+1} \beta_k. \quad (19)$$

Therefore, for any arbitrary  $(r_1, r_2, r_3)$  such that  $c/r_k \in (0, 1)$ ,  $k = 1, 2, 3$ , there exist  $(\beta_0, \beta_1, \beta_2, \beta_3)$ —given by equations (18) and (19)—such that  $r_1, r_2$ , and  $r_3$  are three real fixed points of equation (4). Note that our method of constructing multiple self-fulfilling biases can be generalized. For an arbitrary number of self-fulfilling biases with arbitrary values, we can again let the covariates be uniformly distributed and construct the noise as a higher-order polynomial of the covariates. The coefficients of the polynomial can again be solved via linear equations.

## References

- Agrawal S, Goyal N (2013) Thompson sampling for contextual bandits with linear payoffs. *Proceedings of the 30th Annual International Conference on Machine Learning*, 1220–1228.
- Altonji JG, Elder TE, Taber CR (2005) Selection on observed and unobserved variables: Assessing the effectiveness of Catholic schools. *Journal of Political Economy* 113:151–184.
- Andrews I, Stock J, Sun L (2019) Weak instruments in IV regression: Theory and practice. *Annual Review of Economics* 11:727–753.
- Athey S, Imbens G (2017) The econometrics of randomized experiments. Banerjee AV, Duflo E, eds., *Handbook of Field Experiments*, volume 1, 73–140 (North-Holland).
- Athey S, Wager S (2021) Policy learning with observational data. *Econometrica* 89(1):133–161.
- Azevedo EM, Deng A, Montiel Olea JL, Rao J, Weyl EG (2020) A/B testing with fat tails. *Journal of Political Economy* 128(12):4614–4672.
- Bastani H, Bayati M, Khosravi K (2021) Mostly exploration-free algorithms for contextual bandits. *Management Science* 67(3):1329–1349.
- Belloni A, Chen D, Chernozhukov V, Hansen C (2012) Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica* 80(6):2369–2429.
- Bergemann D, Välimäki J (2008) Bandit problems. Durlauf SN, Blume LE, eds., *The New Palgrave Dictionary of Economics* (Basingstoke, UK: Palgrave Macmillan).
- Blake T, Nosko C, Tadelis S (2015) Consumer heterogeneity and paid search effectiveness: A large-scale field experiment. *Econometrica* 83(1):155–174.
- Blattner L, Nelson S, Spiess J (2021) Unpacking the black box: Regulating algorithmic decisions. Preprint available at *arXiv:2110.03443*.
- Caria S, Gordon G, Kasy M, Quinn S, Shami S, Teytelboym A (2020) An adaptive targeted field experiment: Job search assistance for refugees in Jordan. Preprint available at *SSRN:3689456*.
- Chen X, Pouzo D (2012) Estimation of nonparametric conditional moment models with possibly nonsmooth generalized residuals. *Econometrica* 80(1):277–321.

- Chu W, Li L, Reyzin L, Schapire R (2011) Contextual bandits with linear payoff functions. *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics*, 208–214.
- Cosslett SR (2004) Efficient semiparametric estimation of censored and truncated regressions via a smoothed self-consistency equation. *Econometrica* 72(4):1277–1293.
- Currie JM, MacLeod WB (2020) Understanding doctor decision making: The case of depression treatment. *Econometrica* 88(3):847–878.
- Dimakopoulou M, Zhou Z, Athey S, Imbens G (2019) Balanced linear contextual bandits. *Proceedings of the 33rd AAAI Conference on Artificial Intelligence*, 3445–3453.
- Donald SG, Imbens GW, Newey WK (2009) Choosing instrumental variables in conditional moment restriction models. *Journal of Econometrics* 152(1):28–36.
- Gibbons R, Henderson R (2012) What do managers do?: Exploring persistent performance differences among seemingly similar enterprises. Gibbons R, Roberts J, eds., *The Handbook of Organizational Economics*, 680–731 (Princeton University Press).
- Goldenshluger A, Zeevi A (2013) A linear response bandit problem. *Stochastic Systems* 3(1):230–261.
- Heckman JJ (1979) Sample selection bias as a specification error. *Econometrica* 47(1):153–161.
- Heckman JJ (1990) Varieties of selection bias. *American Economic Review: Papers and Proceedings* 80(2):313–318.
- Imbens GW, Rubin DB (2015) *Causal Inference for Statistics, Social, and Biomedical Sciences* (Cambridge University Press).
- Kallus N, Zhou A (2021) Minimax-optimal policy learning under unobserved confounding. *Management Science* 67(5):2870–2890.
- Kasy M, Sautmann A (2021) Adaptive treatment assignment in experiments for policy choice. *Econometrica* 89(1):113–132.
- Kawaguchi K (2021) When will workers follow an algorithm? A field experiment with a retail business. *Management Science* 67(3):1670–1695.
- Kawaguchi K, Uetake K, Watanabe Y (2021) Designing context-based marketing: Product recommendations under time pressure. *Management Science*, forthcoming.
- Kitagawa T, Tetenov A (2018) Who should be treated? Empirical welfare maximization methods for treatment choice. *Econometrica* 86(2):591–616.
- Kohavi R, Tang D, Xu Y (2020) *Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing* (Cambridge University Press).
- Koning R, Hasan S, Chatterji A (2019) Experimentation and startup performance: Evidence from A/B testing. Preprint available at *SSRN:3440291*.
- Li D, Raymond L, Bergman P (2020) Hiring as exploration. Preprint available at *SSRN:3630630*.
- Li J, Luo Y, Zhang X (2021) Causal reinforcement learning: An instrumental variable approach. Preprint available at *SSRN:3792824*.

- March JG (1991) Exploration and exploitation in organizational learning. *Organization Science* 2(1):71–87.
- Nambiar M, Simchi-Levi D, Wang H (2019) Dynamic learning and pricing with model misspecification. *Management Science* 65(11):4980–5000.
- Narita Y, Yasui S, Yata K (2019) Efficient counterfactual learning from bandit feedback. *Proceedings of the 33th AAAI Conference on Artificial Intelligence*, 4634–4641 (AAAI Press).
- Narita Y, Yasui S, Yata K (2021) Debiased off-policy evaluation for recommendation systems. *Fifteenth ACM Conference on Recommender Systems (RecSys '21)*, 372–379.
- Oster E (2019) Unobservable selection and coefficient stability: Theory and evidence. *Journal of Business & Economic Statistics* 37(2):187–204.
- Perchet V, Rigollet P, Chassang S, Snowberg E (2016) Batched bandit problems. *Annals of Statistics* 44(2):660–681.
- Rosenbaum PR, Rubin DB (1983) The central role of the propensity score in observational studies for causal effects. *Biometrika* 70(1):41–55.
- Shi C, Wan R, Chernozhukov V, Song R (2021) Deeply-debiased off-policy interval estimation. *Proceedings of the 38th International Conference on Machine Learning*, 9580–9591.
- Slivkins A (2019) Introduction to multi-armed bandits. *Foundations and Trends® in Machine Learning* 12(1-2):1–286.
- Tsybakov AB (2004) Optimal aggregation of classifiers in statistical learning. *Annals of Statistics* 32(1):135–166.
- Zhan R, Hadad V, Hirshberg DA, Athey S (2021) Off-policy evaluation via adaptive weighting with data from contextual bandits. *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2125–2135.
- Zhong Y, Hong LJ, Liu G (2021) Earning and learning with varying cost. *Production and Operations Management*, forthcoming.

## Supplemental Material

### EC.1. Concentration Inequalities

LEMMA EC.1. For any matrix  $\mathbf{M} \in \mathbb{R}^{k \times \ell}$ ,

$$\phi_{\max}(\mathbf{M}) = \sup_{\mathbf{u} \in \mathbb{R}^{\ell}: \|\mathbf{u}\|=1} \|\mathbf{M}\mathbf{u}\| \quad \text{and} \quad \phi_{\min}(\mathbf{M}) = \inf_{\mathbf{u} \in \mathbb{R}^{\ell}: \|\mathbf{u}\|=1} \|\mathbf{M}\mathbf{u}\|,$$

where  $\|\mathbf{u}\|$  denotes the Euclidean norm of  $\mathbf{u}$ .

*Proof.* See page 78 of Golub and Van Loan (2013).  $\square$

LEMMA EC.2. For any matrices  $\mathbf{M} \in \mathbb{R}^{k \times \ell}$  and  $\mathbf{N} \in \mathbb{R}^{\ell \times m}$ , we have

$$\phi_{\max}(\mathbf{MN}) \leq \phi_{\max}(\mathbf{M})\phi_{\max}(\mathbf{N}) \quad \text{and} \quad \phi_{\min}(\mathbf{MN}) \geq \phi_{\min}(\mathbf{M})\phi_{\min}(\mathbf{N}). \quad (\text{EC.1.1})$$

If  $\mathbf{M}$  is an invertible square matrix, then

$$\phi_{\max}(\mathbf{M}^{-1}) = (\phi_{\min}(\mathbf{M}))^{-1} \quad \text{and} \quad \phi_{\min}(\mathbf{M}^{-1}) = (\phi_{\max}(\mathbf{M}))^{-1}. \quad (\text{EC.1.2})$$

*Proof.* The inequalities in (EC.1.1) follow from Theorem 3 of Wang and Xi (1997), and the identities in (EC.1.2) follow from the definition of singular values.  $\square$

LEMMA EC.3. For any matrix  $\mathbf{M} \in \mathbb{R}^{k \times \ell}$ ,  $\|\mathbf{M}\| \leq \sqrt{k\ell}\|\mathbf{M}\|_{\max}$ .

*Proof.* See page 72 of Golub and Van Loan (2013).  $\square$

LEMMA EC.4. Let  $\{\mathbf{M}_i\}_{i=1}^N$  be a sequence of zero-mean independent  $k$ -by- $\ell$  random matrices. For each  $i = 1, \dots, N$ , assume that  $\|\mathbf{M}_i\| \leq b$  almost surely for some constant  $b > 0$ . Then, for all  $\tau \geq 0$ ,

$$\Pr\left(\left\|\frac{1}{N}\sum_{i=1}^N \mathbf{M}_i\right\| \geq \tau\right) \leq 2(k+\ell) \exp\left(\frac{-N\tau^2}{2b^2}\right),$$

Moreover, if  $k = \ell$ , then the leading constant  $2(k + \ell)$  on the right-hand-side of the above inequalities can be reduced to  $2k$ .

*Proof.* This is Hoeffding's inequality for bounded random matrices. See pages 174–175 of Wainwright (2019).  $\square$

COROLLARY EC.1. Let  $\{\mathbf{M}_i\}_{i=1}^N$  be a sequence of independent  $k$ -by- $\ell$  random matrices. For each  $i = 1, \dots, N$ , assume that  $\|\mathbf{M}_i\| \leq b$  almost surely for some constant  $b > 0$ . Then, for all  $\tau \geq 0$ ,

$$\Pr\left(\phi_{\max}\left(\frac{1}{N}\sum_{i=1}^N \mathbf{M}_i\right) \geq \phi_{\max}(\mathbb{E}[\mathbf{M}_1]) + \tau\right) \leq 2(k+\ell) \exp\left(\frac{-N\tau^2}{8b^2}\right), \quad (\text{EC.1.3})$$

$$\Pr\left(\phi_{\min}\left(\frac{1}{N}\sum_{i=1}^N\mathbf{M}_i\right)\leq\phi_{\min}(\mathbb{E}[\mathbf{M}_1])-\tau\right)\leq 2(k+\ell)\exp\left(\frac{-N\tau^2}{8b^2}\right). \quad (\text{EC.1.4})$$

Moreover, if  $k=\ell$ , then the leading constant  $2(k+\ell)$  on the right-hand-side of the above inequalities can be reduced to  $2k$ .

*Proof.* Let  $\widetilde{\mathbf{M}}_i = \mathbf{M}_i - \mathbb{E}[\mathbf{M}_i]$ . Then,  $\|\widetilde{\mathbf{M}}_i\| \leq \|\mathbf{M}_i\| + \|\mathbb{E}[\mathbf{M}_i]\| \leq 2b$ . By Lemmas EC.1 and EC.4,

$$\phi_{\max}\left(\frac{1}{N}\sum_{i=1}^N\widetilde{\mathbf{M}}_i\right) = \left\|\frac{1}{N}\sum_{i=1}^N\widetilde{\mathbf{M}}_i\right\| \leq \tau \quad (\text{EC.1.5})$$

holds with probability at least  $1 - 2(k+\ell)\exp\left(-\frac{N\tau^2}{8b^2}\right)$ .

Given the event (EC.1.5) holds, it follows from Lemma EC.1 that

$$\begin{aligned} \phi_{\max}\left(\frac{1}{N}\sum_{i=1}^N\mathbf{M}_i\right) &= \sup_{\mathbf{u}\in\mathbb{R}^\ell:\|\mathbf{u}\|=1}\left\|\frac{1}{N}\sum_{i=1}^N\mathbf{M}_i\mathbf{u}\right\| \\ &\leq \sup_{\mathbf{u}\in\mathbb{R}^\ell:\|\mathbf{u}\|=1}\|\mathbb{E}[\mathbf{M}_1]\mathbf{u}\| + \sup_{\mathbf{u}\in\mathbb{R}^\ell:\|\mathbf{u}\|=1}\left\|\frac{1}{N}\sum_{i=1}^N\widetilde{\mathbf{M}}_i\mathbf{u}\right\| \\ &\leq \phi_{\max}(\mathbb{E}[\mathbf{M}_1]) + \tau \end{aligned}$$

holds with probability at least  $1 - 2(k+\ell)\exp\left(-\frac{N\tau^2}{8b^2}\right)$ , proving (EC.1.3). Likewise, given the event (EC.1.5) holds,

$$\begin{aligned} \phi_{\min}\left(\frac{1}{N}\sum_{i=1}^N\mathbf{M}_i\right) &= \inf_{\mathbf{u}\in\mathbb{R}^\ell:\|\mathbf{u}\|=1}\left\|\frac{1}{N}\sum_{i=1}^N\mathbf{M}_i\mathbf{u}\right\| \\ &\geq \inf_{\mathbf{u}\in\mathbb{R}^\ell:\|\mathbf{u}\|=1}\|\mathbb{E}[\mathbf{M}_1]\mathbf{u}\| + \sup_{\mathbf{u}\in\mathbb{R}^\ell:\|\mathbf{u}\|=1}\left\|\frac{1}{N}\sum_{i=1}^N\widetilde{\mathbf{M}}_i\mathbf{u}\right\| \\ &\geq \phi_{\min}(\mathbb{E}[\mathbf{M}_1]) - \tau, \end{aligned}$$

which proves (EC.1.4).  $\square$

**LEMMA EC.5.** *Let  $\{X_i\}_{i=1}^N$  be a sequence of i.i.d. zero-mean sub-Gaussian random variables with variance proxy  $\varsigma^2$ . Then, for all  $\tau \geq 0$ ,*

$$\Pr\left(\left|\frac{1}{N}\sum_{i=1}^N X_i\right| \geq \tau\right) \leq 2\exp\left(-\frac{N\tau^2}{2\varsigma^2}\right).$$

*Proof.* This is Hoeffding's inequality for sub-Gaussian random variables. See page 24 of Wainwright (2019).  $\square$

**COROLLARY EC.2.** *Let  $\mathbf{g} \in \mathbb{R}^k$  be a vector of zero-mean sub-Gaussian random variables with variance proxy  $\varsigma^2$ . Let  $\{\mathbf{g}_i\}_{i=1}^N$  be a sequence of i.i.d. random vectors having the same distribution as  $\mathbf{g}$ . Then, for all  $\tau \geq 0$ ,*

$$\Pr\left(\left\|\frac{1}{N}\sum_{i=1}^N\mathbf{g}_i\right\| \geq \tau\right) \leq 2k\exp\left(-\frac{N\tau^2}{2k^2\varsigma^2}\right).$$

*Proof.* Let  $g_{ij}$  denote the  $j$ -th entry of  $\mathbf{g}_i$ , for all  $i = 1, \dots, N$  and  $j = 1, \dots, k$ . Then, for each  $j$ ,  $\{g_{ij}\}_{i=1}^N$  is a sequence of i.i.d. zero-mean sub-Gaussian random variables with variance proxy  $\zeta^2$ . Note that if  $|\frac{1}{N} \sum_{i=1}^N g_{ij}| \leq \frac{\tau}{k}$  for all  $j = 1, \dots, k$ , then  $\|\frac{1}{N} \sum_{i=1}^N \mathbf{g}_i\| \leq \tau$ . Hence,

$$\begin{aligned} \Pr\left(\left\|\frac{1}{N} \sum_{i=1}^N \mathbf{g}_i\right\| \geq \tau\right) &\leq \Pr\left(\left|\frac{1}{N} \sum_{i=1}^N g_{ij}\right| \geq \frac{\tau}{k} \text{ for some } j = 1, \dots, k\right) \\ &\leq \sum_{j=1}^k \Pr\left(\left|\frac{1}{N} \sum_{i=1}^N g_{ij}\right| \geq \frac{\tau}{k}\right) \leq 2k \exp\left(-\frac{N\tau^2}{2k^2\zeta^2}\right), \end{aligned}$$

where the last inequality follows from Lemma EC.5.  $\square$

**DEFINITION EC.1 (SUB-EXPONENTIAL RANDOM VARIABLES).** A random variable  $X$  is said to be *sub-exponential* if its *sub-exponential norm*, defined as  $\|X\|_{\psi_1} := \inf\{w > 0 : \mathbb{E}[e^{|X|/w}] \leq 2\}$ , is finite.

**LEMMA EC.6.** *Let  $\{X_i\}_{i=1}^N$  be a sequence of i.i.d. zero-mean sub-exponential random variables. Then, for all  $\tau \geq 0$ ,*

$$\Pr\left(\left|\frac{1}{N} \sum_{i=1}^N X_i\right| \geq \tau\right) \leq 2 \exp\left(-N C_{\text{Bern}} \cdot \min\left\{\frac{\tau^2}{K^2}, \frac{\tau}{K}\right\}\right),$$

where  $C_{\text{Bern}}$  is an absolute constant and  $K = \|X_1\|_{\psi_1}$ .

*Proof.* This is Bernstein's inequality for sub-exponential random variables. See Theorem 2.8.2 of Vershynin (2018).  $\square$

## EC.2. Coefficient Stabilization Phase: $t = 1, \dots, T_1$

Throughout the rest of the Supplemental Material, we let  $T_1 = C_{T_1} \log(T)$  and  $T_2 = (C_{T_1} + C_{T_2}) \log(T)$  for some constant  $C_{T_2}$  for some constants  $C_{T_1}$  and  $C_{T_2}$ ; we also suppose that Assumption 1 holds and the data are generated by Algorithm 1.

We first define a set of notations before the analysis. For any  $t$ , let

$$\begin{aligned} \tilde{\mathbf{v}}_t &:= \begin{pmatrix} \mathbb{I}(a_t = 1) \mathbf{v}_t \\ \vdots \\ \mathbb{I}(a_t = M) \mathbf{v}_t \end{pmatrix} \in \mathbb{R}^{Mp}, \quad \tilde{\mathbf{z}}_t := \begin{pmatrix} \mathbb{I}(a_t = 1) \mathbf{z}_t \\ \vdots \\ \mathbb{I}(a_t = M) \mathbf{z}_t \end{pmatrix} \in \mathbb{R}^{Mq}, \\ \tilde{\mathbf{V}}_t &:= \begin{pmatrix} \tilde{\mathbf{v}}_1^\top \\ \vdots \\ \tilde{\mathbf{v}}_t^\top \end{pmatrix} \in \mathbb{R}^{t \times (Mp)}, \quad \tilde{\mathbf{Z}}_t := \begin{pmatrix} \tilde{\mathbf{z}}_1^\top \\ \vdots \\ \tilde{\mathbf{z}}_t^\top \end{pmatrix} \in \mathbb{R}^{t \times (Mq)}, \quad \boldsymbol{\epsilon}_t := \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_t \end{pmatrix} \in \mathbb{R}^t, \\ \tilde{\boldsymbol{\Sigma}}_{vz,t} &:= \frac{1}{t} \sum_{s=1}^t \tilde{\mathbf{v}}_s \tilde{\mathbf{z}}_s^\top = \frac{1}{t} \tilde{\mathbf{V}}_t^\top \tilde{\mathbf{Z}}_t \in \mathbb{R}^{(Mp) \times (Mq)}, \end{aligned}$$

$$\begin{aligned}\tilde{\Sigma}_{zz,t} &:= \frac{1}{t} \sum_{s=1}^t \tilde{\mathbf{z}}_s \tilde{\mathbf{z}}_s^\top = \frac{1}{t} \tilde{\mathbf{Z}}_t^\top \tilde{\mathbf{Z}}_t \in \mathbb{R}^{(Mq) \times (Mq)}, \\ \tilde{\mathbf{G}}_t &:= \frac{1}{t} \sum_{s=1}^t \tilde{\mathbf{z}}_s \epsilon_s = \frac{1}{t} \tilde{\mathbf{Z}}_t^\top \boldsymbol{\epsilon}_t \in \mathbb{R}^{2q}, \\ \Sigma_{vz} &:= \mathbb{E}[\mathbf{v}_t \mathbf{z}_t^\top] \in \mathbb{R}^{p \times q}, \quad \Sigma_{zz} := \mathbb{E}[\mathbf{z}_t \mathbf{z}_t^\top] \in \mathbb{R}^{q \times q}, \\ \check{\Sigma}_{vz} &:= \mathbb{E}[\tilde{\mathbf{v}}_t \tilde{\mathbf{z}}_t^\top] \in \mathbb{R}^{(Mp) \times (Mq)}, \quad \check{\Sigma}_{zz} := \mathbb{E}[\tilde{\mathbf{z}}_t \tilde{\mathbf{z}}_t^\top] \in \mathbb{R}^{(Mq) \times (Mq)}.\end{aligned}$$

LEMMA EC.7. *Both  $\check{\Sigma}_{zz}$  and  $\check{\Sigma}_{vz}$  have full rank.*

*Proof.* By condition (iii) of Assumption 1,  $\Sigma_{vz}^* = \mathbb{E}[\mathbf{v}_t^* \mathbf{z}_t^\top]$  is full rank and  $\Sigma_{zz} = \mathbb{E}[\mathbf{z}_t \mathbf{z}_t^\top]$  is positive definite. By definition, for  $t = 1, 2, \dots, T_1$ ,  $a_t$  is randomly selected such that  $\Pr(a_t = i) = \frac{1}{M}$  for all  $i = 1, \dots, M$ . Thus,

$$\check{\Sigma}_{vz} = \frac{1}{M} \begin{pmatrix} \Sigma_{zz} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \Sigma_{zz} \end{pmatrix},$$

and thus  $\check{\Sigma}_{vz}$  is full rank and positive definite.

Because  $\sum_{i=1}^M \mathbb{I}(a_t = i) = 1$  and  $\Sigma_{vz}^*$  is a full-rank matrix of size  $(Mp) \times q$ ,  $\Sigma_{vz} = \left( \mathbf{I}_p \cdots \mathbf{I}_p \right) \Sigma_{vz}^*$  is a full-rank matrix of size  $p \times q$ , where  $\mathbf{I}_p$  is the identity matrix of size  $p \times p$ . As a result,

$$\check{\Sigma}_{vz} = \mathbb{E}[\tilde{\mathbf{v}}_t \tilde{\mathbf{z}}_t^\top] = \begin{pmatrix} \mathbb{E}[\mathbb{I}(a_t = 1) \mathbf{v}_t \mathbf{z}_t^\top] & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \mathbb{E}[\mathbb{I}(a_t = M) \mathbf{v}_t \mathbf{z}_t^\top] \end{pmatrix} = \frac{1}{M} \begin{pmatrix} \Sigma_{vz} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \Sigma_{vz} \end{pmatrix}.$$

Thus,  $\check{\Sigma}_{vz}$  is a full-rank matrix of size  $(Mp) \times (Mq)$ .  $\square$

LEMMA EC.8. *Define*

$$\check{\kappa} := \frac{2\phi_{\max}(\check{\Sigma}_{vz})\phi_{\max}(\check{\Sigma}_{zz})}{(\phi_{\min}(\check{\Sigma}_{vz}))^2\phi_{\min}(\check{\Sigma}_{zz})} \quad (\text{EC.2.1})$$

and

$$\mathcal{A}_{T_1} := \left\{ \|\hat{\boldsymbol{\alpha}}_{T_1} - \boldsymbol{\alpha}\| \leq 2Mq\bar{z}\check{\kappa}C_{T_1}^{-\frac{1}{2}} \right\}. \quad (\text{EC.2.2})$$

Suppose

$$C_{T_1} \geq \left( \frac{40M \max\{\sqrt{pq}\bar{v}\bar{z}, q\bar{z}^2\}}{\min\{\phi_{\min}(\check{\Sigma}_{vz}), \phi_{\min}(\check{\Sigma}_{zz})\}} \right)^2. \quad (\text{EC.2.3})$$

Then,  $\Pr(\mathcal{A}_{T_1}) \geq 1 - 2M(2p + 5q)T^{-2}$  for all  $T \geq T_1$ .

*Proof.* It follows from Lemma EC.7 that  $0 < \phi_{\min}(\check{\Sigma}_{zz}) \leq \phi_{\max}(\check{\Sigma}_{zz}) < \infty$  and  $0 < \phi_{\min}(\check{\Sigma}_{vz}) \leq \phi_{\max}(\check{\Sigma}_{vz}) < \infty$ . Hence,  $\check{\kappa}$  is a positive finite constant.

Because in the coefficient stabilization phase  $a_t$  is randomly selected for all  $t = 1, \dots, T_1$ , we may rewrite the arm-specific 2SLS for the  $M$  arms together equivalently as the joint-2SLS using the notations  $\tilde{\mathbf{V}}_{T_1}$  and  $\tilde{\mathbf{Z}}_{T_1}$ . Hence,

$$\begin{aligned}\hat{\boldsymbol{\alpha}}_{T_1} - \boldsymbol{\alpha} &= \left( \tilde{\mathbf{V}}_{T_1}^\top \mathcal{P}[\tilde{\mathbf{Z}}_{T_1}] \tilde{\mathbf{V}}_{T_1} \right)^{-1} \tilde{\mathbf{V}}_{T_1}^\top \mathcal{P}[\tilde{\mathbf{Z}}_{T_1}] \boldsymbol{\epsilon}_{T_1} \\ &= \left( \tilde{\mathbf{V}}_{T_1}^\top \tilde{\mathbf{Z}}_{T_1} (\tilde{\mathbf{Z}}_{T_1}^\top \tilde{\mathbf{Z}}_{T_1})^{-1} \tilde{\mathbf{Z}}_{T_1}^\top \tilde{\mathbf{V}}_{T_1} \right)^{-1} \tilde{\mathbf{V}}_{T_1}^\top \tilde{\mathbf{Z}}_{T_1} (\tilde{\mathbf{Z}}_{T_1}^\top \tilde{\mathbf{Z}}_{T_1})^{-1} \tilde{\mathbf{Z}}_{T_1}^\top \boldsymbol{\epsilon}_{T_1} \\ &= \left( \tilde{\boldsymbol{\Sigma}}_{vz, T_1} \tilde{\boldsymbol{\Sigma}}_{zz, T_1}^{-1} \tilde{\boldsymbol{\Sigma}}_{vz, T_1}^\top \right)^{-1} \tilde{\boldsymbol{\Sigma}}_{vz, T_1} \tilde{\boldsymbol{\Sigma}}_{zz, T_1}^{-1} \tilde{\mathbf{G}}_{T_1},\end{aligned}$$

where the projection operator  $\mathcal{P}$  maps a matrix  $\mathbf{Z}$  to  $\mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top$ . Therefore,

$$\|\hat{\boldsymbol{\alpha}}_{T_1} - \boldsymbol{\alpha}\| \leq \left\| \left( \tilde{\boldsymbol{\Sigma}}_{vz, T_1} \tilde{\boldsymbol{\Sigma}}_{zz, T_1}^{-1} \tilde{\boldsymbol{\Sigma}}_{vz, T_1}^\top \right)^{-1} \tilde{\boldsymbol{\Sigma}}_{vz, T_1} \tilde{\boldsymbol{\Sigma}}_{zz, T_1}^{-1} \right\| \cdot \|\tilde{\mathbf{G}}_{T_1}\| \quad (\text{EC.2.4})$$

In what follows, we analyze the two norms on the right-hand-side of (EC.2.4) separately.

Because the samples are i.i.d.,

$$\mathbb{E}[\tilde{\boldsymbol{\Sigma}}_{vz, T_1}] = \mathbb{E}[\tilde{\mathbf{v}}_s \tilde{\mathbf{z}}_s^\top] = \check{\boldsymbol{\Sigma}}_{vz} \quad \text{and} \quad \mathbb{E}[\tilde{\boldsymbol{\Sigma}}_{zz, T_1}] = \mathbb{E}[\tilde{\mathbf{z}}_s \tilde{\mathbf{z}}_s^\top] = \check{\boldsymbol{\Sigma}}_{zz}.$$

Moreover, because  $\tilde{\mathbf{v}}_s \tilde{\mathbf{z}}_s^\top$  is of size  $(Mp)$ -by- $(Mq)$ , we know by Lemma EC.3 that

$$\|\tilde{\mathbf{v}}_s \tilde{\mathbf{z}}_s^\top\| \leq \sqrt{(Mp) \times (Mq)} \|\tilde{\mathbf{v}}_s \tilde{\mathbf{z}}_s^\top\|_{\max} \leq M \sqrt{pq} \bar{v} \bar{z}.$$

Likewise,  $\|\tilde{\mathbf{z}}_s \tilde{\mathbf{z}}_s^\top\| \leq Mq \bar{z}^2$ . Hence, if we define the following events:

$$\begin{aligned}\mathcal{S}_{vz, \min} &:= \left\{ \phi_{\min}(\tilde{\boldsymbol{\Sigma}}_{vz, T_1}) \leq \phi_{\min}(\check{\boldsymbol{\Sigma}}_{vz}) - \tau \right\}, \\ \mathcal{S}_{vz, \max} &:= \left\{ \phi_{\max}(\tilde{\boldsymbol{\Sigma}}_{vz, T_1}) \geq \phi_{\max}(\check{\boldsymbol{\Sigma}}_{vz}) + \tau \right\}, \\ \mathcal{S}_{zz, \min} &:= \left\{ \phi_{\min}(\tilde{\boldsymbol{\Sigma}}_{zz, T_1}) \leq \phi_{\min}(\check{\boldsymbol{\Sigma}}_{zz}) - \tau \right\}, \\ \mathcal{S}_{zz, \max} &:= \left\{ \phi_{\max}(\tilde{\boldsymbol{\Sigma}}_{zz, T_1}) \geq \phi_{\max}(\check{\boldsymbol{\Sigma}}_{zz}) + \tau \right\},\end{aligned}$$

then we can apply Corollary EC.1 to conclude that for all  $\tau \geq 0$ ,

$$\begin{aligned}\Pr(\mathcal{S}_{vz, \min}) &\leq 2M(p+q) \exp\left(-\frac{\tau^2 T_1}{8M^2 pq \bar{v}^2 \bar{z}^2}\right), \\ \Pr(\mathcal{S}_{vz, \max}) &\leq 2M(p+q) \exp\left(-\frac{\tau^2 T_1}{8M^2 pq \bar{v}^2 \bar{z}^2}\right), \\ \Pr(\mathcal{S}_{zz, \min}) &\leq 2Mq \exp\left(-\frac{\tau^2 T_1}{8M^2 q^2 \bar{z}^4}\right), \\ \Pr(\mathcal{S}_{zz, \max}) &\leq 2Mq \exp\left(-\frac{\tau^2 T_1}{8M^2 q^2 \bar{z}^4}\right).\end{aligned}$$

Setting  $\tau = \zeta_1 := 4M \max\{\sqrt{pq}\bar{v}\bar{z}, q\bar{z}^2\}C_{T_1}^{-\frac{1}{2}}$  leads to

$$\exp\left(-\frac{\zeta_1^2 T_1}{8M^2 pq \bar{v}^2 \bar{z}^2}\right) \leq T^{-2} \quad \text{and} \quad \exp\left(-\frac{\zeta_1^2 T_1}{8M^2 q^2 \bar{z}^4}\right) \leq T^{-2}.$$

Moreover, by the condition (EC.2.3), we have

$$\zeta_1 = 4M \max\{\sqrt{pq}\bar{v}\bar{z}, q\bar{z}^2\}C_{T_1}^{-\frac{1}{2}} \leq \frac{1}{10} \min\left\{\phi_{\min}(\check{\Sigma}_{vz}), \phi_{\min}(\check{\Sigma}_{zz})\right\}. \quad (\text{EC.2.5})$$

Therefore, upon defining the joint event  $\mathcal{S}_{\phi, T_1} := \mathcal{S}_{vz, \min}^c \cap \mathcal{S}_{vz, \max}^c \cap \mathcal{S}_{zz, \min}^c \cap \mathcal{S}_{zz, \max}^c$ , it follows from the Bonferroni inequality that

$$\begin{aligned} \Pr(\mathcal{S}_{\phi, T_1}) &\geq 1 - [\Pr(\mathcal{S}_{vz, \min}) + \Pr(\mathcal{S}_{vz, \max}) + \Pr(\mathcal{S}_{zz, \min}) + \Pr(\mathcal{S}_{zz, \max})] \\ &\geq 1 - 4M(p + 2q)T^{-2}. \end{aligned} \quad (\text{EC.2.6})$$

Note that conditional on the event  $\mathcal{S}_{\phi, T_1}$  with  $\tau = \zeta_1$ , we have

$$\begin{aligned} \phi_{\min}(\tilde{\Sigma}_{vz, T_1}) &> \phi_{\min}(\check{\Sigma}_{vz}) - \zeta_1 > 0, & \phi_{\max}(\tilde{\Sigma}_{vz, T_1}) &< \phi_{\max}(\check{\Sigma}_{vz}) + \zeta_1, \\ \phi_{\min}(\tilde{\Sigma}_{zz, T_1}) &> \phi_{\min}(\check{\Sigma}_{zz}) - \zeta_1 > 0, & \phi_{\max}(\tilde{\Sigma}_{zz, T_1}) &< \phi_{\max}(\check{\Sigma}_{zz}) + \zeta_1. \end{aligned} \quad (\text{EC.2.7})$$

This implies that

$$\begin{aligned} &\left\| \left( \tilde{\Sigma}_{vz, T_1} \tilde{\Sigma}_{zz, T_1}^{-1} \tilde{\Sigma}_{vz, T_1}^\top \right)^{-1} \tilde{\Sigma}_{vz, T_1} \tilde{\Sigma}_{zz, T_1}^{-1} \right\| \\ &= \phi_{\max} \left( \left( \tilde{\Sigma}_{vz, T_1} \tilde{\Sigma}_{zz, T_1}^{-1} \tilde{\Sigma}_{vz, T_1}^\top \right)^{-1} \tilde{\Sigma}_{vz, T_1} \tilde{\Sigma}_{zz, T_1}^{-1} \right) \\ &\leq \frac{\phi_{\max}(\tilde{\Sigma}_{vz, T_1}) \phi_{\max}(\tilde{\Sigma}_{zz, T_1})}{(\phi_{\min}(\tilde{\Sigma}_{vz, T_1}))^2 \phi_{\min}(\tilde{\Sigma}_{zz, T_1})} \end{aligned} \quad (\text{EC.2.8})$$

$$\leq \frac{(\phi_{\max}(\check{\Sigma}_{vz}) + \zeta_1)(\phi_{\max}(\check{\Sigma}_{zz}) + \zeta_1)}{(\phi_{\min}(\check{\Sigma}_{vz}) - \zeta_1)^2 (\phi_{\min}(\check{\Sigma}_{zz}) - \zeta_1)} \quad (\text{EC.2.9})$$

$$\leq \left(\frac{11}{10}\right)^2 \left(\frac{10}{9}\right)^3 \frac{\phi_{\max}(\check{\Sigma}_{vz}) \phi_{\max}(\check{\Sigma}_{zz})}{(\phi_{\min}(\check{\Sigma}_{vz}))^2 \phi_{\min}(\check{\Sigma}_{zz})} \quad (\text{EC.2.10})$$

$$\leq \check{\kappa}, \quad (\text{EC.2.11})$$

where (EC.2.8) follows from Lemma EC.2, (EC.2.9) from (EC.2.7), (EC.2.10) from (EC.2.5), and (EC.2.11) from (EC.2.1). Therefore,

$$\mathcal{S}_{\phi, T_1} \subseteq \left\{ \left\| \left( \tilde{\Sigma}_{vz, T_1} \tilde{\Sigma}_{zz, T_1}^{-1} \tilde{\Sigma}_{vz, T_1}^\top \right)^{-1} \tilde{\Sigma}_{vz, T_1} \tilde{\Sigma}_{zz, T_1}^{-1} \right\| \leq \check{\kappa} \right\}. \quad (\text{EC.2.12})$$

Note that  $\mathbb{E}[\tilde{\mathbf{G}}_{T_1}] = \mathbb{E}[\tilde{\mathbf{z}}_s \epsilon_s] = \mathbf{0}$ . Moreover, note that  $\tilde{\mathbf{z}}_s \epsilon_s \in \mathbb{R}^{Mq}$  is a vector of zero-mean sub-Gaussian random variables with variance proxy  $\bar{z}^2 \zeta^2$ . Hence, applying Corollary EC.2 to  $\tilde{\mathbf{G}}_{T_1}$  yields that for all  $\tau \geq 0$ ,

$$\Pr\left(\|\tilde{\mathbf{G}}_{T_1}\| \geq \tau\right) \leq 2Mq \exp\left(-\frac{\tau^2 T_1}{2M^2 q^2 \bar{z}^2 \zeta^2}\right).$$

Thus, setting  $\tau = \zeta_2 := 2Mq\bar{z}\zeta C_{T_1}^{-\frac{1}{2}}$  yields

$$\Pr\left(\|\tilde{\mathbf{G}}_{T_1}\| > \zeta_2\right) \leq 2Mq \exp(-2T_1 C_{T_1}^{-1}) = 2Mq \exp(-2\log(T)) = 2MqT^{-2}. \quad (\text{EC.2.13})$$

It follows from (EC.2.4) and (EC.2.12) that

$$\mathcal{S}_{\phi, T_1} \cap \{\|\tilde{\mathbf{G}}_{T_1}\| \leq \zeta_2\} \subseteq \{\|\hat{\boldsymbol{\alpha}}_{T_1} - \boldsymbol{\alpha}\| \leq \zeta_2 \check{\kappa}\} = \mathcal{A}_{T_1}.$$

Hence,

$$\Pr(\mathcal{A}_{T_1}) \geq \Pr\left(\mathcal{S}_{\phi, T_1} \cap \{\|\tilde{\mathbf{G}}_{T_1}\| \leq \zeta_2\}\right) \geq 1 - 2M(2p + 5q)T^{-2},$$

where the second inequality follows from (EC.2.6) and (EC.2.13).  $\square$

### EC.3. Covariance Stabilization Phase and Beyond: $t > T_1$

For any  $t$ , let

$$\tilde{\mathbf{v}}_t := \begin{pmatrix} \mathbb{I}(a_t = 1)\mathbf{v}_t \\ \vdots \\ \mathbb{I}(a_t = M)\mathbf{v}_t \end{pmatrix} \in \mathbb{R}^{Mp}, \quad \mathbf{v}_t^* := \begin{pmatrix} \mathbb{I}(a_t^* = 1)\mathbf{v}_t \\ \vdots \\ \mathbb{I}(a_t^* = M)\mathbf{v}_t \end{pmatrix} \in \mathbb{R}^{Mp}.$$

For any  $t_1 < t_2$ , let

$$\begin{aligned} \tilde{\mathbf{V}}_{t_1:t_2} &:= \begin{pmatrix} \tilde{\mathbf{v}}_{t_1+1}^\top \\ \vdots \\ \tilde{\mathbf{v}}_{t_2}^\top \end{pmatrix} \in \mathbb{R}^{(t_2-t_1) \times (Mp)}, \quad \tilde{\mathbf{V}}_{t_1:t_2}^* := \begin{pmatrix} (\mathbf{v}_{t_1+1}^*)^\top \\ \vdots \\ (\mathbf{v}_{t_2}^*)^\top \end{pmatrix} \in \mathbb{R}^{(t_2-t_1) \times (Mp)}, \\ \mathbf{Z}_{t_1:t_2} &:= \begin{pmatrix} \mathbf{z}_{t_1+1}^\top \\ \vdots \\ \mathbf{z}_{t_2}^\top \end{pmatrix} \in \mathbb{R}^{(t_2-t_1) \times q}, \quad \mathbf{R}_{t_1:t_2} := \begin{pmatrix} R_{t_1+1} \\ \vdots \\ R_{t_2} \end{pmatrix} \in \mathbb{R}^{t_2-t_1}, \quad \boldsymbol{\epsilon}_{t_1:t_2} := \begin{pmatrix} \epsilon_{t_1+1} \\ \vdots \\ \epsilon_{t_2} \end{pmatrix} \in \mathbb{R}^{t_2-t_1}, \\ \hat{\boldsymbol{\Sigma}}_{vz, t_1:t_2} &:= \frac{1}{t_2-t_1} \sum_{s=t_1+1}^{t_2} \tilde{\mathbf{v}}_s \mathbf{z}_s^\top = \frac{1}{t_2-t_1} \tilde{\mathbf{V}}_{t_1:t_2}^\top \mathbf{Z}_{t_1:t_2} \in \mathbb{R}^{(Mp) \times q}, \\ \hat{\boldsymbol{\Sigma}}_{vz, t_1:t_2}^* &:= \frac{1}{t_2-t_1} \sum_{s=t_1+1}^{t_2} \mathbf{v}_s^* \mathbf{z}_s^\top = \frac{1}{t_2-t_2} (\mathbf{V}_{t_1:t_2}^*)^\top \mathbf{Z}_{t_1:t_2} \in \mathbb{R}^{(Mp) \times q}, \\ \hat{\boldsymbol{\Sigma}}_{zz, t_1:t_2} &:= \frac{1}{t_2-t_1} \sum_{s=t_1+1}^{t_2} \mathbf{z}_s \mathbf{z}_s^\top = \frac{1}{t_2-t_1} \mathbf{Z}_{t_1:t_2}^\top \mathbf{Z}_{t_1:t_2} \in \mathbb{R}^{q \times q}, \\ \hat{\mathbf{G}}_{t_1:t_2} &:= \frac{1}{t_2-t_1} \sum_{s=t_1+1}^{t_2} \mathbf{z}_s \epsilon_s = \frac{1}{t_2-t_1} \mathbf{Z}_{t_1:t_2}^\top \boldsymbol{\epsilon}_{t_1:t_2} \in \mathbb{R}^q, \\ \boldsymbol{\Sigma}_{vz}^* &:= \mathbb{E}[\mathbf{v}_t^* \mathbf{z}_t^\top] \in \mathbb{R}^{(Mp) \times q}, \quad \boldsymbol{\Sigma}_{zz} := \mathbb{E}[\mathbf{z}_t \mathbf{z}_t^\top] \in \mathbb{R}^{q \times q}. \end{aligned}$$

Note that for  $t > T_1$ , the 2SLS estimator  $\hat{\alpha}_t$  in Algorithm 1 satisfies

$$\begin{aligned}\hat{\alpha}_t - \alpha &= (\tilde{\mathbf{V}}_{T_1:t}^\top \mathcal{P}[\mathbf{Z}_{T_1:t}] \tilde{\mathbf{V}}_{T_1:t})^{-1} \tilde{\mathbf{V}}_{T_1:t}^\top \mathcal{P}[\mathbf{Z}_{T_1:t}] \epsilon_{T_1:t} \\ &= \left( \hat{\Sigma}_{vz, T_1:t} \hat{\Sigma}_{zz, T_1:t}^{-1} \hat{\Sigma}_{vz, T_1:t}^\top \right)^{-1} \hat{\Sigma}_{vz, T_1:t} \hat{\Sigma}_{zz, T_1:t}^{-1} \hat{\mathbf{G}}_{T_1:t}.\end{aligned}\quad (\text{EC.3.1})$$

Moreover, recall that in Algorithm 1 we use the following estimator for the variance of noise terms:

$$\hat{\sigma}_t := (t - T_1)^{-\frac{1}{2}} \|\mathbf{R}_{T_1:t} - \tilde{\mathbf{V}}_{T_1:t} \hat{\alpha}_t\| \in \mathbb{R}_+.\quad (\text{EC.3.2})$$

In order to establish a probabilistic bound on  $\|\hat{\alpha}_t - \alpha\|$ , we establish below bounds on a set of events that characterize how accurate the estimators such as  $\hat{\Sigma}_{vz, T_1:t}$ ,  $\hat{\Sigma}_{zz, T_1:t}$ , and  $\hat{\sigma}_t^2$  are relative to their corresponding true values.

Specifically, we fix arbitrary constants  $C_{\mathcal{A}}, \eta, \Delta > 0$ , and define for any  $t > T_1$  the following events:

$$\begin{aligned}\mathcal{A}_t &:= \left\{ \|\hat{\alpha}_t - \alpha\| \leq C_{\mathcal{A}} \sqrt{\frac{\log(T)}{t - T_1}} \right\}, & \mathcal{B}_t &:= \left\{ \|\hat{\Sigma}_{vz, T_1:t} - \Sigma_{vz}^*\| \leq \eta \right\}, \\ \mathcal{C}_t &:= \left\{ \|\hat{\Sigma}_{zz, T_1:t} - \Sigma_{zz}\| \leq \eta \right\}, & \mathcal{D}_t &:= \left\{ |\hat{\sigma}_t^2 - \sigma^2| \leq \Delta \right\}.\end{aligned}$$

In Lemmas EC.9—EC.13 below, we will establish probabilistic bounds of these events and their various intersections.

LEMMA EC.9. *Define*

$$\bar{\kappa} := \frac{2\phi_{\max}(\Sigma_{vz}^*)\phi_{\max}(\Sigma_{zz})}{(\phi_{\min}(\Sigma_{vz}^*))^2\phi_{\min}(\Sigma_{zz})}.\quad (\text{EC.3.3})$$

Suppose that the constants  $C_{\mathcal{A}}$  and  $\eta$  in the definitions of  $\mathcal{A}_t$ ,  $\mathcal{B}_t$  and  $\mathcal{C}_t$  satisfy

$$C_{\mathcal{A}} \geq 2q\bar{z}\bar{\kappa},\quad (\text{EC.3.4})$$

$$\eta \leq \frac{1}{10} \min\left\{ \phi_{\min}(\Sigma_{vz}^*), \phi_{\min}(\Sigma_{zz}) \right\}.\quad (\text{EC.3.5})$$

Then, for all  $t = T_1 + 1, \dots, T$ ,

$$\Pr(\mathcal{A}_t^c \cap \mathcal{B}_t \cap \mathcal{C}_t) \leq 2qT^{-2}.$$

*Proof.* It follows from Assumption 1 that  $0 < \phi_{\min}(\Sigma_{zz}) \leq \phi_{\max}(\Sigma_{zz}) < \infty$  and  $0 < \phi_{\min}(\Sigma_{vz}^*) \leq \phi_{\max}(\Sigma_{vz}^*) < \infty$ . Hence,  $\bar{\kappa}$  is a positive finite constant.

We apply an analysis that is similar to that used in the proof of Lemma EC.8. Specifically, it follows from (EC.3.1) that for all  $t > T_1$ ,

$$\|\hat{\alpha}_t - \alpha\| \leq \left\| \left( \hat{\Sigma}_{vz, T_1:t} \hat{\Sigma}_{zz, T_1:t}^{-1} \hat{\Sigma}_{vz, T_1:t}^\top \right)^{-1} \hat{\Sigma}_{vz, T_1:t} \hat{\Sigma}_{zz, T_1:t}^{-1} \right\| \cdot \|\hat{\mathbf{G}}_{T_1:t}\|.\quad (\text{EC.3.6})$$

We analyze below the two norms on the right-hand-side of (EC.3.6) separately.

Note that conditional on the event  $\mathcal{B}_t \cap \mathcal{C}_t$  with  $\eta < \min\{\phi_{\min}(\boldsymbol{\Sigma}_{vz}^*), \phi_{\min}(\boldsymbol{\Sigma}_{zz})\}$ , which is implied by condition (EC.3.5), we have

$$\begin{aligned} \phi_{\min}(\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}) &> \phi_{\min}(\boldsymbol{\Sigma}_{vz}^*) - \eta > 0, & \phi_{\max}(\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}) &< \phi_{\max}(\boldsymbol{\Sigma}_{vz}^*) + \eta, \\ \phi_{\min}(\widehat{\boldsymbol{\Sigma}}_{zz, T_1:t}) &> \phi_{\min}(\boldsymbol{\Sigma}_{zz}) - \eta > 0, & \phi_{\max}(\widehat{\boldsymbol{\Sigma}}_{zz, T_1:t}) &< \phi_{\max}(\boldsymbol{\Sigma}_{zz}) + \eta. \end{aligned} \quad (\text{EC.3.7})$$

This implies that

$$\begin{aligned} \left\| \left( \widehat{\boldsymbol{\Sigma}}_{vz, T_1:t} \widehat{\boldsymbol{\Sigma}}_{zz, T_1:t}^{-1} \widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}^\top \right)^{-1} \widehat{\boldsymbol{\Sigma}}_{vz, T_1:t} \widehat{\boldsymbol{\Sigma}}_{zz, T_1:t}^{-1} \right\| &= \phi_{\max} \left( \left( \widehat{\boldsymbol{\Sigma}}_{vz, T_1:t} \widehat{\boldsymbol{\Sigma}}_{zz, T_1:t}^{-1} \widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}^\top \right)^{-1} \widehat{\boldsymbol{\Sigma}}_{vz, T_1:t} \widehat{\boldsymbol{\Sigma}}_{zz, T_1:t}^{-1} \right) \\ &\leq \frac{\phi_{\max}(\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}) \phi_{\max}(\widehat{\boldsymbol{\Sigma}}_{zz, T_1:t})}{(\phi_{\min}(\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}))^2 \phi_{\min}(\widehat{\boldsymbol{\Sigma}}_{zz, T_1:t})} \end{aligned} \quad (\text{EC.3.8})$$

$$\leq \frac{(\phi_{\max}(\boldsymbol{\Sigma}_{vz}^*) + \eta)(\phi_{\max}(\boldsymbol{\Sigma}_{zz}) + \eta)}{(\phi_{\min}(\boldsymbol{\Sigma}_{vz}^*) - \eta)^2 (\phi_{\min}(\boldsymbol{\Sigma}_{zz}) - \eta)} \quad (\text{EC.3.9})$$

$$\leq \left(\frac{11}{10}\right)^2 \left(\frac{10}{9}\right)^3 \frac{\phi_{\max}(\boldsymbol{\Sigma}_{vz}^*) \phi_{\max}(\boldsymbol{\Sigma}_{zz})}{(\phi_{\min}(\boldsymbol{\Sigma}_{vz}^*))^2 \phi_{\min}(\boldsymbol{\Sigma}_{zz})} \quad (\text{EC.3.10})$$

$$\leq \bar{\kappa}, \quad (\text{EC.3.11})$$

where (EC.3.8) follows from Lemma EC.2, (EC.3.9) from (EC.3.7), (EC.3.10) from (EC.3.5), and (EC.3.11) from (EC.3.3). Therefore,

$$\mathcal{B}_t \cap \mathcal{C}_t \subseteq \left\{ \left\| \left( \widehat{\boldsymbol{\Sigma}}_{vz, T_1:t} \widehat{\boldsymbol{\Sigma}}_{zz, T_1:t}^{-1} \widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}^\top \right)^{-1} \widehat{\boldsymbol{\Sigma}}_{vz, T_1:t} \widehat{\boldsymbol{\Sigma}}_{zz, T_1:t}^{-1} \right\| \leq \bar{\kappa} \right\}. \quad (\text{EC.3.12})$$

Next, note that  $\mathbb{E}[\widehat{\mathbf{G}}_{T_1:t}] = \mathbb{E}[\mathbf{z}_s \epsilon_s] = \mathbf{0}$ . Moreover, it is straightforward to see that  $\mathbf{z}_s \epsilon_s \in \mathbb{R}^q$  is a vector of zero-mean sub-Gaussian random variables with variance proxy  $\bar{z}^2 \zeta^2$ . Hence, applying Corollary EC.2 to  $\widehat{\mathbf{G}}_{T_1:t}$  yields that for all  $\tau \geq 0$ ,

$$\Pr\left(\|\widehat{\mathbf{G}}_{T_1:t}\| \geq \tau\right) \leq 2q \exp\left(-\frac{\tau^2(t - T_1)}{2q^2 \bar{z}^2 \zeta^2}\right). \quad (\text{EC.3.13})$$

Thus, setting  $\tau = \zeta_3 := 2q\bar{z}\zeta\sqrt{\frac{\log(T)}{t - T_1}}$  yields

$$\Pr\left(\|\widehat{\mathbf{G}}_{T_1:t}\| > \zeta_3\right) \leq 2q \exp(-2\log(T)) = 2qT^{-2}. \quad (\text{EC.3.14})$$

It follows from (EC.3.6) and (EC.3.12) that

$$\left\{ \|\widehat{\mathbf{G}}_{T_1:t}\| \leq \zeta_3 \right\} \cap \mathcal{B}_t \cap \mathcal{C}_t \subseteq \left\{ \|\widehat{\boldsymbol{\alpha}}_t - \boldsymbol{\alpha}\| \leq \zeta_3 \bar{\kappa} \right\} \subseteq \left\{ \|\widehat{\boldsymbol{\alpha}}_t - \boldsymbol{\alpha}\| \leq C_{\mathcal{A}} \sqrt{\frac{\log(T)}{t - T_1}} \right\} = \mathcal{A}_t,$$

where the second step follows from the assumption on  $C_{\mathcal{A}}$  in (EC.3.4) as well as the definitions of  $\zeta_3$  and  $\bar{\kappa}$ . Hence,

$$\left\{ \|\widehat{\mathbf{G}}_{T_1:t}\| \leq \zeta_3 \right\} \cap \mathcal{B}_t \cap \mathcal{C}_t \subseteq \mathcal{A}_t \cap \mathcal{B}_t \cap \mathcal{C}_t,$$

and thus

$$\{\|\widehat{\mathbf{G}}_{\mathbf{T}_1:t}\| > \zeta_3\} \cap \mathcal{B}_t \cap \mathcal{C}_t \supseteq \mathcal{A}_t^c \cap \mathcal{B}_t \cap \mathcal{C}_t.$$

It follows that

$$\Pr(\mathcal{A}_t^c \cap \mathcal{B}_t \cap \mathcal{C}_t) \leq \Pr(\|\widehat{\mathbf{G}}_{\mathbf{T}_1:t}\| > \zeta_3) \leq 2qT^{-2},$$

where the second inequality follows from (EC.3.14).  $\square$

We now define for any  $t > \mathbf{T}_2$  the following joint events:

$$\mathcal{A}_{\mathbf{T}_2:t} = \bigcap_{s=\mathbf{T}_2+1}^t \mathcal{A}_s, \quad \mathcal{B}_{\mathbf{T}_2:t} = \bigcap_{s=\mathbf{T}_2+1}^t \mathcal{B}_s, \quad \mathcal{C}_{\mathbf{T}_2:t} = \bigcap_{s=\mathbf{T}_2+1}^t \mathcal{C}_s, \quad \mathcal{D}_{\mathbf{T}_2:t} := \bigcap_{s=\mathbf{T}_2+1}^t \mathcal{D}_s.$$

LEMMA EC.10. *Suppose  $C_{\mathbf{T}_2} \geq 6q^2\bar{z}^4\eta^{-2}$ . Then, for all  $t = \mathbf{T}_2 + 1, \dots, T$ ,*

$$\Pr(\mathcal{C}_t) \geq 1 - 2qT^{-3} \quad \text{and} \quad \Pr(\mathcal{C}_{\mathbf{T}_2:t}) \geq 1 - 2q(t - \mathbf{T}_2)T^{-3}.$$

*Proof.* Note that by the definition of  $\widehat{\Sigma}_{zz, \mathbf{T}_1:t}$ ,

$$\mathbb{E}[\widehat{\Sigma}_{zz, \mathbf{T}_1:t}] = \frac{1}{\mathbf{T}_1 - t} \sum_{s=\mathbf{T}_1+1}^t \mathbb{E}[\mathbf{z}_s \mathbf{z}_s^\top] = \Sigma_{zz}.$$

In addition,  $\|\mathbf{z}_s \mathbf{z}_s^\top\| \leq q\bar{z}^2$ . Hence, it follows from Lemma EC.4 that for all  $t \geq \mathbf{T}_2$ ,

$$\begin{aligned} \Pr(\mathcal{C}_t^c) &= \Pr\left(\|\widehat{\Sigma}_{zz, \mathbf{T}_1:t} - \Sigma_{zz}\| > \eta\right) \leq 2q \exp\left(-\frac{\eta^2(t - \mathbf{T}_1)}{2q^2\bar{z}^4}\right) \\ &\leq 2q \exp\left(-\frac{\eta^2(\mathbf{T}_2 - \mathbf{T}_1)}{2q^2\bar{z}^4}\right) = 2q \exp\left(-\frac{\eta^2 C_{\mathbf{T}_2} \log(T)}{2q^2\bar{z}^4}\right) \\ &\leq 2q \exp(-3 \log(T)) = 2qT^{-3}, \end{aligned}$$

where the last inequality holds because  $C_{\mathbf{T}_2} \geq 6q^2\bar{z}^4\eta^{-2}$ . Therefore,  $\Pr(\mathcal{C}_t) \geq 1 - 2qT^{-3}$ ; further, by the Bonferroni inequality,

$$\Pr(\mathcal{C}_{\mathbf{T}_2:t}) \geq 1 - \sum_{s=\mathbf{T}_2+1}^t \Pr(\mathcal{C}_s^c) \geq 1 - 2q(t - \mathbf{T}_2)T^{-3}. \quad \square$$

LEMMA EC.11. *Suppose*

$$C_{\mathbf{T}_2} \geq \max\left\{\frac{3}{C_{\text{Bern}} \min\left\{\frac{\Delta^2}{4K_\epsilon^2}, \frac{\Delta}{2K_\epsilon}\right\}}, \frac{6Mpq^2\bar{v}^2\bar{z}^2\zeta^2\bar{\kappa}^2}{\min\left\{\frac{\Delta^2}{128\sigma^2}, \frac{\Delta}{64}\right\}}\right\}, \quad (\text{EC.3.15})$$

where  $C_{\text{Bern}}$  is the constant in Lemma EC.6 and  $K_\epsilon = \|\epsilon_1^2 - \sigma^2\|_{\psi_1} < \infty$ . Suppose that the constant  $\eta$  in the definitions of  $\mathcal{B}_t$  and  $\mathcal{C}_t$  satisfies

$$\eta \leq \frac{1}{10} \min\left\{\phi_{\min}(\Sigma_{vz}^*), \phi_{\min}(\Sigma_{zz})\right\}. \quad (\text{EC.3.16})$$

Then, for all  $t = \mathbf{T}_2 + 1, \dots, T$ ,

$$\Pr(\mathcal{B}_t \cap \mathcal{C}_t \cap \mathcal{D}_t^c) \leq (2q + 2)T^{-3}.$$

*Proof.* Denote the estimated residual by  $\hat{\epsilon}_t := R_t - \tilde{\mathbf{v}}_t^\top \hat{\boldsymbol{\alpha}}_t$ , and define for all  $t > T_1$ ,

$$\begin{aligned} \hat{\boldsymbol{\epsilon}}_{T_1:t} &:= \begin{pmatrix} \hat{\epsilon}_{T_1} \\ \vdots \\ \hat{\epsilon}_t \end{pmatrix} = \mathbf{R}_{T_1:t} - \tilde{\mathbf{V}}_{T_1:t} \hat{\boldsymbol{\alpha}}_t = \boldsymbol{\epsilon}_{T_1:t} + \tilde{\mathbf{V}}_{T_1:t} \boldsymbol{\alpha} - \tilde{\mathbf{V}}_{T_1:t} \hat{\boldsymbol{\alpha}}_t \\ &= \boldsymbol{\epsilon}_{T_1:t} - \tilde{\mathbf{V}}_{T_1:t} \mathbf{M}_{T_1:t} \hat{\mathbf{G}}_{T_1:t}, \end{aligned} \quad (\text{EC.3.17})$$

where  $\mathbf{M}_{T_1:t} := \left( \hat{\boldsymbol{\Sigma}}_{vz, T_1:t} \hat{\boldsymbol{\Sigma}}_{zz, T_1:t}^{-1} \hat{\boldsymbol{\Sigma}}_{vz, T_1:t}^\top \right)^{-1} \hat{\boldsymbol{\Sigma}}_{vz, T_1:t} \hat{\boldsymbol{\Sigma}}_{zz, T_1:t}^{-1}$ , and the last equality follows from (EC.3.1).

Then, it follows from (EC.3.2) that

$$\begin{aligned} \hat{\sigma}_t^2 &= \frac{1}{t - T_1} \|\hat{\boldsymbol{\epsilon}}_{T_1:t}\|^2 = \frac{1}{t - T_1} \|\boldsymbol{\epsilon}_{T_1:t} - \tilde{\mathbf{V}}_{T_1:t} \mathbf{M}_{T_1:t} \hat{\mathbf{G}}_{T_1:t}\|^2 \\ &= \frac{1}{t - T_1} \left( \|\boldsymbol{\epsilon}_{T_1:t}\|^2 - 2\boldsymbol{\epsilon}_{T_1:t}^\top \tilde{\mathbf{V}}_{T_1:t} \mathbf{M}_{T_1:t} \hat{\mathbf{G}}_{T_1:t} + \|\tilde{\mathbf{V}}_{T_1:t} \mathbf{M}_{T_1:t} \hat{\mathbf{G}}_{T_1:t}\|^2 \right) \\ &\leq \frac{1}{t - T_1} \left( \|\boldsymbol{\epsilon}_{T_1:t}\|^2 + 2\|\boldsymbol{\epsilon}_{T_1:t}\| \|\tilde{\mathbf{V}}_{T_1:t} \mathbf{M}_{T_1:t} \hat{\mathbf{G}}_{T_1:t}\| + \|\tilde{\mathbf{V}}_{T_1:t} \mathbf{M}_{T_1:t} \hat{\mathbf{G}}_{T_1:t}\|^2 \right). \end{aligned} \quad (\text{EC.3.18})$$

Therefore, to examine the event  $\mathcal{D}_t^c = \{|\hat{\sigma}_t^2 - \sigma^2| \leq \Delta\}$ , it suffices to analyze the following two events separately: (i)  $\frac{1}{t - T_1} \|\boldsymbol{\epsilon}_{T_1:t}\|^2$  is in the proximity of  $\sigma^2$ , and (ii)  $\frac{1}{t - T_1} \|\tilde{\mathbf{V}}_{T_1:t} \mathbf{M}_{T_1:t} \hat{\mathbf{G}}_{T_1:t}\|^2$  is small.

We first analyze the former. Define  $\mathcal{E}_t := \left\{ \left| \frac{1}{t - T_1} \|\boldsymbol{\epsilon}_{T_1:t}\|^2 - \sigma^2 \right| \leq \frac{\Delta}{2} \right\}$ . Note that

$$\frac{1}{t - T_1} \|\boldsymbol{\epsilon}_{T_1:t}\|^2 - \sigma^2 = \frac{1}{t - T_1} \sum_{s=T_1+1}^t (\epsilon_s^2 - \sigma^2),$$

and  $\mathbb{E}[\epsilon_s^2] = \sigma^2$ . Moreover, since  $\epsilon_s$  is sub-Gaussian, it follows from Lemma 2.7.6 of Vershynin (2018) that  $\epsilon_s^2 - \sigma^2$  is sub-exponential. Hence, by Lemma EC.6,

$$\begin{aligned} \Pr(\mathcal{E}_t^c) &\leq 2 \exp\left(- (t - T_1) C_{\text{Bern}} \min\left\{ \frac{\Delta^2}{4K_\epsilon^2}, \frac{\Delta}{2K_\epsilon} \right\}\right) \\ &\leq 2 \exp\left(- C_{T_2} \log(T) C_{\text{Bern}} \min\left\{ \frac{\Delta^2}{4K_\epsilon^2}, \frac{\Delta}{2K_\epsilon} \right\}\right) \\ &\leq 2 \exp(-3 \log(T)) = 2T^{-3}, \end{aligned} \quad (\text{EC.3.19})$$

where the last inequality holds because of the condition (EC.3.15).

Now, we analyze the term  $\frac{1}{t - T_1} \|\tilde{\mathbf{V}}_{T_1:t} \mathbf{M}_{T_1:t} \hat{\mathbf{G}}_{T_1:t}\|^2$ . Note that by Lemma EC.3,

$$\frac{1}{t - T_1} \|\tilde{\mathbf{V}}_{T_1:t}\|^2 \leq \frac{1}{t - T_1} (\sqrt{Mp(t - T_1)} \|\tilde{\mathbf{V}}_{T_1:t}\|_{\max})^2 \leq Mp\bar{v}^2. \quad (\text{EC.3.20})$$

It follows from (EC.3.11) that conditional on the event  $\mathcal{B}_t \cap \mathcal{C}_t$  with  $\eta_t$  satisfying (EC.3.16),

$$\|\mathbf{M}_{T_1:t}\|^2 = \left\| \left( \hat{\boldsymbol{\Sigma}}_{vz, T_1:t} \hat{\boldsymbol{\Sigma}}_{zz, T_1:t}^{-1} \hat{\boldsymbol{\Sigma}}_{vz, T_1:t}^\top \right)^{-1} \hat{\boldsymbol{\Sigma}}_{vz, T_1:t} \hat{\boldsymbol{\Sigma}}_{zz, T_1:t}^{-1} \right\|^2 \leq \bar{\kappa}^2. \quad (\text{EC.3.21})$$

Note that

$$\frac{1}{t - T_1} \|\tilde{\mathbf{V}}_{T_1:t} \mathbf{M}_{T_1:t} \hat{\mathbf{G}}_{T_1:t}\|^2 \leq \frac{1}{t - T_1} \|\tilde{\mathbf{V}}_{T_1:t}\|^2 \|\mathbf{M}_{T_1:t}\|^2 \|\hat{\mathbf{G}}_{T_1:t}\|^2.$$

It follows from (EC.3.20) and (EC.3.21) that conditional on the event  $\mathcal{B}_t \cap \mathcal{C}_t \cap \{\|\hat{\mathbf{G}}_{T_1:t}\| \leq \zeta_4\}$  with  $\zeta_4 := \sqrt{6q\bar{z}\zeta} \sqrt{\frac{\log(T)}{t - T_1}}$ , we have

$$\begin{aligned} \frac{1}{t - T_1} \|\tilde{\mathbf{V}}_{T_1:t} \mathbf{M}_{T_1:t} \hat{\mathbf{G}}_{T_1:t}\|^2 &\leq Mp\bar{v}^2 \zeta_4^2 \bar{\kappa}^2 = 6Mpq^2 \bar{v}^2 \bar{z}^2 \zeta^2 \bar{\kappa}^2 \frac{\log(T)}{t - T_1} \\ &\leq \frac{6Mpq^2 \bar{v}^2 \bar{z}^2 \zeta^2 \bar{\kappa}^2}{C_{T_2}} \leq \zeta_5 \end{aligned} \quad (\text{EC.3.22})$$

where  $\zeta_5 := \min\{\frac{\Delta^2}{128\sigma^2}, \frac{\Delta}{64}\}$ . Note that the second inequality holds because  $t - T_1 \geq T_2 - T_1 \geq C_{T_2} \log(T)$ , and the last inequality because of the condition (EC.3.15). Consequently,

$$\{\|\hat{\mathbf{G}}_{T_1:t}\| \leq \zeta_4\} \cap \mathcal{B}_t \cap \mathcal{C}_t \subseteq \left\{ \frac{1}{t - T_1} \|\tilde{\mathbf{V}}_{T_1:t} \mathbf{M}_{T_1:t} \hat{\mathbf{G}}_{T_1:t}\|^2 \leq \zeta_5 \right\} := \mathcal{F}_t. \quad (\text{EC.3.23})$$

Note that conditional on the event  $\mathcal{E}_t \cap \mathcal{F}_t$ , we have

$$\begin{aligned} |\hat{\sigma}_t^2 - \sigma^2| &\leq \left| \frac{1}{t - T_1} \|\epsilon_{T_1:t}\|^2 - \sigma^2 \right| + \left| \hat{\sigma}_t^2 - \frac{1}{t - T_1} \|\epsilon_{T_1:t}\|^2 \right| \\ &\leq \left| \frac{1}{t - T_1} \|\epsilon_{T_1:t}\|^2 - \sigma^2 \right| + 2 \left( \frac{1}{t - T_1} \|\epsilon_{T_1:t}\|^2 \cdot \frac{1}{t - T_1} \|\tilde{\mathbf{V}}_{T_1:t} \mathbf{M}_{T_1:t} \hat{\mathbf{G}}_{T_1:t}\|^2 \right)^{1/2} \\ &\quad + \frac{1}{t - T_1} \|\tilde{\mathbf{V}}_{T_1:t} \mathbf{M}_{T_1:t} \hat{\mathbf{G}}_{T_1:t}\|^2 \\ &\leq \frac{\Delta}{2} + 2\sqrt{\left(\sigma^2 + \frac{\Delta}{2}\right)\zeta_5} + \zeta_5, \end{aligned} \quad (\text{EC.3.24})$$

where the second inequality follows from (EC.3.18). Further, by the definition of  $\zeta_5$ , we have

$$\left(\sigma^2 + \frac{\Delta}{2}\right)\zeta_5 \leq \sigma^2 \cdot \frac{\Delta^2}{128\sigma^2} + \frac{\Delta}{2} \cdot \frac{\Delta}{64} = \frac{\Delta^2}{64}.$$

It then follows from (EC.3.24) that  $|\hat{\sigma}_t^2 - \sigma^2| \leq \frac{\Delta}{2} + \frac{\Delta}{4} + \frac{\Delta}{64} < \Delta$ . Hence,  $\mathcal{E}_t \cap \mathcal{F}_t \subseteq \mathcal{D}_t$ , which in conjunction with (EC.3.23) implies that

$$\mathcal{E}_t \cap \{\|\hat{\mathbf{G}}_{T_1:t}\| \leq \zeta_4\} \cap \mathcal{B}_t \cap \mathcal{C}_t \subseteq \mathcal{E}_t \cap \mathcal{F}_t \subseteq \mathcal{D}_t,$$

and thus  $\mathcal{E}_t \cap \{\|\hat{\mathbf{G}}_{T_1:t}\| \leq \zeta_4\} \cap \mathcal{B}_t \cap \mathcal{C}_t \subseteq \mathcal{B}_t \cap \mathcal{C}_t \cap \mathcal{D}_t$ . Hence,

$$\mathcal{B}_t \cap \mathcal{C}_t \cap \mathcal{D}_t^c \subseteq \left(\mathcal{E}_t \cap \{\|\hat{\mathbf{G}}_{T_1:t}\| \leq \zeta_4\}\right)^c \cap \mathcal{B}_t \cap \mathcal{C}_t \subseteq \mathcal{E}_t^c \cup \{\|\hat{\mathbf{G}}_{T_1:t}\| > \zeta_4\}. \quad (\text{EC.3.25})$$

Note that by (EC.3.13), we have

$$\Pr\left(\|\hat{\mathbf{G}}_{T_1:t}\| > \zeta_4\right) \leq 2q \exp\left(-\frac{\zeta_4^2(t - T_1)}{2q^2 \bar{z}^2 \zeta^2}\right) \leq 2qT^{-3}. \quad (\text{EC.3.26})$$

Combining (EC.3.19), (EC.3.25), and (EC.3.26) yields

$$\Pr(\mathcal{B}_t \cap \mathcal{C}_t \cap \mathcal{D}_t^c) \leq \Pr(\mathcal{E}_t^c) + \Pr\left(\|\hat{\mathbf{G}}_{T_1:t}\| > \zeta_4\right) \leq (2q + 2)T^{-3}. \quad \square$$

LEMMA EC.12. *Define*

$$\hat{\kappa} := \frac{2\phi_{\max}(\boldsymbol{\Sigma}_{zz})}{(\phi_{\min}(\boldsymbol{\Sigma}_{vz}^*))^2} \quad (\text{EC.3.27})$$

*Suppose*

$$C_{T_1} \geq 1024M^4p^2q^3\bar{v}^4\bar{z}^4\zeta^2\hat{\kappa}^2L^2\eta^{-2}, \quad (\text{EC.3.28})$$

$$C_{T_2} \geq \eta^{-2} \max \left\{ 128(M-1)^2pq\bar{v}^2\bar{z}^2, 256M^2L^2p^2q\bar{v}^4\bar{z}^2 \left( C_A + \sqrt{\theta\hat{\kappa}(\sigma^2 + \Delta)} \right)^2 \right\}. \quad (\text{EC.3.29})$$

*Suppose that the constant  $\eta$  in the definitions of  $\mathcal{B}_t$  and  $\mathcal{C}_t$  satisfies (EC.3.16). Then, for all  $t = T_2 + 1, \dots, T$ ,*

$$\Pr(\mathcal{B}_t^c \cap \mathcal{A}_{T_1} \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1} \cap \mathcal{D}_{T_2:t-1}) \leq 2(Mp + q + 1)T^{-2}.$$

*Proof.* A subtlety in analyzing the event  $\mathcal{B}_t^c = \{\|\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t} - \boldsymbol{\Sigma}_{vz}^*\| > \eta\}$  for  $t > T_1$  is that  $\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}$  is not an unbiased estimator of  $\boldsymbol{\Sigma}_{vz}^*$ —that is,  $\mathbb{E}[\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}] \neq \boldsymbol{\Sigma}_{vz}^*$ . This is because  $\boldsymbol{\Sigma}_{vz}^*$  is defined with respect to the (unknown) optimal policy (i.e., the optimal arm  $a_s^*$  is taken for all  $s = T_1 + 1, \dots, t$ ); whereas  $\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}$  is defined with respect to the actual policy that is used to generate the data. To address the issue, we define

$$\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}^* := \frac{1}{t - T_1} \sum_{s=T_1+1}^t \mathbf{v}_s^* \mathbf{z}_s^T \in \mathbb{R}^{(Mp) \times q},$$

and analyze separately the two components of the following decomposition:

$$\|\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t} - \boldsymbol{\Sigma}_{vz}^*\| \leq \underbrace{\|\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t} - \widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}^*\|}_{I_1} + \underbrace{\|\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}^* - \boldsymbol{\Sigma}_{vz}^*\|}_{I_2}. \quad (\text{EC.3.30})$$

For  $I_2$ , note that  $\mathbb{E}[\mathbf{v}_t^* \mathbf{z}_t^T] = \boldsymbol{\Sigma}_{vz}^*$  and that, by Lemma EC.3,

$$\|\mathbf{v}_t^* \mathbf{z}_t^T - \boldsymbol{\Sigma}_{vz}^*\| \leq \|\mathbf{v}_t^* \mathbf{z}_t^T\| + \|\boldsymbol{\Sigma}_{vz}^*\| \leq 2\sqrt{Mpq\bar{v}\bar{z}}.$$

It then follows from Lemma EC.4 that for all  $t \geq T_2$ ,

$$\begin{aligned} \Pr\left(\|\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}^* - \boldsymbol{\Sigma}_{vz}^*\| > \frac{\eta}{2}\right) &\leq 2(Mp + q) \exp\left(-\frac{(t - T_1)\eta^2/4}{2(2\sqrt{Mpq\bar{v}\bar{z}})^2}\right) \\ &\leq 2(Mp + q) \exp\left(-\frac{\eta^2 C_{T_2} \log(T)}{32Mpq\bar{v}^2\bar{z}^2}\right) \\ &\leq 2(Mp + q) \exp(-2\log(T)) = 2(Mp + q)T^{-2}, \end{aligned} \quad (\text{EC.3.31})$$

where the second inequality holds because  $t - T_1 \geq T_2 - T_1 = C_{T_2} \log(T)$ , and the third holds because of the condition (EC.3.29).

We now consider  $I_1$ . Note that

$$\begin{aligned} \|\widehat{\Sigma}_{vz, T_1:t} - \widehat{\Sigma}_{vz, T_1:t}^*\| &= \left\| \frac{1}{t - T_1} \sum_{s=T_1+1}^t \begin{pmatrix} (\mathbb{I}(a_s = 1) - \mathbb{I}(a_s^* = 1))\mathbf{v}_s \\ \vdots \\ (\mathbb{I}(a_s = M) - \mathbb{I}(a_s^* = M))\mathbf{v}_s \end{pmatrix} \mathbf{z}_s^\top \right\| \\ &\leq \frac{2}{t - T_1} \sum_{s=T_1+1}^t \mathbb{I}(a_s \neq a_s^*) \|\mathbf{v}_s \mathbf{z}_s^\top\| \leq \frac{2\sqrt{pq\bar{v}\bar{z}}}{t - T_1} \sum_{s=T_1+1}^t \mathbb{I}(a_s \neq a_s^*), \end{aligned} \quad (\text{EC.3.32})$$

where the first inequality holds because there are at most non-zero entries in the vector  $(\mathbb{I}(a_s = 1) - \mathbb{I}(a_s^* = 1), \dots, \mathbb{I}(a_s = M) - \mathbb{I}(a_s^* = M))^\top$ ; the last inequality follows from Lemma EC.3. We show next that conditional on some other events,

$$\{a_s = j\} \subseteq \left\{ |\mathbf{v}_s^\top \boldsymbol{\delta}_{j, a_s^*}| \leq \frac{\eta}{8ML\sqrt{pq\bar{v}\bar{z}}} \right\}, \quad \forall s \geq T_1 + 1, \forall j \neq a_s^*,$$

where  $\boldsymbol{\delta}_{i,j} = \boldsymbol{\alpha}_i - \boldsymbol{\alpha}_j$  and  $L$  is the constant defined in the margin condition of Assumption 1. We consider two cases: (i)  $s = T_1 + 1, \dots, T_2$  and (ii)  $s = T_2 + 1, \dots, T$ . Let  $\widehat{\boldsymbol{\delta}}_{i,j,s} := \widehat{\boldsymbol{\alpha}}_{i,s} - \widehat{\boldsymbol{\alpha}}_{j,s}$  for each  $s$ .

**Case (i):**  $s = T_1 + 1, \dots, T_2$ . Recall that for  $s = T_1 + 1, \dots, T_2$ , the selected arm is  $a_s = \arg \max_{i=1, \dots, M} \{\mathbf{v}_s^\top \widehat{\boldsymbol{\alpha}}_{i, T_1}\}$ . Hence, for any  $j \neq a_s^*$ , if  $a_s = j$ , then  $\mathbf{v}_s^\top \widehat{\boldsymbol{\delta}}_{j, a_s^*, T_1} = \mathbf{v}_s^\top \widehat{\boldsymbol{\alpha}}_{j, T_1} - \mathbf{v}_s^\top \widehat{\boldsymbol{\alpha}}_{a_s^*, T_1} > 0$ . Thus,

$$\begin{aligned} 0 &\geq \mathbf{v}_s^\top \boldsymbol{\delta}_{j, a_s^*} > \mathbf{v}_s^\top \boldsymbol{\delta}_{j, a_s^*} - \mathbf{v}_s^\top \widehat{\boldsymbol{\delta}}_{j, a_s^*, T_1} = \mathbf{v}_s^\top [(\boldsymbol{\alpha}_j - \widehat{\boldsymbol{\alpha}}_{j, T_1}) - (\boldsymbol{\alpha}_{a_s^*} - \widehat{\boldsymbol{\alpha}}_{a_s^*, T_1})] \\ &\geq -2\|\mathbf{v}_s\| \|\boldsymbol{\alpha} - \widehat{\boldsymbol{\alpha}}_{T_1}\| \geq -2\sqrt{p\bar{v}} \|\boldsymbol{\alpha} - \widehat{\boldsymbol{\alpha}}_{T_1}\|. \end{aligned}$$

Further, conditional on  $\mathcal{A}_{T_1} = \{\|\widehat{\boldsymbol{\alpha}}_{T_1} - \boldsymbol{\alpha}\| \leq 2Mq\bar{z}\check{\kappa}C_{T_1}^{-\frac{1}{2}}\}$  defined in (EC.2.2), we have

$$|\mathbf{v}_s^\top \boldsymbol{\delta}_{j, a_s^*}| \leq 2\sqrt{p\bar{v}} \|\boldsymbol{\alpha} - \widehat{\boldsymbol{\alpha}}_{T_1}\| \leq 4M\sqrt{pq\bar{v}\bar{z}}\check{\kappa}C_{T_1}^{-\frac{1}{2}} \leq \frac{\eta}{8ML\sqrt{pq\bar{v}\bar{z}}},$$

where the last inequality follows from the condition (EC.3.28). Hence,

$$\mathcal{A}_{T_1} \cap \{a_s = j\} \subseteq \left\{ |\mathbf{v}_s^\top \boldsymbol{\delta}_{j, a_s^*}| \leq \frac{\eta}{8ML\sqrt{pq\bar{v}\bar{z}}} \right\}, \quad \forall s = T_1 + 1, \dots, T_2, \forall j \neq a_s^*. \quad (\text{EC.3.33})$$

**Case (ii):**  $s = T_2 + 1, \dots, T$ . Recall that for  $s \geq T_2 + 1$ , the arm selected at time  $s$  is

$$a_s = \arg \max_{i=1, \dots, M} \left\{ \mathbf{v}_s^\top \widehat{\boldsymbol{\alpha}}_{i, s-1} + \widehat{\sigma}_{s-1} \sqrt{2\theta \log(s - T_1) \cdot \mathbf{v}_s^\top \widehat{\boldsymbol{\Omega}}_{i, s-1} \mathbf{v}_s} \right\}$$

where  $\widehat{\boldsymbol{\Omega}}_{i, s-1} \in \mathbb{R}^{p \times p}$  is the  $i$ -th diagonal  $p \times p$  block of  $\widehat{\boldsymbol{\Omega}}_{s-1} \in \mathbb{R}^{(Mp) \times (Mp)}$  defined as

$$\begin{aligned} \widehat{\boldsymbol{\Omega}}_{s-1} &= (\widetilde{\mathbf{V}}_{T_1:s-1}^\top \mathcal{P}[\mathbf{Z}_{T_1:s-1}] \widetilde{\mathbf{V}}_{T_1:s-1})^{-1} \\ &= \frac{1}{s - 1 - T_1} \left( \widehat{\Sigma}_{vz, T_1:s-1} \widehat{\Sigma}_{zz, T_1:s-1}^{-1} \widehat{\Sigma}_{vz, T_1:s-1}^\top \right)^{-1}. \end{aligned}$$

Thus, for any  $j \neq a_s^*$ , if  $a_s = j$ , then

$$\mathbf{v}_s^\top \widehat{\boldsymbol{\delta}}_{j,a_s^*,s-1} + \widehat{\sigma}_{s-1} \sqrt{2\theta \log(s - T_1)} \left( \sqrt{\mathbf{v}_s^\top \widehat{\boldsymbol{\Omega}}_{j,s-1} \mathbf{v}_s} - \sqrt{\mathbf{v}_s^\top \widehat{\boldsymbol{\Omega}}_{a_s^*,s-1} \mathbf{v}_s} \right) > 0.$$

Since  $\mathbf{v}_s^\top \boldsymbol{\delta}_{j,a_s^*} \leq 0$ , we have

$$\begin{aligned} 0 &\geq \mathbf{v}_s^\top \boldsymbol{\delta}_{j,a_s^*} > \underbrace{\mathbf{v}_s^\top (\boldsymbol{\delta}_{j,a_s^*} - \widehat{\boldsymbol{\delta}}_{j,a_s^*,s-1})}_{W_{1,j}} \\ &\quad - \underbrace{\widehat{\sigma}_{s-1} \sqrt{2\theta \log(s - T_1)} \left( \sqrt{\mathbf{v}_s^\top \widehat{\boldsymbol{\Omega}}_{j,s-1} \mathbf{v}_s} - \sqrt{\mathbf{v}_s^\top \widehat{\boldsymbol{\Omega}}_{a_s^*,s-1} \mathbf{v}_s} \right)}_{W_{2,j}}. \end{aligned} \quad (\text{EC.3.34})$$

For  $W_{1,j}$ , note that conditional on  $\mathcal{A}_{s-1} = \left\{ \|\widehat{\boldsymbol{\alpha}}_{s-1} - \boldsymbol{\alpha}\| \leq C_{\mathcal{A}} \sqrt{\frac{\log(T)}{s-1-T_1}} \right\}$ , we have

$$\begin{aligned} |W_{1,j}| &= |\mathbf{v}_s^\top [(\boldsymbol{\alpha}_j - \widehat{\boldsymbol{\alpha}}_{j,s-1}) - (\boldsymbol{\alpha}_{a_s^*} - \widehat{\boldsymbol{\alpha}}_{a_s^*,s-1})]| \leq 2 \|\mathbf{v}_s\| \|\boldsymbol{\alpha} - \widehat{\boldsymbol{\alpha}}_{s-1}\| \\ &\leq 2\sqrt{p\bar{v}} C_{\mathcal{A}} \sqrt{\frac{\log(T)}{s-1-T_1}} \end{aligned} \quad (\text{EC.3.35})$$

$$\leq 2\sqrt{p\bar{v}} C_{\mathcal{A}} C_{T_2}^{-\frac{1}{2}}. \quad (\text{EC.3.36})$$

For  $W_{2,j}$ , note that

$$\begin{aligned} &\left( \sqrt{\mathbf{v}_s^\top \widehat{\boldsymbol{\Omega}}_{j,s-1} \mathbf{v}_s} - \sqrt{\mathbf{v}_s^\top \widehat{\boldsymbol{\Omega}}_{a_s^*,s-1} \mathbf{v}_s} \right)^2 \\ &\leq \left| \sqrt{\mathbf{v}_s^\top \widehat{\boldsymbol{\Omega}}_{j,s-1} \mathbf{v}_s} - \sqrt{\mathbf{v}_s^\top \widehat{\boldsymbol{\Omega}}_{a_s^*,s-1} \mathbf{v}_s} \right| \cdot \left| \sqrt{\mathbf{v}_s^\top \widehat{\boldsymbol{\Omega}}_{j,s-1} \mathbf{v}_s} + \sqrt{\mathbf{v}_s^\top \widehat{\boldsymbol{\Omega}}_{a_s^*,s-1} \mathbf{v}_s} \right| \\ &\leq |\mathbf{v}_s^\top (\widehat{\boldsymbol{\Omega}}_{j,s-1} - \widehat{\boldsymbol{\Omega}}_{a_s^*,s-1}) \mathbf{v}_s| \\ &\leq \|\mathbf{v}_s\|^2 \|\widehat{\boldsymbol{\Omega}}_{j,s-1} - \widehat{\boldsymbol{\Omega}}_{a_s^*,s-1}\| \leq \|\mathbf{v}_s\|^2 (\|\widehat{\boldsymbol{\Omega}}_{j,s-1}\| + \|\widehat{\boldsymbol{\Omega}}_{a_s^*,s-1}\|) \\ &\leq 2\|\mathbf{v}_s\|^2 \|\widehat{\boldsymbol{\Omega}}_{s-1}\| \leq 2p\bar{v}^2 \|\widehat{\boldsymbol{\Omega}}_{s-1}\|, \end{aligned} \quad (\text{EC.3.37})$$

where the second to the last inequality holds because  $\widehat{\boldsymbol{\Omega}}_{i,s-1}$  is a submatrix of  $\widehat{\boldsymbol{\Omega}}_{s-1}$ , so the norm of the former is upper bounded by that of the latter; see Golub and Van Loan (2013, page 72). Further, conditional on the event  $\mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} = \left\{ \|\widehat{\boldsymbol{\Sigma}}_{vz, T_1:s-1} - \boldsymbol{\Sigma}_{vz}^*\| \leq \eta \right\} \cap \left\{ \|\widehat{\boldsymbol{\Sigma}}_{zz, T_1:s-1} - \boldsymbol{\Sigma}_{zz}\| \leq \eta \right\}$ , we have

$$\begin{aligned} (s-1-T_1) \|\widehat{\boldsymbol{\Omega}}_{s-1}\| &= \phi_{\max}((s-1-T_1) \widehat{\boldsymbol{\Omega}}_{s-1}) \\ &= \phi_{\max} \left( \left( \widehat{\boldsymbol{\Sigma}}_{vz, T_1:s-1} \widehat{\boldsymbol{\Sigma}}_{zz, T_1:s-1}^{-1} \widehat{\boldsymbol{\Sigma}}_{vz, T_1:s-1}^\top \right)^{-1} \right) \\ &\leq \frac{\phi_{\max}(\widehat{\boldsymbol{\Sigma}}_{zz, T_1:s-1})}{(\phi_{\min}(\widehat{\boldsymbol{\Sigma}}_{vz, T_1:s-1}))^2} \leq \frac{\phi_{\max}(\boldsymbol{\Sigma}_{zz}) + \eta}{(\phi_{\min}(\boldsymbol{\Sigma}_{vz}^*) - \eta)^2} \\ &\leq \left( \frac{11}{10} \right) \left( \frac{10}{9} \right)^2 \frac{\phi_{\max}(\boldsymbol{\Sigma}_{zz})}{(\phi_{\min}(\boldsymbol{\Sigma}_{vz}^*))^2} \leq \mathring{\kappa}, \end{aligned} \quad (\text{EC.3.38})$$

where the first inequality follows from Lemma EC.2, the third from the condition (EC.3.16), and the last from the definition of  $\hat{\kappa}$  in (EC.3.27). Plugging (EC.3.38) in (EC.3.37) yields

$$\mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} \subseteq \left\{ \sqrt{\mathbf{v}_s^\top \hat{\boldsymbol{\Omega}}_{j,s-1} \mathbf{v}_s} - \sqrt{\mathbf{v}_s^\top \hat{\boldsymbol{\Omega}}_{a_s^*,s-1} \mathbf{v}_s} \leq \sqrt{\frac{2p\bar{v}^2 \hat{\kappa}}{s-1-T_1}} \right\} \quad (\text{EC.3.39})$$

In addition, note that

$$\mathcal{D}_{s-1} = \{|\hat{\sigma}_{s-1}^2 - \sigma^2| \leq \Delta\} \subseteq \{|\hat{\sigma}_{s-1}^2| \leq \sigma^2 + \Delta\}. \quad (\text{EC.3.40})$$

Plugging (EC.3.39) and (EC.3.40) into the definition of  $W_{2,j}$  in (EC.3.34), we have that conditional on  $\mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} \cap \mathcal{D}_{s-1}$ ,

$$\begin{aligned} |W_{2,j}| &\leq \sqrt{\sigma^2 + \Delta} \cdot \sqrt{2\theta \log(s-T_1)} \cdot \sqrt{\frac{2p\bar{v}^2 \hat{\kappa}}{s-1-T_1}} \\ &\leq 2\sqrt{\theta p \bar{v}^2 \hat{\kappa} (\sigma^2 + \Delta)} \sqrt{\frac{\log(T)}{s-1-T_1}} \end{aligned} \quad (\text{EC.3.41})$$

$$\leq 2\sqrt{\theta p \bar{v}^2 \hat{\kappa} (\sigma^2 + \Delta)} C_{T_2}^{-\frac{1}{2}}, \quad (\text{EC.3.42})$$

where the last inequality holds because  $s-1-T_1 \geq T_2 - T_1 = C_{T_2} \log(T)$ .

Plugging (EC.3.36) and (EC.3.42) into (EC.3.34) yields that conditional on  $\mathcal{A}_{s-1} \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} \cap \mathcal{D}_{s-1}$ ,

$$|\mathbf{v}_s^\top \boldsymbol{\delta}_{j,a_s^*}| \leq |W_{1,j}| + |W_{2,j}| \leq 2\sqrt{p\bar{v}} \left( C_{\mathcal{A}} + \sqrt{\theta \hat{\kappa} (\sigma^2 + \Delta)} \right) C_{T_2}^{-\frac{1}{2}} \leq \frac{\eta}{8ML\sqrt{pq\bar{v}\bar{z}}},$$

where the second inequality holds because of the condition (EC.3.29). Therefore,

$$\begin{aligned} &\mathcal{A}_{s-1} \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} \cap \mathcal{D}_{s-1} \cap \{a_s = j\} \\ &\subseteq \left\{ |\mathbf{v}_s^\top \boldsymbol{\delta}_{j,a_s^*}| \leq \frac{\eta}{8ML\sqrt{pq\bar{v}\bar{z}}} \right\}, \quad \forall s = T_2 + 1, \dots, T, \forall j \neq a_s^*. \end{aligned} \quad (\text{EC.3.43})$$

Consequently, combining (EC.3.33) and (EC.3.43) yields that

$$\begin{aligned} &\mathcal{A}_{T_1} \cap \mathcal{A}_{s-1} \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} \cap \mathcal{D}_{s-1} \cap \{a_s = j\} \\ &\subseteq \left\{ |\mathbf{v}_s^\top \boldsymbol{\delta}_{j,a_s^*}| \leq \frac{\eta}{8ML\sqrt{pq\bar{v}\bar{z}}} \right\}, \quad \forall s = T_1 + 1, \dots, T, \forall j \neq a_s^*. \end{aligned}$$

Now, let us return the task of bounding  $I_1$ , which is defined in (EC.3.30). By (EC.3.32), conditional on  $\mathcal{A}_{T_1} \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1} \cap \mathcal{D}_{T_2:t-1}$ , we have

$$\|\hat{\boldsymbol{\Sigma}}_{vz, T_1:t} - \hat{\boldsymbol{\Sigma}}_{vz, T_1:t}^*\| \leq \frac{2\sqrt{pq\bar{v}\bar{z}}}{t-T_1} \sum_{s=T_1+1}^t \mathbb{I}(a_s \neq a_s^*) \leq \frac{2\sqrt{pq\bar{v}\bar{z}}}{t-T_1} \sum_{s=T_1+1}^t \sum_{j \neq a_s^*} \xi_{j,s},$$

for all  $t = T_2 + 1, \dots, T$ , where  $\xi_{j,s} := \mathbb{I}(|\mathbf{v}_s^\top \boldsymbol{\delta}_{j,a_s^*}| \leq \frac{\eta}{8ML\sqrt{pq\bar{v}\bar{z}}})$ . Hence,

$$\begin{aligned}
& \Pr\left(\left\{\|\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t} - \widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}^*\| \geq \frac{\eta}{2}\right\} \cap \mathcal{A}_{T_1} \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1} \cap \mathcal{D}_{T_2:t-1}\right) \\
& \leq \Pr\left(\frac{1}{t-T_1} \sum_{s=T_1+1}^t \sum_{j \neq a_s^*} \xi_{j,s} \geq \frac{\eta}{4\sqrt{pq\bar{v}\bar{z}}}\right) \\
& \leq \Pr\left(\left|\frac{1}{t-T_1} \sum_{s=T_1+1}^t \sum_{j \neq a_s^*} (\xi_{j,s} - \mathbb{E}[\xi_{j,s}])\right| \geq \frac{\eta}{4\sqrt{pq\bar{v}\bar{z}}} - \sum_{j \neq a_s^*} \mathbb{E}[\xi_{j,s}]\right) \\
& \leq \Pr\left(\left|\frac{1}{t-T_1} \sum_{s=T_1+1}^t \sum_{j \neq a_s^*} (\xi_{j,s} - \mathbb{E}[\xi_{j,s}])\right| \geq \frac{\eta}{8\sqrt{pq\bar{v}\bar{z}}}\right) \tag{EC.3.44}
\end{aligned}$$

where the last step holds because

$$\sum_{j \neq a_s^*} \mathbb{E}[\xi_{j,s}] = \sum_{j \neq a_s^*} \Pr\left(|\mathbf{v}_s^\top \boldsymbol{\delta}_{j,a_s^*}| \leq \frac{\eta}{8ML\sqrt{pq\bar{v}\bar{z}}}\right) \leq \sum_{j \neq a_s^*} \frac{\eta}{8M\sqrt{pq\bar{v}\bar{z}}} \leq \frac{\eta}{8\sqrt{pq\bar{v}\bar{z}}},$$

by the margin condition in Assumption 1.

Note that  $\{\sum_{j \neq a_s^*} \xi_{j,s} : s = T_1 + 1, \dots, T\}$  are i.i.d. bounded random variables taking values on  $[0, M-1]$ . Thus, they are sub-Gaussian with variance proxy  $\frac{M-1}{4}$ ; see, e.g., Wainwright (2019, page 24). Then, setting  $\zeta_6 := \frac{\eta}{8\sqrt{pq\bar{v}\bar{z}}}$  and applying Lemma EC.5 leads to

$$\begin{aligned}
\Pr\left(\left|\frac{1}{t-T_1} \sum_{s=T_1+1}^t (\xi_s - \mathbb{E}[\xi_s])\right| \geq \zeta_6\right) & \leq 2 \exp\left(\frac{-2\zeta_6^2(t-T_1)}{(M-1)^2}\right) \\
& \leq 2 \exp\left(\frac{-\eta^2 C_{T_2} \log(T)}{32(M-1)^2 pq\bar{v}\bar{z}^2}\right) \leq 2T^{-4}, \tag{EC.3.45}
\end{aligned}$$

for all  $t = T_2 + 1, \dots, T$ , where the last inequality follows from the condition (EC.3.29). It then follows from (EC.3.44) that

$$\begin{aligned}
& \Pr\left(\left\{\|\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t} - \widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}^*\| \geq \frac{\eta}{2}\right\} \cap \mathcal{A}_{T_1} \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1} \cap \mathcal{D}_{T_2:t-1}\right) \\
& \leq 2T^{-4}, \quad \forall t = T_2 + 1, \dots, T. \tag{EC.3.46}
\end{aligned}$$

Combining (EC.3.30), (EC.3.31), and (EC.3.46) results in

$$\begin{aligned}
& \Pr\left(\left\{\|\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t} - \boldsymbol{\Sigma}_{vz}^*\| \geq \eta\right\} \cap \mathcal{A}_{T_1} \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1} \cap \mathcal{D}_{T_2:t-1}\right) \\
& \leq \Pr\left(\left\{\|\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t} - \widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}^*\| \geq \frac{\eta}{2}\right\} \cap \mathcal{A}_{T_1} \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1} \cap \mathcal{D}_{T_2:t-1}\right) \\
& \quad + \Pr\left(\|\widehat{\boldsymbol{\Sigma}}_{vz, T_1:t}^* - \boldsymbol{\Sigma}_{vz}^*\| > \frac{\eta}{2}\right) \\
& \leq 2(Mp+q)T^{-2} + 2T^{-4} \leq 2(Mp+q+1)T^{-2},
\end{aligned}$$

for all  $t = T_2 + 1, \dots, T$ , which concludes the proof.  $\square$

We now summarize the conditions of Lemmas EC.8–EC.12 as follows.

ASSUMPTION EC.1. *The constants  $C_{T_1}$  and  $C_{T_2}$  satisfy*

$$C_{T_1} \geq \max \left\{ \left( \frac{40M \max\{\sqrt{pq\bar{v}\bar{z}}, q\bar{z}^2\}}{\min\{\phi_{\min}(\check{\Sigma}_{vz}), \phi_{\min}(\check{\Sigma}_{zz})\}} \right)^2, 1024M^4 p^2 q^3 \bar{v}^4 \bar{z}^4 \check{\zeta}^2 \check{\kappa}^2 L^2 \eta^{-2} \right\},$$

$$C_{T_2} \geq \max \left\{ \frac{3}{C_{\text{Bern}} \min\left\{\frac{\Delta^2}{4K_\epsilon^2}, \frac{\Delta}{2K_\epsilon}\right\}}, \frac{6Mpq^2 \bar{v}^2 \bar{z}^2 \zeta^2 \bar{\kappa}^2}{\min\left\{\frac{\Delta^2}{128\sigma^2}, \frac{\Delta}{64}\right\}}, 6q^2 \bar{z}^4 \eta^{-2}, \right.$$

$$\left. 128(M-1)^2 pq \bar{v}^2 \bar{z}^2 \eta^{-2}, 256M^2 L^2 p^2 q \bar{v}^4 \bar{z}^2 \left(C_{\mathcal{A}} + \sqrt{\theta_{\check{\kappa}}(\sigma^2 + \Delta)}\right)^2 \eta^{-2} \right\},$$

where  $\check{\kappa}$ ,  $\bar{\kappa}$ , and  $\hat{\kappa}$  are defined in (EC.2.1), (EC.3.3), and (EC.3.27), respectively;  $C_{\text{Bern}}$  and  $K_\epsilon$  are constants defined in Lemma EC.11,  $C_{\mathcal{A}} \geq 2q\bar{z}\bar{\kappa}$ , and  $\eta \leq \frac{1}{10} \min\{\phi_{\min}(\Sigma_{vz}^*), \phi_{\min}(\Sigma_{zz})\}$ .

LEMMA EC.13. *Suppose  $C_{T_1}$  and  $C_{T_2}$  satisfy Assumption EC.1. Then, for all  $t = T_2 + 1, \dots, T$ ,*

$$\Pr(\mathcal{A}_{T_2:t}) \geq 1 - [6Mp + (10M + 10)q + 3]T^{-1},$$

$$\Pr(\mathcal{B}_{T_2:t}) \geq 1 - [6Mp + (10M + 10)q + 3]T^{-1}.$$

*Proof.* We first summarize the results from Lemmas EC.8–EC.12: for all  $t = T_2 + 1, \dots, T$ ,

$$\Pr(\mathcal{A}_{T_1}) \geq 1 - 2M(2p + 5q)T^{-2}, \tag{EC.3.47}$$

$$\Pr(\mathcal{A}_t^c \cap \mathcal{B}_t \cap \mathcal{C}_t) \leq 2qT^{-2}, \tag{EC.3.48}$$

$$\Pr(\mathcal{C}_t) \geq 1 - 2qT^{-3}, \tag{EC.3.49}$$

$$\Pr(\mathcal{C}_{T_2:t}) \geq 1 - 2q(t - T_2)T^{-3} \geq 1 - 2qT^{-2} \tag{EC.3.50}$$

$$\Pr(\mathcal{B}_t \cap \mathcal{C}_t \cap \mathcal{D}_t^c) \leq (2q + 2)T^{-3}, \tag{EC.3.51}$$

$$\Pr(\mathcal{B}_t^c \cap \mathcal{A}_{T_1} \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1} \cap \mathcal{D}_{T_2:t-1}) \leq 2(Mp + q + 1)T^{-2}. \tag{EC.3.52}$$

By (EC.3.48) and (EC.3.49), we have

$$\begin{aligned} \Pr(\mathcal{A}_t^c \cap \mathcal{B}_t) &= \Pr(\mathcal{A}_t^c \cap \mathcal{B}_t \cap \mathcal{C}_t) + \Pr(\mathcal{A}_t^c \cap \mathcal{B}_t \cap \mathcal{C}_t^c) \\ &\leq \Pr(\mathcal{A}_t^c \cap \mathcal{B}_t \cap \mathcal{C}_t) + \Pr(\mathcal{C}_t^c) \leq 2q(T^{-2} + T^{-3}) \leq 4qT^{-2}. \end{aligned} \tag{EC.3.53}$$

By applying (EC.3.51) for  $s = T_2 + 1, \dots, t$ , we have

$$\begin{aligned} \Pr(\mathcal{B}_{T_2:t} \cap \mathcal{C}_{T_2:t} \cap \mathcal{D}_{T_2:t}^c) &= \Pr(\mathcal{B}_{T_2:t} \cap \mathcal{C}_{T_2:t} \cap (\cup_{s=T_2+1}^t \mathcal{D}_s^c)) \\ &\leq \sum_{s=T_2+1}^t \Pr(\mathcal{B}_{T_2:t} \cap \mathcal{C}_{T_2:t} \cap \mathcal{D}_s^c) \leq \sum_{s=T_2+1}^t \Pr(\mathcal{B}_s \cap \mathcal{C}_s \cap \mathcal{D}_s^c) \end{aligned}$$

$$\leq (t - T_2)(2q + 2)T^{-3} \leq (2q + 2)T^{-2}. \quad (\text{EC.3.54})$$

By (EC.3.47) and (EC.3.52), we have

$$\begin{aligned} & \Pr(\mathcal{B}_t^c \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1} \cap \mathcal{D}_{T_2:t-1}) \\ & \leq \Pr(\mathcal{B}_t^c \cap \mathcal{A}_{T_1} \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1} \cap \mathcal{D}_{T_2:t-1}) + \Pr(\mathcal{A}_{T_1}^c) \\ & \leq 2(Mp + q + 1)T^{-2} + 2M(2p + 5q)T^{-2} = [6Mp + (10M + 2)q + 1]T^{-2}. \end{aligned} \quad (\text{EC.3.55})$$

Define  $\mathcal{G}_t := \mathcal{A}_{T_2:t} \cap \mathcal{B}_{T_2:t}$  for all  $t = T_2 + 1, \dots, T$ . Let  $\Pr(\mathcal{G}_{T_2}) = 1$  for convenience. Moreover, note that

$$\begin{aligned} & \Pr(\mathcal{B}_t^c \cap \mathcal{G}_{t-1}) = \Pr(\mathcal{B}_t^c \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1}) \\ & = \Pr(\mathcal{B}_t^c \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1}) + \Pr(\mathcal{B}_t^c \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1}^c) \\ & \leq \Pr(\mathcal{B}_t^c \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1} \cap \mathcal{D}_{T_2:t-1}) \\ & \quad + \Pr(\mathcal{B}_t^c \cap \mathcal{A}_{T_2:t-1} \cap \mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1} \cap \mathcal{D}_{T_2:t-1}^c) + \Pr(\mathcal{C}_{T_2:t-1}^c) \\ & \leq [6Mp + (10M + 2)q + 1]T^{-2} + \Pr(\mathcal{B}_{T_2:t-1} \cap \mathcal{C}_{T_2:t-1} \cap \mathcal{D}_{T_2:t-1}^c) + 2qT^{-2} \\ & \leq [6Mp + (10M + 6)q + 3]T^{-2}, \end{aligned} \quad (\text{EC.3.56})$$

where the second inequality follows from (EC.3.50) and (EC.3.55), and the third from (EC.3.54).

Therefore,

$$\begin{aligned} \Pr(\mathcal{G}_t) & = \Pr(\mathcal{G}_{t-1} \cap \mathcal{A}_t \cap \mathcal{B}_t) = \Pr(\mathcal{G}_{t-1}) - \Pr(\mathcal{G}_{t-1} \cap (\mathcal{A}_t^c \cup \mathcal{B}_t^c)) \\ & = \Pr(\mathcal{G}_{t-1}) - \Pr(\mathcal{G}_{t-1} \cap \mathcal{A}_t^c \cap \mathcal{B}_t) - \Pr(\mathcal{G}_{t-1} \cap \mathcal{B}_t^c) \\ & \geq \Pr(\mathcal{G}_{t-1}) - \Pr(\mathcal{A}_t^c \cap \mathcal{B}_t) - \Pr(\mathcal{G}_{t-1} \cap \mathcal{B}_t^c) \\ & \geq \Pr(\mathcal{G}_{t-1}) - [10Mp + (10M + 10)q + 3]T^{-2}, \end{aligned} \quad (\text{EC.3.57})$$

where the last inequality follows from (EC.3.53) and (EC.3.56).

Note that  $\Pr(\mathcal{G}_{T_2}) = 1$ , applying (EC.3.57) iteratively, we have that for all  $t > T_2$ ,

$$\Pr(\mathcal{G}_t) \geq 1 - (t - T_2)[10Mp + (10M + 10)q + 3]T^{-2} \geq 1 - [10Mp + (10M + 10)q + 3]T^{-1}.$$

The proof is concluded by noting that  $\Pr(\mathcal{A}_{T_2:t}) \geq \Pr(\mathcal{G}_t)$  and  $\Pr(\mathcal{B}_{T_2:t}) \geq \Pr(\mathcal{G}_t)$ .  $\square$

## EC.4. Regret Analysis for IV-Greedy

We state a more specific version of Theorem 1 as follows.

**THEOREM EC.1.** *Let  $T_1 = C_{T_1} \log(T)$  and  $T_2 = (C_{T_1} + C_{T_2}) \log(T)$  for some constants  $C_{T_1}$  and  $C_{T_2}$  that satisfy Assumption EC.1. Let  $\bar{\delta} := \max_{1 \leq i \neq j \leq M} \|\boldsymbol{\delta}_{i,j}\|$ , where  $\boldsymbol{\delta}_{i,j} := \boldsymbol{\alpha}_i - \boldsymbol{\alpha}_j$ . Then, under Assumption 1, the regret of IV-Greedy satisfies that, for all  $T > T_2$ ,*

$$\begin{aligned} \text{Regret}(T) &\leq [\sqrt{p\bar{v}}\bar{\delta}(C_{T_1} + C_{T_2}) + 16(M-1)Lpq^3\bar{v}^2\bar{z}^2\zeta^2\bar{\kappa}^2] \log(T) \\ &\quad + \sqrt{p\bar{v}}\bar{\delta}(M-1)[12Mp + (20M+24)q + 7]. \end{aligned}$$

*Proof of Theorem EC.1.* Note that for all  $s \geq 1$ ,

$$R_s^{\pi^*} - R_s^{\pi} = \mathbf{v}_s^{\top} \boldsymbol{\alpha}_{a_s^*} - \mathbf{v}_s^{\top} \boldsymbol{\alpha}_{a_s} = \sum_{j \neq a_s^*} \mathbb{I}(a_s = j) |\mathbf{v}_s^{\top} \boldsymbol{\delta}_{j,a_s^*}|. \quad (\text{EC.4.1})$$

We first focus on the case that  $s = T_2 + 2, \dots, T$ . By applying (EC.3.34) in the online Supplemental Material with  $\theta = 0$  (which corresponds to IV-Greedy), we know that for any  $j \neq a_s^*$ , conditional on the event  $\{a_s = j\}$ ,

$$\begin{aligned} |\mathbf{v}_s^{\top} \boldsymbol{\delta}_{j,a_s^*}| &\leq |\mathbf{v}_s^{\top} (\boldsymbol{\delta}_{j,a_s^*} - \widehat{\boldsymbol{\delta}}_{j,a_s^*,s-1})| = |\mathbf{v}_s^{\top} [(\boldsymbol{\alpha}_j - \widehat{\boldsymbol{\alpha}}_{j,s-1}) - (\boldsymbol{\alpha}_{a_s^*} - \widehat{\boldsymbol{\alpha}}_{a_s^*,s-1})]| \\ &\leq 2\|\mathbf{v}_s\| \|\widehat{\boldsymbol{\alpha}}_{s-1} - \boldsymbol{\alpha}\| \leq 2\sqrt{p\bar{v}} \|\widehat{\boldsymbol{\alpha}}_{s-1} - \boldsymbol{\alpha}\|. \end{aligned} \quad (\text{EC.4.2})$$

Further, recall from (EC.3.1) that

$$\widehat{\boldsymbol{\alpha}}_{s-1} - \boldsymbol{\alpha} = \left( \widehat{\boldsymbol{\Sigma}}_{vz, T_1:s-1} \widehat{\boldsymbol{\Sigma}}_{zz, T_1:s-1}^{-1} \widehat{\boldsymbol{\Sigma}}_{vz, T_1:s-1}^{\top} \right)^{-1} \widehat{\boldsymbol{\Sigma}}_{vz, T_1:s-1} \widehat{\boldsymbol{\Sigma}}_{zz, T_1:s-1}^{-1} \widehat{\mathbf{G}}_{T_1:s-1},$$

and from (EC.3.12) in the online Supplemental Material that

$$\mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} \subseteq \left\{ \left\| \left( \widehat{\boldsymbol{\Sigma}}_{vz, T_1:s-1} \widehat{\boldsymbol{\Sigma}}_{zz, T_1:s-1}^{-1} \widehat{\boldsymbol{\Sigma}}_{vz, T_1:s-1}^{\top} \right)^{-1} \widehat{\boldsymbol{\Sigma}}_{vz, T_1:s-1} \widehat{\boldsymbol{\Sigma}}_{zz, T_1:s-1}^{-1} \right\| \leq \bar{\kappa} \right\}.$$

It then follows from (EC.4.2) that for any  $j \neq a_s^*$ ,

$$\{a_s = j\} \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} \subseteq \left\{ |\mathbf{v}_s^{\top} \boldsymbol{\delta}_{j,a_s^*}| \leq 2\sqrt{p\bar{v}}\bar{\kappa} \|\widehat{\mathbf{G}}_{T_1:s-1}\| \right\}. \quad (\text{EC.4.3})$$

Hence, for any  $j \neq a_s^*$ ,

$$\begin{aligned} &\mathbb{E}[\mathbb{I}(a_s = j) |\mathbf{v}_s^{\top} \boldsymbol{\delta}_{j,a_s^*}|] \\ &= \mathbb{E}[\mathbb{I}(\{a_s = j\} \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1}) |\mathbf{v}_s^{\top} \boldsymbol{\delta}_{j,a_s^*}|] + \mathbb{E}[\mathbb{I}(\{a_s = j\} \cap (\mathcal{B}_{s-1} \cap \mathcal{C}_{s-1})^c) |\mathbf{v}_s^{\top} \boldsymbol{\delta}_{j,a_s^*}|] \\ &\leq \mathbb{E}[\mathbb{I}(\{a_s = j\} \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1}) |\mathbf{v}_s^{\top} \boldsymbol{\delta}_{j,a_s^*}|] + \mathbb{E}[\mathbb{I}((\mathcal{B}_{s-1} \cap \mathcal{C}_{s-1})^c) |\mathbf{v}_s^{\top} \boldsymbol{\delta}_{j,a_s^*}|] \\ &\leq \mathbb{E} \left[ \mathbb{I} \left( |\mathbf{v}_s^{\top} \boldsymbol{\delta}_{j,a_s^*}| \leq 2\sqrt{p\bar{v}}\bar{\kappa} \|\widehat{\mathbf{G}}_{T_1:s-1}\| \right) |\mathbf{v}_s^{\top} \boldsymbol{\delta}_{j,a_s^*}| \right] + 2\sqrt{p\bar{v}} \|\boldsymbol{\delta}_{j,a_s^*}\| \Pr(\mathcal{B}_{s-1}^c \cup \mathcal{C}_{s-1}^c) \end{aligned}$$

$$\leq 2\sqrt{p\bar{v}\bar{\kappa}} \underbrace{\mathbb{E} \left[ \mathbb{I} \left( |\mathbf{v}_s^\top \boldsymbol{\delta}_{j,a_s^*}| \leq 2\sqrt{p\bar{v}\bar{\kappa}} \|\widehat{\mathbf{G}}_{\mathbf{T}_1:s-1}\| \right) \|\widehat{\mathbf{G}}_{\mathbf{T}_1:s-1}\| \right]}_{U_{1,j}} + 2\sqrt{p\bar{v}\bar{\delta}} \underbrace{\Pr(\mathcal{B}_{s-1}^c \cup \mathcal{C}_{s-1}^c)}_{U_2},$$

where the second inequality follows from (EC.4.3), and that  $|\mathbf{v}_s^\top \boldsymbol{\delta}_{j,a_s^*}| \leq \|\mathbf{v}_s\| \|\boldsymbol{\delta}_{j,a_s^*}\| \leq 2\sqrt{p\bar{v}\bar{\delta}}$ . By applying (EC.4.1), we have

$$\mathbb{E}[R_s^{\pi^*} - R_s^\pi] = \sum_{j \neq a_s^*} \mathbb{E}[\mathbb{I}(a_s = j) |\mathbf{v}_s^\top \boldsymbol{\delta}_{j,a_s^*}|] \leq \sum_{j \neq a_s^*} (U_{1,j} + U_2). \quad (\text{EC.4.4})$$

For  $U_2$ , we may apply Lemmas EC.10 and EC.13 to conclude that

$$\begin{aligned} U_2 &\leq \Pr(\mathcal{B}_{s-1}^c) + \Pr(\mathcal{C}_{s-1}^c) \leq \Pr(\mathcal{B}_{\mathbf{T}_2:s-1}^c) + \Pr(\mathcal{C}_{s-1}^c) \\ &\leq [6Mp + (10M + 10)q + 3]T^{-1} + 2qT^{-3} \leq [6Mp + (10M + 12)q + 3]T^{-1}. \end{aligned} \quad (\text{EC.4.5})$$

For  $U_{1,j}$ , we have

$$\begin{aligned} U_{1,j} &= \mathbb{E} \left[ \mathbb{E} \left[ \mathbb{I} \left( |\mathbf{v}_s^\top \boldsymbol{\delta}_{j,a_s^*}| \leq 2\sqrt{p\bar{v}\bar{\kappa}} \|\widehat{\mathbf{G}}_{\mathbf{T}_1:s-1}\| \right) \|\widehat{\mathbf{G}}_{\mathbf{T}_1:s-1}\| \mid \|\widehat{\mathbf{G}}_{\mathbf{T}_1:s-1}\| \right] \right] \\ &\leq \mathbb{E} \left[ \Pr \left( |\mathbf{v}_s^\top \boldsymbol{\delta}_{j,a_s^*}| \leq 2\sqrt{p\bar{v}\bar{\kappa}} \|\widehat{\mathbf{G}}_{\mathbf{T}_1:s-1}\| \mid \|\widehat{\mathbf{G}}_{\mathbf{T}_1:s-1}\| \right) \|\widehat{\mathbf{G}}_{\mathbf{T}_1:s-1}\| \right] \\ &\leq \mathbb{E} \left[ 2L\sqrt{p\bar{v}\bar{\kappa}} \|\widehat{\mathbf{G}}_{\mathbf{T}_1:s-1}\|^2 \right], \end{aligned} \quad (\text{EC.4.6})$$

where the second inequality follows from the margin condition of Assumption 1,

Let  $\mathbf{p}(w)$  denote the probability density function of  $\|\widehat{\mathbf{G}}_{\mathbf{T}_1:s-1}\|$ . Then,

$$\begin{aligned} \mathbb{E} \left[ \|\widehat{\mathbf{G}}_{\mathbf{T}_1:s-1}\|^2 \right] &= \int_0^\infty w^2 \mathbf{p}(w) dw = \int_0^\infty \mathbf{p}(w) \int_0^w 2\tau d\tau dw \\ &= \int_0^\infty 2\tau \int_\tau^\infty \mathbf{p}(w) dw d\tau = \int_0^\infty 2\tau \Pr \left( \|\widehat{\mathbf{G}}_{\mathbf{T}_1:s-1}\| \geq \tau \right) d\tau, \end{aligned} \quad (\text{EC.4.7})$$

where the third equality follows from Fubini's theorem. We then apply (EC.3.13):

$$\int_0^\infty 2\tau \Pr \left( \|\widehat{\mathbf{G}}_{\mathbf{T}_1:s-1}\| \geq \tau \right) d\tau \leq \int_0^\infty 4q\tau \exp \left( -\frac{\tau^2(s - \mathbf{T}_2 - 1)}{2q^2\bar{z}^2\zeta^2} \right) d\tau = \frac{4q^3\bar{z}^2\zeta^2}{s - \mathbf{T}_2 - 1}. \quad (\text{EC.4.8})$$

Combining (EC.4.6), (EC.4.7), and (EC.4.8), we have that for any  $j \neq a_s^*$ ,

$$U_{1,j} \leq \frac{8L\sqrt{p}q^3\bar{v}\bar{z}^2\zeta^2\bar{\kappa}}{s - \mathbf{T}_2 - 1}. \quad (\text{EC.4.9})$$

Plugging in (EC.4.5) and (EC.4.9) into (EC.4.4), we have

$$\mathbb{E}[R_s^{\pi^*} - R_s^\pi] \leq (M - 1) \left[ \frac{16Lpq^3\bar{v}\bar{z}^2\zeta^2\bar{\kappa}^2}{s - \mathbf{T}_2 - 1} + 2\sqrt{p\bar{v}\bar{\delta}}[6Mp + (10M + 12)q + 3]T^{-1} \right]. \quad (\text{EC.4.10})$$

Note that  $\sum_{s=\mathbf{T}_2+2}^T \frac{1}{s-\mathbf{T}_2-1} \leq \int_1^T \frac{du}{u} = \log(T)$ . Then, summing (EC.4.10) over  $s = \mathbf{T}_2 + 2, \dots, T$  yields

$$\sum_{s=\mathbf{T}_2+2}^T \mathbb{E}[R_s^{\pi^*} - R_s^\pi]$$

$$\leq (M-1)[16Lpq^3\bar{v}^2\bar{z}^2\zeta^2\bar{\kappa}^2\log(T) + 2\sqrt{p\bar{v}\bar{\delta}}[6Mp + (10M+12)q + 3]]. \quad (\text{EC.4.11})$$

Moreover, for the first  $T_2 + 1$  periods, we can simply take advantage of the fact that each reward is bounded and  $T_2$  is of order  $\mathcal{O}(\log(T))$ . Specifically, by (EC.4.1),

$$\begin{aligned} \sum_{s=1}^{T_2+1} \mathbb{E}[R_s^{\pi^*} - R_s^\pi] &= \sum_{s=1}^{T_2+1} \mathbb{E}[\mathbb{I}(a_s \neq a_s^*) |\mathbf{v}_s^\top \boldsymbol{\delta}_{a_s, a_s^*}|] \leq \sum_{s=1}^{T_2+1} \mathbb{E}[\|\mathbf{v}_s\| \|\boldsymbol{\delta}_{a_s, a_s^*}\|] \\ &\leq \sqrt{p\bar{v}\bar{\delta}}(T_2 + 1) = \sqrt{p\bar{v}\bar{\delta}}[(C_{T_1} + C_{T_2})\log(T) + 1]. \end{aligned} \quad (\text{EC.4.12})$$

Combining (EC.4.11) and (EC.4.12) completes the proof.  $\square$

## EC.5. Regret Analysis for IV-UCB

We state a more specific version of Theorem 2 as follows.

**THEOREM EC.2.** *Let  $T_1 = C_{T_1} \log(T)$  and  $T_2 = (C_{T_1} + C_{T_2}) \log(T)$  for some constants  $C_{T_1}$  and  $C_{T_2}$  that satisfy Assumption EC.1. Let  $\bar{\delta} := \max_{1 \leq i \neq j \leq M} \|\boldsymbol{\delta}_{i,j}\|$ , where  $\boldsymbol{\delta}_{i,j} := \boldsymbol{\alpha}_i - \boldsymbol{\alpha}_j$ . Then, under Assumption 1, the regret of IV-UCB satisfies that, for all  $T > T_2$ ,*

$$\begin{aligned} \text{Regret}(T) &\leq 4Lp\bar{v}^2(M-1) \left( C_{\mathcal{A}} + \sqrt{\theta\hat{\kappa}(\sigma^2 + \Delta)} \right)^2 (\log(T))^2 + \sqrt{p\bar{v}\bar{\delta}}(C_{T_1} + C_{T_2})\log(T) \\ &\quad + \sqrt{p\bar{v}\bar{\delta}}[(M-1)[24Mp + (40M+48)q + 16] + 1]. \end{aligned}$$

*Proof of Theorem EC.2.* We first focus on the case that  $s = T_2 + 2, \dots, T$ . By applying (EC.3.34) with  $\theta > 0$  (which corresponds to the IV-UCB algorithm), we know that for any  $j \neq a_s^*$ , conditional on the event  $\{a_s = j\}$ ,

$$\begin{aligned} \|\mathbf{v}_s^\top \boldsymbol{\delta}_{j, a_s^*}\| &\leq \underbrace{|\mathbf{v}_s^\top (\boldsymbol{\delta}_{j, a_s^*} - \hat{\boldsymbol{\delta}}_{j, a_s^*, s-1})|}_{W_{1,j}} \\ &\quad + \underbrace{\hat{\sigma}_{s-1} \sqrt{2\theta \log(s - T_1)} \left( \sqrt{\mathbf{v}_s^\top \hat{\boldsymbol{\Omega}}_{j, s-1} \mathbf{v}_s} - \sqrt{\mathbf{v}_s^\top \hat{\boldsymbol{\Omega}}_{a_s^*, s-1} \mathbf{v}_s} \right)}_{W_{2,j}}. \end{aligned} \quad (\text{EC.5.1})$$

Further, recall from (EC.3.35) that

$$\{a_s = j\} \cap \mathcal{A}_{s-1} \subseteq \left\{ |W_{1,j}| \leq 2\sqrt{p\bar{v}}C_{\mathcal{A}} \sqrt{\frac{\log(T)}{s-1-T_1}} \right\}; \quad (\text{EC.5.2})$$

recall also from (EC.3.41) that

$$\{a_s = j\} \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} \cap \mathcal{D}_{s-1} \subseteq \left\{ |W_{2,j}| \leq 2\sqrt{\theta p\bar{v}^2\hat{\kappa}(\sigma^2 + \Delta)} \sqrt{\frac{\log(T)}{s-1-T_1}} \right\}. \quad (\text{EC.5.3})$$

It follows from (EC.5.1), (EC.5.2), and (EC.5.3) that for any  $j \neq a_s^*$ ,

$$\begin{aligned} & \{a_s = j\} \cap \mathcal{A}_{s-1} \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} \cap \mathcal{D}_{s-1} \\ & \subseteq \left\{ |\mathbf{v}_s^\top \boldsymbol{\delta}_{j, a_s^*}| \leq \ddot{\kappa} \sqrt{\frac{\log(T)}{s-1-\mathbf{T}_1}} \right\} := \mathcal{H}_{j,s}, \end{aligned} \quad (\text{EC.5.4})$$

where  $\ddot{\kappa} := 2\sqrt{p\bar{v}} \left( C_A + \sqrt{\theta \ddot{\kappa} (\sigma^2 + \Delta)} \right)$ . By the margin condition in Assumption 1,

$$\Pr(\mathcal{H}_{j,s}) \leq L\ddot{\kappa} \sqrt{\frac{\log(T)}{s-1-\mathbf{T}_1}}. \quad (\text{EC.5.5})$$

Hence,

$$\mathbb{E}[\mathbb{I}(a_s = j) | \mathbf{v}_s^\top \boldsymbol{\delta}_{j, a_s^*}] \quad (\text{EC.5.6})$$

$$\begin{aligned} &= \mathbb{E}[\mathbb{I}(\{a_s = j\} \cap \mathcal{A}_{s-1} \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} \cap \mathcal{D}_{s-1}) | \mathbf{v}_s^\top \boldsymbol{\delta}_{j, a_s^*}] \\ & \quad + \mathbb{E}[\mathbb{I}(\{a_s = j\} \cap (\mathcal{A}_{s-1} \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} \cap \mathcal{D}_{s-1})^c) | \mathbf{v}_s^\top \boldsymbol{\delta}_{j, a_s^*}] \\ & \leq \mathbb{E}[\mathbb{I}(\mathcal{H}_{j,s}) | \mathbf{v}_s^\top \boldsymbol{\delta}_{j, a_s^*}] + 2\sqrt{p\bar{v}}\bar{\delta} \Pr((\mathcal{A}_{s-1} \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} \cap \mathcal{D}_{s-1})^c), \end{aligned} \quad (\text{EC.5.7})$$

where the inequality follows from (EC.5.4).

For the first term on the right-hand-side of (EC.5.7), we have

$$\mathbb{E}[\mathbb{I}(\mathcal{H}_{j,s}) | \mathbf{v}_s^\top \boldsymbol{\delta}_{j, a_s^*}] \leq \ddot{\kappa} \sqrt{\frac{\log(T)}{s-1-\mathbf{T}_1}} \cdot \Pr(\mathcal{H}_{j,s}) \leq \frac{L\ddot{\kappa}^2 \log(T)}{s-1-\mathbf{T}_1}, \quad (\text{EC.5.8})$$

where the two inequalities follow from (EC.5.4) and (EC.5.5), respectively.

For the second term on the right-hand-side of (EC.5.7), we let  $\mathcal{J}_{s-1} := \mathcal{A}_{s-1} \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1}$  and note that

$$\begin{aligned} (\mathcal{J}_{s-1} \cap \mathcal{D}_{s-1})^c &= \mathcal{J}_{s-1}^c \cup \mathcal{D}_{s-1}^c \\ &= \mathcal{J}_{s-1}^c \cup [(\mathcal{D}_{s-1}^c \cap \mathcal{J}_{s-1}^c) \cup (\mathcal{D}_{s-1}^c \cap \mathcal{J}_{s-1})] \\ &= \mathcal{J}_{s-1}^c \cup (\mathcal{D}_{s-1}^c \cap \mathcal{J}_{s-1}) \\ &\subseteq \mathcal{A}_{s-1}^c \cup \mathcal{B}_{s-1}^c \cup \mathcal{C}_{s-1}^c \cup (\mathcal{D}_{s-1}^c \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1}). \end{aligned}$$

Thus,

$$\begin{aligned} & \Pr((\mathcal{A}_{s-1} \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1} \cap \mathcal{D}_{s-1})^c) \\ & \leq \Pr(\mathcal{A}_{s-1}^c) + \Pr(\mathcal{B}_{s-1}^c) + \Pr(\mathcal{C}_{s-1}^c) + \Pr(\mathcal{D}_{s-1}^c \cap \mathcal{B}_{s-1} \cap \mathcal{C}_{s-1}) \\ & \leq [6Mp + (10M + 10)q + 3]T^{-1} + [6Mp + (10M + 10)q + 3]T^{-1} \end{aligned}$$

$$\begin{aligned}
& + 2qT^{-3} + (2q + 2)T^{-3} \\
& \leq [12Mp + (20M + 24)q + 8]T^{-1}, \tag{EC.5.9}
\end{aligned}$$

where the second inequality follows from Lemmas EC.10, EC.11, and EC.13.

Plugging (EC.5.8) and (EC.5.9) into (EC.5.7) and applying (EC.4.1), we have

$$\begin{aligned}
\mathbb{E}[R_s^{\pi^*} - R_s^\pi] &= \sum_{j \neq a_s^*} \mathbb{E}[\mathbb{I}(a_s = j) |\mathbf{v}_s^\top \boldsymbol{\delta}_{j, a_s^*}|] \\
&\leq (M - 1) \left( \frac{L\ddot{\kappa}^2 \log(T)}{s - 1 - \mathbf{T}_1} + 2\sqrt{p\bar{v}\bar{\delta}}[12Mp + (20M + 24)q + 8]T^{-1} \right). \tag{EC.5.10}
\end{aligned}$$

Note that  $\sum_{s=\mathbf{T}_2+2}^T \frac{1}{s-\mathbf{T}_2-1} \leq \int_1^T \frac{du}{u} = \log(T)$ . Then, summing (EC.5.10) over  $s = \mathbf{T}_2 + 2, \dots, T$ , we have

$$\sum_{s=\mathbf{T}_2+2}^T \mathbb{E}[R_s^{\pi^*} - R_s^\pi] \leq (M - 1) \left[ L\ddot{\kappa}^2 (\log(T))^2 + 2\sqrt{p\bar{v}\bar{\delta}}[12Mp + (20M + 24)q + 8] \right].$$

The regret for the first  $\mathbf{T}_2 + 1$  periods is calculated in (EC.4.12). Hence, the proof is completed by plugging the definition of  $\ddot{\kappa}$  into the above inequality.  $\square$

## EC.6. Proofs for Statistical Inference

We state a more specific version of Theorem 3 as follows.

**THEOREM EC.3.** *Let  $\mathbf{T}_1 = C_{\mathbf{T}_1} \log(T)$  and  $\mathbf{T}_2 = (C_{\mathbf{T}_1} + C_{\mathbf{T}_2}) \log(T)$  for some constants  $C_{\mathbf{T}_1}$  and  $C_{\mathbf{T}_2}$  that satisfy Assumption EC.1. Then, under Assumption 1, the estimator  $\hat{\boldsymbol{\alpha}}_T$  of Algorithm 1 satisfies*

$$\sqrt{T - \mathbf{T}_1} (\hat{\boldsymbol{\alpha}}_T - \boldsymbol{\alpha}) \rightsquigarrow \mathcal{N}(\mathbf{0}, \sigma^2 (\boldsymbol{\Sigma}_{vz}^* \boldsymbol{\Sigma}_{zz}^{-1} (\boldsymbol{\Sigma}_{vz}^*)^\top)^{-1}), \tag{EC.6.1}$$

as  $T \rightarrow \infty$ , where  $\boldsymbol{\Sigma}_{vz}^* = \mathbb{E}[\mathbf{v}_t^* \mathbf{z}_t^\top]$  and  $\boldsymbol{\Sigma}_{zz} = \mathbb{E}[\mathbf{z}_t \mathbf{z}_t^\top]$ .

*Proof of Theorem EC.3.* By (EC.3.1), for all  $T > \mathbf{T}_1$ ,

$$\hat{\boldsymbol{\alpha}}_T - \boldsymbol{\alpha} = \left( \hat{\boldsymbol{\Sigma}}_{vz, \mathbf{T}_1:T} \hat{\boldsymbol{\Sigma}}_{zz, \mathbf{T}_1:T}^{-1} \hat{\boldsymbol{\Sigma}}_{vz, \mathbf{T}_1:T}^\top \right)^{-1} \hat{\boldsymbol{\Sigma}}_{vz, \mathbf{T}_1:T} \hat{\boldsymbol{\Sigma}}_{zz, \mathbf{T}_1:T}^{-1} \hat{\mathbf{G}}_{\mathbf{T}_1:T}.$$

It follows from the definition of  $\hat{\mathbf{G}}_{\mathbf{T}_1:T}$  and the central limit theorem that as  $T \rightarrow \infty$ ,

$$\sqrt{T - \mathbf{T}_1} \hat{\mathbf{G}}_{\mathbf{T}_1:T} = \frac{1}{\sqrt{T - \mathbf{T}_1}} \sum_{s=\mathbf{T}_1+1}^T \mathbf{z}_s \epsilon_s \rightsquigarrow \mathcal{N}(\mathbf{0}, \sigma^2 \boldsymbol{\Sigma}_{zz}).$$

By Lemmas EC.10 and EC.13, for arbitrarily small  $\tau > 0$ , for  $T$  large enough (dependent on  $\tau$ ), we have that:

$$\Pr\left(\|\hat{\boldsymbol{\Sigma}}_{vz, \mathbf{T}_1:T} - \boldsymbol{\Sigma}_{vz}^*\| \leq \tau\right) \geq 1 - [6Mp + (10M + 10)q + 3]T^{-1},$$

$$\Pr\left(\|\widehat{\Sigma}_{zz, \tau_1:T} - \Sigma_{zz}\| \leq \tau\right) \geq 1 - 2qT^{-3}.$$

Therefore, as  $T \rightarrow \infty$ ,

$$\widehat{\Sigma}_{vz, \tau_1:T} - \Sigma_{vz}^* \rightarrow_p \mathbf{0} \quad \text{and} \quad \widehat{\Sigma}_{zz, \tau_1:T} - \Sigma_{zz} \rightarrow_p \mathbf{0}.$$

It then follows from the continuous mapping theorem that, as  $T \rightarrow \infty$ ,

$$\begin{aligned} \sqrt{T - \tau_1}(\widehat{\alpha}_T - \alpha) &\rightsquigarrow (\Sigma_{vz}^* \Sigma_{zz}^{-1} (\Sigma_{vz}^*)^\top)^{-1} \Sigma_{vz}^* \Sigma_{zz}^{-1} \cdot \mathcal{N}(\mathbf{0}, \sigma^2 \Sigma_{zz}) \\ &\stackrel{d}{=} \mathcal{N}(\mathbf{0}, \sigma^2 (\Sigma_{vz}^* \Sigma_{zz}^{-1} (\Sigma_{vz}^*)^\top)^{-1}). \quad \square \end{aligned}$$

*Proof of Corollary 1.* By Lemmas EC.10 and EC.13, for arbitrarily small  $\tau > 0$ , for  $T$  large enough (dependent on  $\tau$ ), we have that:

$$\begin{aligned} \Pr\left(\|\widehat{\Sigma}_{vz, \tau_1:T} - \Sigma_{vz}^*\| \leq \tau\right) &\geq 1 - [6Mp + (10M + 10)q + 3]T^{-1}, \\ \Pr\left(\|\widehat{\Sigma}_{zz, \tau_1:T} - \Sigma_{zz}\| \leq \tau\right) &\geq 1 - 2qT^{-3}. \end{aligned}$$

Therefore, as  $T \rightarrow \infty$ ,

$$\widehat{\Sigma}_{vz, \tau_1:T} - \Sigma_{vz}^* \rightarrow_p \mathbf{0} \quad \text{and} \quad \widehat{\Sigma}_{zz, \tau_1:T} - \Sigma_{zz} \rightarrow_p \mathbf{0}.$$

It implies that

$$\widehat{\Omega}_T = \left(\widehat{\Sigma}_{vz, \tau_1:T} \widehat{\Sigma}_{vz, \tau_1:T}^{-1} (\widehat{\Sigma}_{vz, \tau_1:T})^\top\right)^{-1} \rightarrow_p \Omega^*.$$

Similarly, by Lemma EC.11, for arbitrarily small  $\eta$  and large enough  $T$  (dependent on  $\eta$ ), we have that

$$\Pr(\mathcal{B}_T \cap \mathcal{C}_T \cap \mathcal{D}_T^c) \leq (2q + 2)T^{-3}.$$

By Lemma EC.10 and Lemma EC.13, we know that for  $T$  large enough,

$$\begin{aligned} \Pr(\mathcal{B}_T \cap \mathcal{C}_T) &\geq 1 - \Pr(\mathcal{B}_T^c) - \Pr(\mathcal{C}_T^c) \\ &\geq 1 - [6Mp + (10M + 10)q + 3]T^{-1} - 2qT^{-3}. \end{aligned}$$

Therefore,

$$\begin{aligned} \Pr(\mathcal{D}_T^c) &= \Pr(\mathcal{B}_T \cap \mathcal{C}_T \cap \mathcal{D}_T^c) + \Pr((\mathcal{B}_T \cap \mathcal{C}_T)^c \cap \mathcal{D}_T^c) \\ &\leq (2q + 2)T^{-3} + [6Mp + (10M + 10)q + 3]T^{-1} - 2qT^{-3}. \end{aligned}$$

Hence,  $\Pr(\mathcal{D}_T) \rightarrow 1$  as  $T \rightarrow \infty$ . Since  $\eta$  can be arbitrarily small, it implies that  $\widehat{\sigma}_T^2 \rightarrow_p \sigma^2$ . Combining all the above, we have that

$$\widehat{\sigma}_T^2 \widehat{\Omega}_T \rightarrow_p \sigma^2 \Omega^*.$$

By Theorem 3, we have that  $\sqrt{T - \tau_1}(\widehat{\alpha}_T - \alpha) \rightsquigarrow \mathcal{N}(0, \sigma^2 \Omega^*)$  as  $T \rightarrow \infty$ . Since  $\widehat{\sigma}_T^2 \widehat{\Omega}_T$  is a consistent estimator of  $\sigma^2 \Omega^*$ , the statements in Corollary 1 hold.  $\square$

## References

- Golub GH, Van Loan CF (2013) *Matrix Computations* (Johns Hopkins University Press), 4th edition.
- Vershynin R (2018) *High-Dimensional Probability: An Introduction with Applications in Data Science* (Cambridge University Press).
- Wainwright MJ (2019) *High-Dimensional Statistics: A Non-Asymptotic Viewpoint* (Cambridge University Press).
- Wang BY, Xi BY (1997) Some inequalities for singular values of matrix products. *Linear Algebra and its Applications* 264:109–115.