# Stochastic Optimization Algorithms for Instrumental Variable Regression with Streaming Data

Xuxing Chen[♯][*]     Abhishek Roy[♯][†]     Yifan Hu[‡]     Krishnakumar Balasubramanian[§]

May 31, 2024

## Abstract

We develop and analyze algorithms for instrumental variable regression by viewing the problem as a conditional stochastic optimization problem. In the context of least-squares instrumental variable regression, our algorithms neither require matrix inversions nor mini-batches and provides a fully online approach for performing instrumental variable regression with streaming data. When the true model is linear, we derive rates of convergence in expectation, that are of order $\mathcal{O}(\log T/T)$ and $\mathcal{O}(1/T^{1-\iota})$ for any $\iota > 0$, respectively under the availability of two-sample and one-sample oracles, respectively, where $T$ is the number of iterations. Importantly, under the availability of the two-sample oracle, our procedure avoids explicitly modeling and estimating the relationship between confounder and the instrumental variables, demonstrating the benefit of the proposed approach over recent works based on reformulating the problem as minimax optimization problems. Numerical experiments are provided to corroborate the theoretical results.

## 1 Introduction

Instrumental variable analysis is widely used in fields like econometrics, health care, social science, and online advertisement to estimate the causal effect of a random variable, $X$, on an outcome variable, $Y$, when an unobservable confounder influences both. By identifying an instrumental variable correlated with the variable $X$ but unrelated to the confounders, researchers can isolate the exogenous variation in $X$ and estimate a causal relationship between $X$ and $Y$. In the context of regression, Instrumental Variable Regression (IVaR) addresses endogeneity issues when an independent variable is correlated with the error term in the regression model, leveraging an instrument variable $Z$ such that $Y$ is independent of $X|Z$. In this paper, we focus on the following statistical model:

$$Y = g_{\theta^*}(X) + \epsilon_1 \quad \text{with} \quad X = h_{\gamma^*}(Z) + \epsilon_2 \tag{1}$$

where $X \in \mathbb{R}^{d_x}$ and $\epsilon_1$ are correlated and $\epsilon_2$ is a centered unobserved noise (independent of $Z \in \mathbb{R}^{d_z}$), leading to confounding in the model between $X$ and $Y \in \mathbb{R}$. Here $\epsilon_1$ and $\epsilon_2$ are dependent, and $\theta^*$ and $\gamma^*$ are true parameters for the respective function $g$ and $h$. Our goal is to design efficient algorithms that recovers $\theta^*$ from the data.

---

[*]Department of Mathematics, University of California, Davis. Email: xuxchen@ucdavis.edu.

[†]Halıcıoğlu Data Science Institute, University of California, San Diego. Email: a2roy@ucsd.edu.

[‡]College of Management of Technology, EPFL, Department of Computer Science, ETH Zurich, Switzerland. Email: yifan.hu@epfl.ch. YH is supported by NCCR Automation of Swiss National Science Foundation.

[§]Department of Statistics, University of California, Davis. Email: kbala@ucdavis.edu. KB is supported by NSF grant DMS-2053918.

[♯]XC and AR contributed equally to this work.

Traditionally, IVaR algorithms are based on two-stage estimation procedures, where we first regress $Z$ and $X$ to obtain an estimator $\widehat{X}$, and then regress $\widehat{X}$ and $Y$, with the essence that $\widehat{X}$ is independent of $Y$, and thus eliminating the aforementioned endogeneity of the unknown confounder. A vast literature has devoted to understanding the two-stage approaches (Hall and Horowitz, 2005; Darolles et al., 2011; Hartford et al., 2017), with the parametric two-stage least-squares (2SLS) procedure being the most canonical one (Angrist and Imbens, 1995). The main drawback of this approach is that the second-stage regression problem is affected by the estimation error from the regression problem corresponding to first stage. In fact, Angrist and Pischke (2009) call the first stage regression as "forbidden regression", due to the concerns in estimating a nuisance parameter.

Considering the squared loss function, Muandet et al. (2020) formulate the IVaR problem as a conditional stochastic optimization problem (Hu et al., 2020b):

$$\min_{g \in \mathcal{G}} F(g) := \mathbb{E}_Z \mathbb{E}_{Y|Z}[(Y - \mathbb{E}_{X|Z}[g(X)])^2]. \tag{2}$$

However, Muandet et al. (2020) did not solve problem (2) efficiently, and resort to reformulating (2) further as a minimax optimization problem. Indeed, they mention explicitly in their work that "*it remains cumbersome to solve* (2) *directly because of the inner expectation*". Then, they leverage the Fenchel conjugate of the squared loss, leading to a minimax optimization with maximization over a continuous functional space. Following Dai et al. (2017), Muandet et al. (2020) propose to use reproducing kernel Hilbert space (RKHS) to handle the maximization over continuous functional space. See also Lewis and Syrgkanis (2018); Bennett et al. (2019); Dikkala et al. (2020); Liao et al. (2020); Bennett et al. (2023) for similar minimax approaches. The issue with such an approach is that approximating the dual variable via maximization over continuous functional space inevitably introduces approximation error. Hence, although there is no explicit nuisance parameter estimation step like in the two-stage approach, there is an implicit one, which makes the minimax approach less appealing as an alternate to the two-stage procedures.

In this work, contrary to the claim made in Muandet et al. (2020) that problem (2) is cumbersome to solve, we design and analyze efficient streaming algorithms to directly solve the conditional stochastic optimization problem in (2). Direct application of methods from Hu et al. (2020b) for solving (2) is possible, yet their approach utilizes nested sampling, i.e., for each sample of $Z$, Hu et al. (2020b) generate a batch of samples of $X$ from $\mathbb{P}(Z|X)$, to reduce the bias in estimating the composition of non-linear loss function with conditional expectations. Thus their methods are not suitable for the streaming setting that we are interested in. Considering (2), we first parameterize the function class $\mathcal{G} := \{g(\theta; X) \mid \theta \in \mathbb{R}^{d_\theta}\}$. Now, defining $F(g) := F(\theta)$, we observe that the gradient $\nabla F(\theta)$ admits the following form

$$\nabla F(\theta) = \mathbb{E}_Z[(\mathbb{E}_{X|Z}[g(\theta; X)] - \mathbb{E}_{Y|Z}[Y])\nabla_\theta \mathbb{E}_{X|Z}[g(\theta; X)]], \tag{3}$$

which implies that one does not need the nested sampling technique to reduce the bias. However, the presence of product of two conditional expectations $\mathbb{E}_{X|Z}[g(\theta; X)]$ still causes significant challenges in developing stochastic estimators of the above gradient in the streaming setting. In this work, we overcome this challenge and develop two algorithms that are applicable to the streaming data setting avoiding the need for generating batches of samples of $X$ from $\mathbb{P}(Z|X)$.

**Contributions.** We make the following contributions in this work.

- **Two-sample oracles:** Our first algorithm leverages the observation that if we have access to a two-sample oracle that outputs *two* samples $X$ and $X'$ that are independent conditioned on the instrument $Z$, we can immediately construct an unbiased stochastic gradient estimator of the gradient in (3). Based on this crucial observation, we propose the *Two-Sample One-stage Stochastic Gradient IVaR* (TOSG-IVaR) method (Algorithm 1) that avoids explicitly having to estimate or model the relationship between $Z$ and $X$

thereby overcoming the "forbidden regression" problem.. Under standard statistical model assumptions, for the case when $g$ is a linear model, we establish rates of convergence of order $\mathcal{O}(\log T/T)$ for the proposed method, where $T$ is the overall number of iterations; see Theorem 1.

- **One-sample oracles:** In the case when we do not have the aforementioned two-sample oracle, we estimate the stochastic gradient in (3) by using the streaming data to estimate one of the conditional expectations, and the corresponding prediction to estimate the other, resulting in the *One-Sample Two-stage Stochastic Gradient IVaR* (OTSG-IVaR) method (Algorithm 2). Assuming further that the $X$ depends linearly on the instrument $Z$, we establish a rate of convergence of order $\mathcal{O}(1/T^{1-\iota})$, for any $\iota > 0$; see Theorem 2.

## 1.1 Literature Review

**IVaR analysis.** Instrumental variable analysis has a long history, starting from the early works by Wright (1928) and Reiersøl (1945). Several works considered the aforementioned two-stage procedure for IVaR; a summary could be found in the work by Angrist and Pischke (2009). Nonparametric approaches based on wavelets, splines, reproducing kernels and deep neural networks could be found, for example, in the works by Hartford et al. (2017); Singh et al. (2019); Bennett et al. (2019); Muandet et al. (2020); Mastouri et al. (2021); Xu et al. (2021); Zhu et al. (2022); Peixoto et al. (2024). Another popular approach for IVaR is via Generalized Method of Moments (GMM); see, for example, Chen and Pouzo (2012); Bennett et al. (2019); Dikkala et al. (2020) for an overview. Such approaches essentially reformulate the problem as a minimax problem and hence suffer from the aforementioned "forbidden regression" problem.

**Identifiability conditions for IVaR.** Several works in the literature have also focused on establishing the identifiability conditions for IVaR in the parametric and the nonparametric setting. Regardless of the procedure used, they are invariably based on certain source conditions motivated by the inverse problems literature (see, for example, Carrasco et al. (2007); Chen and Reiss (2011); Bennett et al. (2023)) or the related problem of completeness conditions, which posits that the conditional expectation operator is one-to-one (Babii and Florens, 2017; Liao et al., 2020). Semi-parametric identifiability is also considered recently in the work of Cui et al. (2023). Our focus in this work is not focused on the identifiability; for the formulation (2) that we consider, Muandet et al. (2020) provide necessary conditions for identifiability that we adopt.

**Stochastic optimization with nested expectations.** Recently, much attention in the stochastic optimization literature has focused on optimizing a nested composition of $T$ expectation functions. Sample average approximation algorithms in this context are considered in the works of Ermoliev and Norkin (2013) and Hu et al. (2020a). Optimal iterative stochastic optimization algorithms for the case of $T = 2$ were by derived by Ghadimi et al. (2020). For the general $T \geq 1$ case, Wang et al. (2017) provided sub-optimal rates, whereas Balasubramanian et al. (2022) derived optimal rates; see also Zhang and Xiao (2021) and Chen et al. (2021) for related works under stronger assumptions, and Ruszczynski (2021) for similar asymptotic results. While the above works required certain independence assumptions regarding the randomness across the different compositions, Hu et al. (2020b, 2024) studied the case of $T = 2$ where the the randomness are generically dependent. They termed this problem setting as conditional stochastic optimization, which is the framework that the IVaR problem in (2) falls in. Compared to prior works, for e.g., Ghadimi et al. (2020) and Balasubramanian et al. (2022), in order to handle the dependency between the levels, Hu et al. (2020b) require mini-batches in each iteration, making their algorithm not immediately applicable to the purely streaming setting. In this work, we show that despite the problem (2) being a conditional stochastic optimization problem, mini-batches are not required due the additional favorable quadratic structure available in IVaR.

**Streaming IVaR.** Venkatraman et al. (2016); Della Vecchia and Basu (2024) analyzed streaming versions of 2SLS in the online[1] and adversarial settings. Focusing on linear models, Venkatraman et al. (2016) provide preliminary asymptotic analysis assuming access to efficient *no-regret learners*, while Della Vecchia and Basu (2024) provide regret bounds under the strong assumption that the instrument is almost surely bounded. Furthermore, our algorithms have significantly improved per-iteration and memory complexity compared to Della Vecchia and Basu (2024); see Sections A and B for details. Chen et al. (2023) developed stochastic optimization algorithms for the GMM formulation and provide asymptotic analysis. Their algorithm requires access to an offline dataset for initialization and is hence not fully online. The above works (i) do not focus on avoiding the forbidden regression problem and (ii) do not view IVaR via the *conditional stochastic optimization* lens, like we do.

## 2 Two-sample One-stage Stochastic Gradient Method for IVaR

Recall that our goal is to solve the objective function given in (2). By Muandet et al. (2020, Theorem 4), the optimal solution of (2) gives the true underlying causal relationship under the following assumption.

**Assumption 2.1.** *(Identifiability Assumption)*

- *The conditional distribution $\mathbb{P}_{Z|X}$ is continuous in $Z$ for any value of $X$.*
- *The function class $\mathcal{G} := \{g(\theta; X) \mid \theta \in \mathbb{R}^{d_\theta}\}$ is correctly specified, i.e., it includes the true underlying relationship between $X$ and $Y$.*

Notice that both assumptions are standard in the IVaR literature (Newey and Powell, 2003; Chen and Pouzo, 2012; Muandet et al., 2020), and makes the objective in (2) is the meaningful for IVaR. However, Muandet et al. (2020) resort to reformulating the objective function in (2) as a minimax optimization problem as described in Section 1. While their original motivation was to avoid two-state estimation procedure and avoid the "forbidden regression", their minimax reformulation ends up having to solve a complicated approximation of the original objective resulting in having to characterize the approximation error which is non-trivial.

**Algorithm and Analysis.** Our aim in this work is to directly solve the original problem in (2), leveraging the structure provided by the quadratic loss. Given the gradient formulation in (3), a natural way to build unbiased gradient estimator is to generate $X$ and $X'$, two independent samples of $X$ from the conditional distributions $\mathbb{P}_{X|Z}$, for a given realization of $Z$ and generate one sample of $Y$ from the conditional distribution $\mathbb{P}_{Y|X}$. Then, an unbiased gradient estimator is

$$v(\theta) = (g(\theta; X) - Y)\nabla_\theta g(\theta; X'). \tag{4}$$

This could be plugged into the standard stochasic gradient descent algorithm, which give us the Two-sample Stochastic Gradient Method for IVR (TSG-IVaR) method illustrated in Algorithm 1. In particular, the algorithm never requires estimating (or modeling) the relationship between $X$ and $Z$ as needed in the two-stage procedure (Angrist and Pischke, 2009) and the minimax formulation based procedures (Muandet et al., 2020; Lewis and Syrgkanis, 2018; Bennett et al., 2019; Dikkala et al., 2020; Liao et al., 2020; Bennett et al., 2023). Furthermore, this viewpoint not only provides a novel algorithm for performing IV regression, but also provides a novel data collection mechanism for the practical implementation of IVaR. In addition, such a two-sample gradient method is not very restrictive when the instrumental variable $Z$ takes value in a discrete set. In this case, to implement the two-sample oracle, it is enough simply pick two sets of samples

---

[1]Their notion of online is from the literature on *online learning* (Shalev-Shwartz et al., 2012).

**Algorithm 1** Two-sample One-stage Stochastic Gradient-IVaR (`TOSG-IVaR`)

---

**Input:** $\sharp$ of iterations $T$, stepsizes $\{\alpha_t\}_{t=1}^T$, initial iterate $\theta_1$.
1: **for** $t = 1$ to $T$ **do**
2:     Sample $Z_t$, sample independently $X_t$ and $X_t'$ from $\mathbb{P}_{X|Z_t}$, and sample $Y_t$ from $\mathbb{P}_{Y|X_t}$.
3:     Update $\theta_t$
$$\theta_{t+1} = \theta_t - \alpha_{t+1}(g(\theta_t; X_t) - Y_t)\nabla_\theta g(\theta_t; X_t').$$
4: **end for**
**Output:** $\theta_T$.

---

$(X, Y, Z)$ and $(X', Y', Z)$ for which $Z$ has repeated observations (which is possible when $Z$ is a discrete random variable) from a pre-collected dataset. To demonstrate the convergence rate of Algorithm 1, we first consider the case when $g$ is a linear function, i.e., $g(\theta; X) = X^\top \theta$. We make the following assumptions.

**Assumption 2.2.** *Suppose there exists $\mu > 0$ such that $\mathbb{E}_Z\left[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\right] \succeq \mu I$.*

**Assumption 2.3.** *Let $(\vartheta_1, \vartheta_2, \vartheta_3, \vartheta_4) \in \mathbb{R}_+^4$. For any $Z$, $X'$ and $X$ i.i.d. generated from $\mathbb{P}_{Z|X}$, and $Y$ generated from $\mathbb{P}_{Y|X}$, and for constants $C_x, C_y, C_{xx}, C_{yx} > 0$, we have*

$$\mathbb{E}\left[\|X'X^\top - \mathbb{E}_{X|Z}[X]\mathbb{E}_{X|Z}[X]^\top\|^2\right] \leq C_x d_x^{\vartheta_1}, \tag{5}$$

$$\mathbb{E}\left[\|YX - \mathbb{E}_{Y|Z}[Y]\mathbb{E}_{X|Z}[X]\|^2\right] \leq C_y d_x^{\vartheta_2}, \tag{6}$$

$$\mathbb{E}\left[\|\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top - \mathbb{E}_Z\left[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\right]\|^2\right] \leq C_{xx} d_z^{\vartheta_3}, \tag{7}$$

$$\mathbb{E}\left[\|\mathbb{E}_{Y|Z}[Y] \cdot \mathbb{E}_{X|Z}[X] - \mathbb{E}_Z\left[\mathbb{E}_{Y|Z}[Y] \cdot \mathbb{E}_{X|Z}[X]\right]\|^2\right] \leq C_{yx} d_z^{\vartheta_4}, \tag{8}$$

*where $\|\cdot\|$ denotes the Euclidean norm and operator norm for a vector and matrix respectively.*

The above assumptions are mild moment assumptions required on the involved random variables. The following result demonstrates that Assumptions 2.2 and 2.3 are naturally satisfied under even under non-linear modeling assumption on (1). We defer its proof to Section D.3.

**Lemma 1.** *Suppose there exist $\theta_* \in \mathbb{R}^{d_x}$, $\gamma_* \in \mathbb{R}^{d_z \times d_x}$, a non-linear map $\phi : \mathbb{R}^{d_x} \to \mathbb{R}^{d_x}$, and a positive semi-definite matrix $\Sigma \in \mathbb{R}^{d_z \times d_z}$ such that*

$$\mathbb{E}_Z\left[\phi(\gamma_*^\top Z) \cdot \phi(\gamma_*^\top Z)^\top\right] \succeq \mu I, \ \mathbb{E}[\|\phi(\gamma_*^\top Z)\|^2] = \mathcal{O}(d_x),$$
$$Z \sim \mathcal{N}(0, \Sigma), \ X = \phi(\gamma_*^\top Z) + \epsilon_2, \ Y = \theta_*^\top X + \epsilon_1, \ \epsilon_2 \sim \mathcal{N}(0, \sigma_{\epsilon_2}^2 I_{d_x}), \ \epsilon_1 \sim \mathcal{N}(0, \sigma_{\epsilon_1}^2), \tag{9}$$

*where $\epsilon_1, \epsilon_2$ are independent of $Z$ and*

$$\mathbb{E}\left[\epsilon_1^2 \|\epsilon_2\|^2\right] \leq \sigma_{\epsilon_1, \epsilon_2}^2 d_x, \ \mathbb{E}\left[\|\phi(\gamma_*^\top Z) \cdot \phi(\gamma_*^\top Z)^\top - \mathbb{E}[\phi(\gamma_*^\top Z) \cdot \phi(\gamma_*^\top Z)^\top]\|^2\right] \leq C d_z, \tag{10}$$

*then Assumptions 2.2 and 2.3 hold with $\vartheta_1 = \vartheta_2 = 2$ and $\vartheta_3 = \vartheta_4 = 1$. if $\phi$ is an identity map, then the conditions involving $\phi$ become $\gamma_*^\top \Sigma \gamma_* \succeq \mu I$, $\mathrm{tr}(\gamma_*^\top \Sigma \gamma_*) = \mathcal{O}(d_x)$, $\mathbb{E}\left[\|ZZ^\top - \Sigma\|^2\right] \leq C d_z$.*

**Assumption 2.4.** *The tuple $(Z_t, X_t, X_t', Y_t)$ is independent and identically distributed, across $t$.*

The above assumption is standard in the stochastic approximation, statistics and econometrics literature. It could be further relaxed to Markovian-type dependency assumptions, following techniques in the works of Duchi et al. (2012); Sun et al. (2018); Even (2023); Roy et al. (2022); we leave a detailed examination of the Markovian streaming setup as future work. Under the above assumptions, we have the following result demonstrating the last-iterate global convergence of Algorithm 1.

**Theorem 1.** *Suppose Assumptions 2.2, 2.3, and 2.4 hold. In Algorithm 1, defining $\sigma_1^2 := 2C_x d_x^{\vartheta_1} + 2C_{xx} d_z^{\vartheta_3}$ and $\sigma_2^2 := C_y d_x^{\vartheta_2} + C_{yx} d_z^{\vartheta_4}$, set $\alpha_t \equiv \alpha = \frac{\log T}{\mu T} \leq \frac{\mu}{\mu^2 + 3\sigma_1^2}$. Then, we have*

$$\mathbb{E}\big[\|\theta_T - \theta_*\|^2\big] \leq \frac{\mathbb{E}\big[\|\theta_0 - \theta_*\|^2\big]}{T} + \frac{3\|\theta_*\|^2(\sigma_1^2 + \sigma_2^2)\log T}{\mu^2 T}.$$

**Proof techniques.** In the analysis of Theorem 1, the following decomposition (see (18) for the derivation) plays a crucial role:

$$\theta_{t+1} - \theta_* = A_t + \alpha_{t+1} B_t,$$
$$A_t = \theta_t - \alpha_{t+1}\mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\Big]\theta_t + \alpha_{t+1}\mathbb{E}_Z\Big[\mathbb{E}_{Y|Z}[Y] \cdot \mathbb{E}_{X|Z}[X]\Big] - \theta_*,$$
$$B_t = -\Big(X_t' X_t^\top - \mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\Big]\Big)\theta_t + \Big(Y_t X_t' - \mathbb{E}_Z\Big[\mathbb{E}_{Y|Z}[Y] \cdot \mathbb{E}_{X|Z}[X]\Big]\Big),$$

where $A_t$ corresponds to deterministic component, and $B_t$ corresponds to the stochastic component arising due to the use of stochastic gradients. Standard assumptions on the variance of the stochastic gradient made in the stochastic optimization literature include the uniformly bounded variance assumption (Lan, 2020) and the expected smoothness condition (Khaled and Richtárik, 2020). In the IVaR setup, such standard assumptions do not hold as $\theta_t$ potentially can be unbounded and thus the gradient estimator can be unbounded. Hence, we establish our results under natural statistical assumptions arising in the context of the IVaR problem, which form the main novelty in our analysis. Furthermore, compared to Muandet et al. (2020), notice that we use two samples of $X$ from the conditional distribution $\mathbb{P}_{X|Z}$ and achieve an $\widetilde{\mathcal{O}}(1/T)$ last iterate convergence rate to the global optimal solution, which is the true underlying causal relationship under Assumption 2.1. In comparison, Muandet et al. (2020) only provide asymptotic convergence result to the optimal solution of an approximation problem.

**Additional discussion.** It is interesting to explore other losses beyond squared loss (for example to handle classification setting (Centorrino and Florens, 2021)), potentially using the Multilevel Monte Carlo (MLMC) based stochastic gradient estimators. While Hu et al. (2021), develops such algorithms, the main challenge is about how to avoid mini-batches required in their work leveraging the problem structure in instrumental variable analysis. Furthermore, in the case when $g(\theta; X)$ is parametrized by a non-linear models, for instance, a neural network, we provide local convergence guarantees under additional stronger conditions made typically in the stochastic optimization literature.

**Assumption 2.5.** *Let the following assumptions hold:*

- *Function $F(\theta)$ is $\ell$-smooth.*
- *The iterates $\{\theta_t\}_{t=1}^{T+1}$ generated by Algorithm 1 are in a compact set A.*
- *The random objects $X|Z$ and $Y|Z$ have bounded variance for any $Z$, i.e., there exist $\sigma > 0$ such that*

$$\mathbb{E}\left[\|X - \mathbb{E}\left[X \mid Z\right]\|^2 \mid Z\right] \leq \sigma^2, \ \mathbb{E}\left[\|Y - \mathbb{E}\left[Y \mid Z\right]\|^2 \mid Z\right] \leq \sigma^2.$$

**Proposition 1.** *Suppose Assumptions 2.1, 2.4, and 2.5 hold. Choosing $\alpha_t \equiv \alpha = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$, for Algorithm 1 we have*

$$\min_{1 \leq t \leq T} \mathbb{E}\left[\|\nabla F(\theta_t)\|^2\right] = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right).$$

The proof of the proposition is immediate. Note that under Assumption 2.5, we can deduce that the unbiased gradient estimator $v(\theta) = (g(\theta; X) - Y)\nabla_\theta g(\theta; X')$ has a bounded variance since

$$\begin{aligned}
\text{Var}(v(\theta)) =& \text{Var}(g(\theta; X) - Y)\text{Var}(\nabla_\theta g(\theta; X')) \\
&+ \text{Var}(g(\theta; X) - Y)\left(\mathbb{E}\left[\nabla_\theta g(\theta; X')\right]\right)^2 + \text{Var}(\nabla_\theta g(\theta; X'))\left(\mathbb{E}\left[g(\theta; X) - Y\right]\right)^2 \le \sigma_v^2,
\end{aligned}$$

where the variance and expectation are taken conditioning on $Z$ and $\theta$, and $\sigma_v > 0$ is a constant that only depends on $\sigma$, function $g$ and the compact set $A$ in Assumption 2.5. Then one can directly follow the analysis of non-convex stochastic optimization (see, for example, Ghadimi and Lan (2013, Theorem 2.1)) to obtain Proposition 1. Relaxing the Assumption 2.5 (typically made in the stochastic optimization literature) with more natural assumptions on the statistical model and obtaining a result as in Theorem 1 for the non-convex setting is left as future work.

# 3   One-sample Two-stage Stochastic Gradient Method for IVaR

We now examine designing streaming IVaR algorithm with access to the classical one-sample oracle, i.e., we observe a streaming set of samples $(X_t, Y_t, Z_t)$ at each time point $t$. Note that in this case, using the same $X_t$ (instead of $X_t'$) in (4) makes the stochastic gradient estimator biased.

**Intuition.** Consider the case of linear models, i.e., $Y = \theta_*^\top X + \epsilon_1$ with $X = \gamma_*^\top Z + \epsilon_2$, where $\theta_* \in \mathbb{R}^{d_x \times 1}$, and $\gamma_* \in \mathbb{R}^{d_z \times d_x}$, as also considered in Lemma 1. Recall the true gradient in (3) and the stochastic gradient estimator of Algorithm 1 in (4). Since we no longer have $X_t'$, we replace the term $X_t'$ with the predicted mean of $X_t$ given $Z_t$. Suppose that $\gamma_*$ is known. We specifically replace $\nabla_{\theta_t} g(\theta_t; X_t') = X_t'$ by $\mathbb{E}_{|Z_t}[X_t] = \gamma_*^\top Z_t$. In such a case, indeed we have an unbiased gradient estimator:

$$\begin{aligned}
\mathbb{E}_t\left[\gamma_*^\top Z_t(X_t^\top \theta_t - Y_t)\right] &= \mathbb{E}_t\left[\mathbb{E}_{|Z_t}[X_t]\left(\mathbb{E}_{|Z_t}[X_t]^\top \theta_t - \mathbb{E}_{|Z_t}[Y_t]\right)\right] \\
&= \mathbb{E}_t\left[\gamma_*^\top Z_t Z_t^\top \gamma_*(\theta_t - \theta_*)\right] = \gamma_*^\top \Sigma_Z \gamma_*(\theta_t - \theta_*) = \nabla_\theta F(\theta_t),
\end{aligned}$$

where $\mathbb{E}_t[\cdot]$ is the conditional expectation w.r.t the filtration defined on $\{\gamma_1, \theta_1, \gamma_2, \theta_2, \cdots, \gamma_t, \theta_t\}$.

In reality, $\gamma_*$ is unknown beforehand. Hence, we estimate $\gamma_*$ using some online procedure and replace $\nabla_{\theta_t} g(\theta_t; X_t')$ by $\gamma_t^\top Z_t$ instead of $\gamma_*^\top Z_t$. It leads to the following updates:

$$\theta_{t+1} = \theta_t - \alpha_{t+1}\gamma_t^\top Z_t(X_t^\top \theta_t - Y_t), \qquad \gamma_{t+1} = \gamma_t - \beta_{t+1}Z_t(Z_t^\top \gamma_t - X_t^\top). \tag{11}$$

A closer inspection reveals that the updates in (11) can diverge until $\gamma_t$ is close enough to $\gamma_*$. It is easy to see this fact from the following expansion of $\theta_{t+1} - \theta_*$. We have

$$\begin{aligned}
\theta_{t+1} - \theta_* =& \widehat{Q}_t(\theta_t - \theta_*) + \alpha_{t+1}(\gamma_t - \gamma_*)^\top \Sigma_{ZY} + \alpha_{t+1}D_t\theta_* + \alpha_{t+1}\gamma_t^\top \xi_{Z_t}\gamma_*(\theta_t - \theta_*) \\
&+ \alpha_{t+1}\gamma_t^\top \xi_{Z_t}\gamma_*\theta_* + \alpha_{t+1}\gamma_t^\top \xi_{Z_tY_t} - \alpha_{t+1}\gamma_t^\top Z_t\epsilon_{2,t}^\top \theta_t,
\end{aligned} \tag{12}$$

where

$$\xi_{Z_t} = \Sigma_Z - Z_t Z_t^\top, \quad \xi_{Z_tY_t} = \Sigma_{ZY} - Z_t Y_t, \quad \widehat{Q}_t := \left(I - \alpha_{t+1}\gamma_t^\top \Sigma_Z \gamma_*\right).$$

However, the matrix $\gamma_t^\top \Sigma_Z \gamma_*$ may not be positive semi-definite, even if $\Sigma_Z$ is positive definite. Thus the negative eigenvalues associated with $\gamma_t^\top \Sigma_Z \gamma_*$ might cause the $\theta_t$ iterates to first diverge, before eventually converging as $\gamma_t$ gets closer to $\gamma_*$. We illustrate this intuition in a simple experiment in Figure 1. To resolve this issue, we propose Algorithm 2, where we replace $g(\theta_t, X_t) = X_t^\top \theta_t$ with $Z_t^T \gamma_t \theta_t$ in
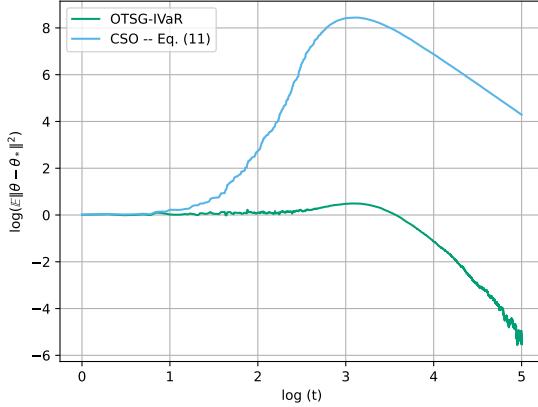
Figure 1: (11) can initially diverge before converging eventually, leading to a worse performance in practical settings compared to Algorithm 2. See Appendix C.2 for the experimental setup.

.

---

**Algorithm 2** One-Sample Two-stage Stochastic Gradient-IVarR (`OTSG-IVaR`)

**Input:** Stepsizes $\{\alpha_t\}_t$, $\{\beta_t\}_t$, initial iterates $\gamma_1, \theta_1$.

1: **for** $t = 1, 2, \cdots$ **do**
2:    Sample $Z_t$, sample $X_t$ from $\mathbb{P}_{X|Z_t}$, Sample $Y_t$ from $\mathbb{P}_{Y|X_t}$.
3:    Update

$$\theta_{t+1} = \theta_t - \alpha_{t+1} \gamma_t^\top Z_t (Z_t^\top \gamma_t \theta_t - Y_t), \tag{13}$$

$$\gamma_{t+1} = \gamma_t - \beta_{t+1} Z_t (Z_t^\top \gamma_t - X_t^\top). \tag{14}$$

4: **end for**

---

(11). With such a modification, in the corresponding decomposition for $\theta_{t+1} - \theta_*$ (see (40)), we have $\widehat{Q}_t = \left(I - \alpha_{t+1} \gamma_t^\top \Sigma_Z \gamma_t\right)$, where the matrix product $\gamma_t^\top \Sigma_Z \gamma_t$ is always positive semi-definite. Hence, with a properly chosen stepsize $\alpha_t$ we could quantify the convergence of $\theta_t$ to $\theta_*$ non-asymptotically. Nevertheless, assuming a warm-start condition on $\theta_0$, we also show the convergence of (11), in Appendix E.3 for completeness.

**Algorithm and Analysis.** Based on the intuition, we present Algorithm 2. One could interpret the algorithm as the SGD analogy of the offline 2SLS algorithm (Angrist and Imbens, 1995). It is also related to the framework of non-linear two-stage stochastic approximation algorithms (Doan and Romberg, 2020; Dalal et al., 2018; Mokkadem and Pelletier, 2006); albeit the updates of $\theta_t$ and $\gamma_t$ are coupled since both updates use $Z_t$. Furthermore, the dependency between the randomness between the two stages in the IVaR problem, makes the analysis significantly different and more challenging from the classical analysis of two-stage algorithms (see below Theorem 2 for additional details). Finally, while Algorithm 2 is designed for linear models, the intuition behind the method is also applicable to non-linear models (i.e., between $Z$ and $X$, and $X$ and $Y$). We focus on linear models in this work in order to derive our theoretical results. A detailed treatment of the nonlinear case (for which the analysis is significantly nontrivial) is left for future work. We make the following additional assumptions for the convergence analysis of Algorithm 2.

8

**Assumption 3.1.** *For some constants $C_z, C_{zy} > 0$, we have the following bounds on the fourth moments:*

$$\mathbb{E}\left[\|\Sigma_Z - ZZ^\top\|^4\right] \leq C_z d_z^{\vartheta_5}, \quad \mathbb{E}\left[\|\Sigma_{ZY} - ZY\|^4\right] \leq C_{zy} d_z^{\vartheta_6}, \quad \vartheta := \max\{\vartheta_5, \vartheta_6\}. \tag{15}$$

**Assumption 3.2.** *There exist constants $0 < \mu_Z \leq \lambda_Z < \infty$ such that $\mu_Z I_{d_z} \preceq \Sigma_Z \preceq \lambda_Z I_{d_z}$.*

The above conditions are rather mild moment conditions, similar to Assumption 2.3, and could be easily verified for the linear model setting we consider.

**Assumption 3.3.** *$\{\gamma_t\}_t$ is within a compact set of diameter $C_\gamma d_z^\varkappa$ for some constants $C_\gamma > 0$, $\varkappa \geq 0$.*

We emphasize that Assumption 3.3 is only for the uncoupled sequence $\gamma_t$, which is an SGD sequence for solving a strongly-convex problem. It holds easily in various cases, for example by projecting the iterates onto any compact sets or a sufficiently large ball containing $\gamma^*$. It is also well-known that, without any projection operations, $\{\gamma_t\}_t$ sequence is almost surely bounded Polyak and Juditsky (1992) under our assumptions. Finally, similar assumptions routinely appear in the analysis of SGD algorithms in various related settings; see, for example, Tseng (1998); Gurbuzbalaban et al. (2019); Haochen and Sra (2019); Nagaraj et al. (2019); Ahn et al. (2020); Rajput et al. (2020).

We now present our result on the convergence of $\{\theta_t\}_t$ below in Theorem 2 (see Appendix E.1 for the proof). In comparison to Theorem 1 (regarding Algorithm 1), we highlight that Theorem 2 provides an any-time guarantee, as the total number of iterations is not required in advance by Algorithm 2.

**Theorem 2.** *Suppose Assumptions 2.2, 2.4 (without $X_t'$), 3.1, 3.3, and 3.2 hold. In Algorithm 2, for any $\iota > 0$, set $\alpha_t = C_\alpha t^{-1+\iota/2}$ and $\beta_t = C_\beta t^{-1+\iota/2}$, where $C_\alpha = \min\{0.5 d_z^{-4\varkappa - \vartheta/2} \lambda_Z^{-1} C_\gamma^{-2}, 0.5(\|\gamma_*\|\lambda_Z)^{-2}\}$, and $C_\beta = \mu^2 d_z^{-1-2\varkappa}/128$. Then, we have*

$$\mathbb{E}\left[\|\theta_t - \theta^*\|^2\right] = O\left(\frac{1}{t^{1-\iota}}\right).$$

**Remark 1.** *In Theorem 2, we present the step-size choices for the fastest rate of convergence. In the proof of Theorem 2 (see Appendix E.1), we show that convergence can be guaranteed for a range of step-sizes given by $\alpha_t = C_\alpha t^{-a}$, $\beta_t = C_\beta t^{-b}$, where $1/2 < a, b < 1$, $b > 2 - 2a$ with corresponding rate being $\mathbb{E}\left[\|\theta_t - \theta^*\|^2\right] = O(\max\{t^{-b(2-(1-\iota/2)^{-1})}, t^{-a}\log(2/(\iota - 1))\})$. In particular, one requires $a, b < 1$ to ensure $(\alpha_t - \alpha_{t+1})/\alpha_t = o(\alpha_t)$, and $(\beta_t - \beta_{t+1})/\beta_t = o(\beta_t)$, as is standard in stochastic approximation literature (see, for example, Chen et al. (2020); Polyak and Juditsky (1992)).*

**Proof Techniques.** The major challenge towards the convergence analysis of $\{\theta_t\}_t$ lies in the interaction term $\gamma_t Z_t Z_t^\top \gamma_t \theta_t$ between $\gamma_t$ and $\theta_t$ in (13). This multiplicative interaction term leads to an involved dependence between the noise in the stochastic gradient updates for the two stages. Such a dependence has not been considered in existing analysis of non-linear two time-scale algorithms (Mokkadem and Pelletier, 2006; Maei et al., 2009; Dalal et al., 2018; Doan and Romberg, 2020; Xu and Liang, 2021; Wang et al., 2021; Doan, 2022). In addition, Doan (2022) considers the case when the noise sequence is not only independent of each other but also independent of iterate locations. Furthermore, they assumes (see their Assumption 3) that the condition in Assumption 2.2 holds for all $\gamma$ whereas Assumption 2.2 only needs to hold for $\gamma_*$, that is much milder. Similarly, many works (for example, Assumption 1 in Wang et al. (2021), Assumption 2 in Xu and Liang (2021) and Theorem 2 in Maei et al. (2009)) assume that the iterates of both stages are bounded in a compact set and consequently, and hence the variance of the stochastic gradients are also uniformly bounded.

In our setting, firstly, the stochastic gradient in (13), evaluated at $(\theta_t, \gamma_t)$ is biased:

$$\mathbb{E}_{t, Z_t}\left[\gamma_t^\top Z_t(Z_t^\top \gamma_t \theta_t - Y_t)\right] = \mathbb{E}_{t, Z_t}\left[\gamma_t^\top Z_t(Z_t^\top \gamma_t \theta_t - Z_t^\top \gamma_* \theta_*)\right] = \mathbb{E}_t\left[\gamma_t^\top \Sigma_Z(\gamma_t \theta_t - \gamma_* \theta_*)\right]$$

$$= \gamma_t^\top \Sigma_Z \gamma_t(\theta_t - \theta_*) + \gamma_t^\top \Sigma_Z(\gamma_t - \gamma_*)\theta_* \neq \gamma_*^\top \Sigma_Z \gamma_*(\theta_t - \theta_*) = \nabla_\theta F(\theta_t).$$
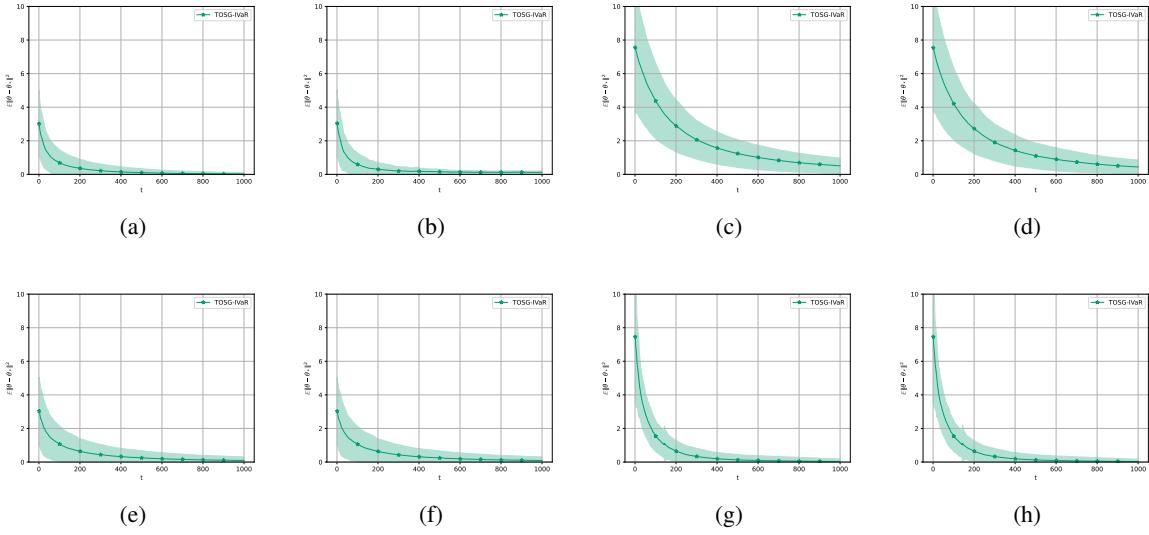
9

Figure 2: $\mathbb{E}[\|\theta_t - \theta_*\|^2]$ of Algorithm 1 under different settings detailed in Section 4.

Furthermore, even under Assumption 3.3, the variance of the stochastic gradient is not (13) uniformly bounded. Overcoming these issues, in addition to the aforementioned dependence between the noise in the stochastic gradient updates for the two stages, forms the major novelty in our analysis. We proceed by noting that if $\gamma_*$, $\Sigma_Z$, and $\Sigma_{ZY}$ were known beforehand, one conduct deterministic gradient updates, i.e., $\widetilde{\theta}_{t+1} = \widetilde{\theta}_t - \alpha_{t+1}\gamma_*^\top \left(\Sigma_Z \gamma_* \widetilde{\theta}_t - \Sigma_{ZY}\right)$, to obtain $\theta_*$. By standard results on gradient descent for strongly convex functions (see, for example, Nesterov (2013)), $\{\widetilde{\theta}_t\}_t$ converges exponentially fast as stated in Lemma 4. Hence, it remains to show that the trajectory of $\theta_t$ converges to the trajectory of $\widetilde{\theta}_t$. That is, defining the sequence $\delta_t := \theta_t - \widetilde{\theta}_t$, our goal is to establish the convergence rate of $\mathbb{E}\left[\|\delta_t\|_2^2\right]$. We first provide an intermediate bound (see Lemma 6) and then progressively sharpen to a tighter bound (see Lemma 7). In doing so, it is also required to show that $\mathbb{E}\left[\|\theta_t\|^4\right]$ is bounded, which we prove in Lemma 5. The proof of Lemma 5 is non-trivial and requires carefully chosen stepsizes satisfying $\sum_{t=1}^\infty (\alpha_t^2 + \alpha_t\sqrt{\beta_t}) < \infty$.

## 4 Numerical Experiments

**Experiments for Algorithm 1 (`TOSG-IVaR`).** We first consider the following problem, in which $(Z, X, Y)$ is generated via

$$Z \sim \mathcal{N}(0, I_{d_z}), \ X = \phi(\gamma_*^\top Z) + c \cdot (h + \epsilon_x), \ Y = \theta_*^\top X + c \cdot (h_1 + \epsilon_y),$$

where $c > 0$ is a scalar to control the variance of the noise vector, and $h_1$ is the first coordinate of $h$. The noise vectors (or scalar) $h, \epsilon_x, \epsilon_y$ are independent of $Z$, and we have $h \sim \mathcal{N}(\mathbf{1}_{d_x}, I_{d_x})$, $\epsilon_x \sim \mathcal{N}(0, I_{d_x})$, $\epsilon_y \sim \mathcal{N}(0, 1)$. In each iteration, one tuple $(X, X', Y)$ is generated and used to update $\theta_t$ according to Algorithm 1. We set $(d_x, d_z) \in \{(4, 8), (8, 16)\}$, $c \in \{0.1, 1.0\}$, and $\phi(s) \in \{s, s^2\}$. We repeat each setting 50 times and report the curves of $\mathbb{E}[\|\theta_t - \theta_*\|^2]$ in Figure 2, where the expectation is computed as the average of $\|\theta_t - \theta_*\|^2$ of all trials, and the shaded region represents the standard deviation. The first row and the second row correspond to $\phi(s) = s$ and $\phi(s) = s^2$ respectively. Here, $c = 0.1$ for odd columns and $c = 1.0$ for even columns. We have $(d_x, d_z) = (4, 8)$ for the first two columns and $(d_x, d_z) = (8, 16)$ for the last two columns. Empirically, we can observe that our Algorithm 1 performs well across all different settings.
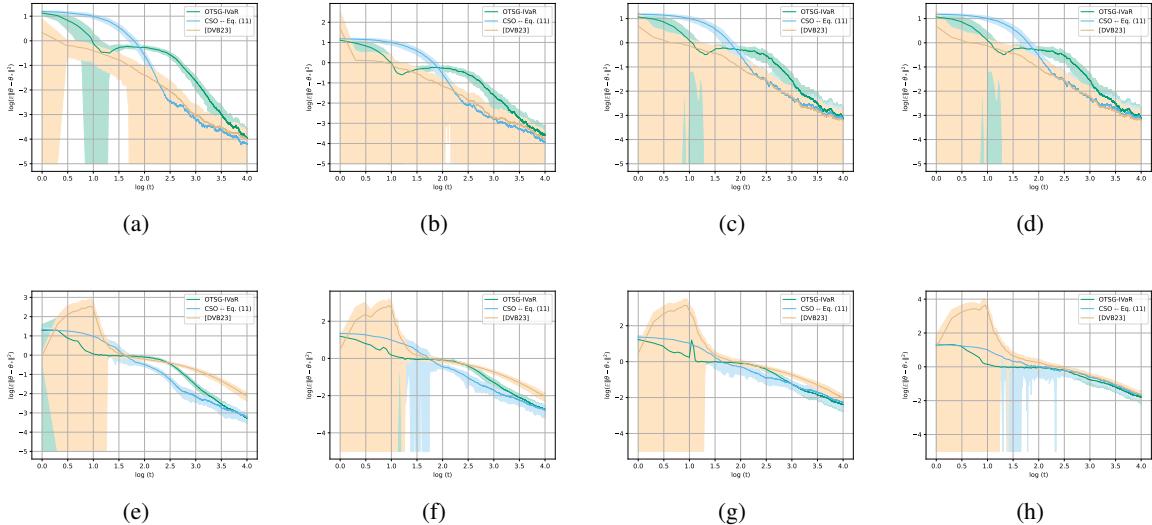
Figure 3: Comparison of $\mathbb{E}[\|\theta_t - \theta_*\|^2]$ (log-log scale) for Algorithm 2, Eq. 11 and Della Vecchia and Basu (2024).

**Experiments for Algorithm 2 (OTSG-IVaR).** Next, we compare our Algorithm 2 as well as its variant and Algorithm 1 in Della Vecchia and Basu (2024). We write "OTSG-IVaR", "CSO – Eq. (11)" and "[DVB23]" to represent Algorithm 2, Algorithm 2 with the updates replaced by (11) and Algorithm 1 in Della Vecchia and Basu (2024) (see Appendix A). We follow simulation settings similar to Della Vecchia and Basu (2024):

$$Y = \theta_*^\top X + \nu, \qquad X = \gamma_*^\top Z + \epsilon, \qquad \epsilon = \sigma_\epsilon \mathcal{N}(0, I_{d_x}), \qquad \nu = \rho \epsilon_1 + \mathcal{N}(0, 0.25), \tag{16}$$

where $\epsilon_1$ is the first coordinate of $\epsilon$, $\theta_* \in \mathbb{R}^{d_x}$ is a unit vector chosen uniformly randomly, and $\gamma_* \in \mathbb{R}^{d_z \times d_x}$ where $\gamma_{ij} = 0$ for $i \neq j$, and $\gamma_{ij} = 1$ for $i = j$, $i = 1, 2, \cdots, d_x$, and $j = 1, 2, \cdots, d_z$. Here $\rho$ controls the level of endogeneity in the model. We compare the performance of Algorithm 2 with (11), and O2SLS (Della Vecchia and Basu, 2024) for $\rho = 1, 4$, and $\sigma_\epsilon = 0.5, 1$. By varying $\sigma_\epsilon$ we control the correlation between $X$ and $Z$. We consider two settings $(d_x, d_z) = (1, 1)$, and $(d_x, d_z) = (8, 16)$. As performance metric, in Figure 3 we plot $\mathbb{E}[\|\theta_t - \theta_*\|^2]$ where the $\mathbb{E}[\cdot]$ is approximated by averaging over 50 trials, and both axes are in $\log$ scale (base 10). We also show, in Figure 4, the convergence of the test Mean Squared Error (MSE) evaluated over 400 test samples to the best possible Test MSE where $\theta_*$ and $\gamma_*$ are known beforehand. For Figures 3 and 4, the first row and second row corresponds to $(d_x, d_z) = (1, 1)$ and $(d_x, d_z) = (8, 16)$ respectively, and $\sigma_\epsilon = 0.5$ in odd columns and $\sigma_\epsilon = 1.0$ in even columns. We have $\rho = 1.0$ for the first two columns and $\rho = 4.0$ for the last two columns. We can observe that O2SLS has much larger variance in different settings, while our algorithms perform consistently well in all settings.

## 5  Conclusion

We presented streaming algorithms for least-squares IVaR based on directly solving the associated conditional stochastic optimization formulation in (2). Our algorithms have several benefits, including avoidance of mini-batches and matrix inverses. We show that the expected rates of convergences for the proposed algorithms are of order $\mathcal{O}(\log T / T)$ and $\mathcal{O}(1/T^{1-\iota})$, for any $\iota > 0$, under the availability of two-sample and one-sample oracles, respectively. As future work, it is interesting to develop streaming inferential methods for IVaR. Leveraging related works for the vanilla SGD (Polyak and Juditsky, 1992; Anastasiou et al., 2019; Shao and Zhang, 2022; Chen et al., 2020; Zhu et al., 2023) to the setting of Algorithms 1 and 2, pro-
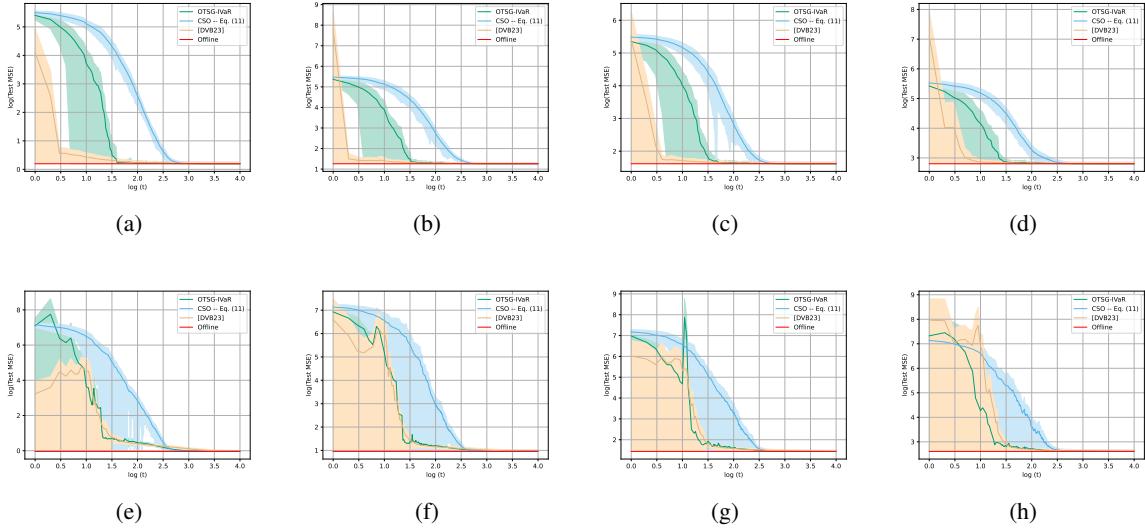
Figure 4: Comparison of Test MSE (log-log scale) for Algorithm 2, Eq. 11 and Della Vecchia and Basu (2024).

vides a concrete direction to establish Central Limit Theorems and develop limiting covariance estimation procedures.

# References

K. Ahn, C. Yun, and S. Sra. SGD with shuffling: optimal rates without component convexity and large epoch requirements. *Advances in Neural Information Processing Systems*, 33:17526–17535, 2020. (Cited on page 9.)

A. Anastasiou, K. Balasubramanian, and M. A. Erdogdu. Normal approximation for stochastic gradient descent via non-asymptotic rates of martingale clt. In *Conference on Learning Theory*, pages 115–137. PMLR, 2019. (Cited on page 11.)

J. D. Angrist and G. W. Imbens. Two-stage least squares estimation of average causal effects in models with variable treatment intensity. *Journal of the American statistical Association*, 90(430):431–442, 1995. (Cited on pages 2 and 8.)

J. D. Angrist and J.-S. Pischke. *Mostly harmless econometrics: An empiricist's companion*. Princeton university press, 2009. (Cited on pages 2, 3, and 4.)

A. Babii and J.-P. Florens. Is completeness necessary? estimation in nonidentified linear models. *arXiv preprint arXiv:1709.03473*, 2017. (Cited on page 3.)

K. Balasubramanian, S. Ghadimi, and A. Nguyen. Stochastic multilevel composition optimization algorithms with level-independent convergence rates. *SIAM Journal on Optimization*, 32(2):519–544, 2022. (Cited on page 3.)

A. Bennett, N. Kallus, and T. Schnabel. Deep generalized method of moments for instrumental variable analysis. *Advances in neural information processing systems*, 32, 2019. (Cited on pages 2, 3, and 4.)

A. Bennett, N. Kallus, X. Mao, W. Newey, V. Syrgkanis, and M. Uehara. Minimax instrumental variable regression and $l_2$ convergence guarantees without identification or closedness. *arXiv preprint arXiv:2302.05404*, 2023. (Cited on pages 2, 3, and 4.)

M. Carrasco, J.-P. Florens, and E. Renault. Linear inverse problems in structural econometrics estimation based on spectral decomposition and regularization. *Handbook of econometrics*, 6:5633–5751, 2007. (Cited on page 3.)

S. Centorrino and J.-P. Florens. Nonparametric instrumental variable estimation of binary response models with continuous endogenous regressors. *Econometrics and Statistics*, 17:35–63, 2021. (Cited on page 6.)

T. Chen, Y. Sun, and W. Yin. Solving stochastic compositional optimization is nearly as easy as solving stochastic optimization. *IEEE Transactions on Signal Processing*, 69:4937–4948, 2021. (Cited on page 3.)

X. Chen and D. Pouzo. Estimation of nonparametric conditional moment models with possibly nonsmooth generalized residuals. *Econometrica*, 80(1):277–321, 2012. (Cited on pages 3 and 4.)

X. Chen and M. Reiss. On rate optimality for ill-posed inverse problems in econometrics. *Econometric Theory*, 27(3):497–521, 2011. (Cited on page 3.)

X. Chen, J. D. Lee, X. T. Tong, and Y. Zhang. Statistical inference for model parameters in stochastic gradient descent. *Annals of Statistics*, 48(1):251–273, 2020. (Cited on pages 9, 11, and 21.)

X. Chen, S. Lee, Y. Liao, M. H. Seo, Y. Shin, and M. Song. SGMM: Stochastic approximation to generalized method of moments. *arXiv preprint arXiv:2308.13564*, 2023. (Cited on page 4.)

Y. Cui, H. Pu, X. Shi, W. Miao, and E. Tchetgen Tchetgen. Semiparametric proximal causal inference. *Journal of the American Statistical Association*, pages 1–12, 2023. (Cited on page 3.)

B. Dai, N. He, Y. Pan, B. Boots, and L. Song. Learning from conditional distributions via dual embeddings. In *Artificial Intelligence and Statistics*, pages 1458–1467. PMLR, 2017. (Cited on page 2.)

G. Dalal, G. Thoppe, B. Szörényi, and S. Mannor. Finite sample analysis of two-timescale stochastic approximation with applications to reinforcement learning. In *Conference On Learning Theory*, pages 1199–1233. PMLR, 2018. (Cited on pages 8 and 9.)

S. Darolles, Y. Fan, J.-P. Florens, and E. Renault. Nonparametric instrumental regression. *Econometrica*, 79(5):1541–1565, 2011. (Cited on page 2.)

R. Della Vecchia and D. Basu. Online instrumental variable regression: Regret analysis and bandit feedback. *arXiv preprint arXiv:2302.09357v1*, 2023. (Cited on pages 17 and 18.)

R. Della Vecchia and D. Basu. Stochastic online instrumental variable regression: Regrets for endogeneity and bandit feedback. *arXiv preprint arXiv:2302.09357v3*, 2024. (Cited on pages 4, 11, 12, 17, and 18.)

N. Dikkala, G. Lewis, L. Mackey, and V. Syrgkanis. Minimax estimation of conditional moment models. *Advances in Neural Information Processing Systems*, 33:12248–12262, 2020. (Cited on pages 2, 3, and 4.)

T. Doan and J. Romberg. Finite-time performance of distributed two-time-scale stochastic approximation. In *Learning for Dynamics and Control*, pages 26–36. PMLR, 2020. (Cited on pages 8 and 9.)

T. T. Doan. Nonlinear two-time-scale stochastic approximation convergence and finite-time performance. *IEEE Transactions on Automatic Control*, 2022. (Cited on page 9.)

J. C. Duchi, A. Agarwal, M. Johansson, and M. I. Jordan. Ergodic mirror descent. *SIAM Journal on Optimization*, 22(4):1549–1578, 2012. (Cited on page 5.)

Y. M. Ermoliev and V. I. Norkin. Sample average approximation method for compound stochastic optimization problems. *SIAM Journal on Optimization*, 23(4):2231–2263, 2013. (Cited on page 3.)

M. Even. Stochastic gradient descent under Markovian sampling schemes. In *International Conference on Machine Learning*, pages 9412–9439. PMLR, 2023. (Cited on page 5.)

S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM journal on optimization*, 23(4):2341–2368, 2013. (Cited on page 7.)

S. Ghadimi, A. Ruszczynski, and M. Wang. A single timescale stochastic approximation method for nested stochastic optimization. *SIAM Journal on Optimization*, 30(1):960–979, 2020. (Cited on page 3.)

M. Gurbuzbalaban, A. Ozdaglar, and P. A. Parrilo. Convergence rate of incremental gradient and incremental newton methods. *SIAM Journal on Optimization*, 29(4):2542–2565, 2019. (Cited on page 9.)

P. Hall and J. L. Horowitz. Nonparametric methods for inference in the presence of instrumental variables. *Annals of statistics*, 33(6):2904–2929, 2005. (Cited on page 2.)

J. Haochen and S. Sra. Random shuffling beats SGD after finite epochs. In *International Conference on Machine Learning*, pages 2624–2633. PMLR, 2019. (Cited on page 9.)

J. Hartford, G. Lewis, K. Leyton-Brown, and M. Taddy. Deep iv: A flexible approach for counterfactual prediction. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1414–1423. JMLR, 2017. (Cited on pages 2 and 3.)

Y. Hu, X. Chen, and N. He. Sample complexity of sample average approximation for conditional stochastic optimization. *SIAM Journal on Optimization*, 30(3):2103–2133, 2020a. (Cited on page 3.)

Y. Hu, S. Zhang, X. Chen, and N. He. Biased stochastic first-order methods for conditional stochastic optimization and applications in meta learning. *Advances in Neural Information Processing Systems*, 33: 2759–2770, 2020b. (Cited on pages 2 and 3.)

Y. Hu, X. Chen, and N. He. On the bias-variance-cost tradeoff of stochastic optimization. *Advances in Neural Information Processing Systems*, 34:22119–22131, 2021. (Cited on page 6.)

Y. Hu, J. Wang, Y. Xie, A. Krause, and D. Kuhn. Contextual stochastic bilevel optimization. *Advances in Neural Information Processing Systems*, 36, 2024. (Cited on page 3.)

A. Khaled and P. Richtárik. Better theory for SGD in the nonconvex world. *arXiv preprint arXiv:2002.03329*, 2020. (Cited on page 6.)

G. Lan. *First-order and stochastic optimization methods for machine learning*, volume 1. Springer, 2020. (Cited on page 6.)

G. Lewis and V. Syrgkanis. Adversarial generalized method of moments. *arXiv preprint arXiv:1803.07164*, 2018. (Cited on pages 2 and 4.)

L. Liao, Y.-L. Chen, Z. Yang, B. Dai, M. Kolar, and Z. Wang. Provably efficient neural estimation of structural equation models: An adversarial approach. *Advances in Neural Information Processing Systems*, 33:8947–8958, 2020. (Cited on pages 2, 3, and 4.)

H. Maei, C. Szepesvari, S. Bhatnagar, D. Precup, D. Silver, and R. S. Sutton. Convergent temporal-difference learning with arbitrary smooth function approximation. *Advances in neural information processing systems*, 22, 2009. (Cited on page 9.)

A. Mastouri, Y. Zhu, L. Gultchin, A. Korba, R. Silva, M. Kusner, A. Gretton, and K. Muandet. Proximal causal learning with kernels: Two-stage estimation and moment restriction. In *International conference on machine learning*, pages 7512–7523. PMLR, 2021. (Cited on page 3.)

A. Mokkadem and M. Pelletier. Convergence rate and averaging of nonlinear two-time-scale stochastic approximation algorithms. *Annals of Applied Probability*, 16(3):1671–1702, 2006. (Cited on pages 8 and 9.)

K. Muandet, A. Mehrjou, S. K. Lee, and A. Raj. Dual instrumental variable regression. *Advances in Neural Information Processing Systems*, 33:2710–2721, 2020. (Cited on pages 2, 3, 4, and 6.)

D. Nagaraj, P. Jain, and P. Netrapalli. SGD without replacement: Sharper rates for general smooth convex functions. In *International Conference on Machine Learning*, pages 4703–4711. PMLR, 2019. (Cited on page 9.)

Y. Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013. (Cited on page 10.)

W. K. Newey and J. L. Powell. Instrumental variable estimation of nonparametric models. *Econometrica*, 71(5):1565–1578, 2003. (Cited on page 4.)

C. H. Papadimitriou. Computational complexity. In *Encyclopedia of computer science*, pages 260–265. 2003. (Cited on page 17.)

C. Peixoto, Y. Saporito, and Y. Fonseca. Nonparametric instrumental variable regression through stochastic approximate gradients. *arXiv preprint arXiv:2402.05639*, 2024. (Cited on page 3.)

B. T. Polyak and A. B. Juditsky. Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization*, 30(4):838–855, 1992. (Cited on pages 9, 11, and 29.)

S. Rajput, A. Gupta, and D. Papailiopoulos. Closing the convergence gap of SGD without replacement. In *International Conference on Machine Learning*, pages 7964–7973. PMLR, 2020. (Cited on page 9.)

O. Reiersøl. *Confluence analysis by means of instrumental sets of variables*. PhD thesis, Almqvist & Wiksell, 1945. (Cited on page 3.)

A. Roy, K. Balasubramanian, and S. Ghadimi. Constrained stochastic nonconvex optimization with state-dependent Markov data. *Advances in Neural Information Processing Systems*, 35:23256–23270, 2022. (Cited on page 5.)

A. Ruszczynski. A stochastic subgradient method for nonsmooth nonconvex multilevel composition optimization. *SIAM Journal on Control and Optimization*, 59(3):2301–2320, 2021. (Cited on page 3.)

S. Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and Trends® in Machine Learning*, 4(2):107–194, 2012. (Cited on page 4.)

Q.-M. Shao and Z.-S. Zhang. Berry–esseen bounds for multivariate nonlinear statistics with applications to m-estimators and stochastic gradient descent algorithms. *Bernoulli*, 28(3):1548–1576, 2022. (Cited on page 11.)

R. Singh, M. Sahani, and A. Gretton. Kernel instrumental variable regression. *Advances in Neural Information Processing Systems*, 32, 2019. (Cited on page 3.)

T. Sun, Y. Sun, and W. Yin. On Markov chain gradient descent. *Advances in neural information processing systems*, 31, 2018. (Cited on page 5.)

P. Tseng. An incremental gradient (-projection) method with momentum term and adaptive stepsize rule. *SIAM Journal on Optimization*, 8(2):506–531, 1998. (Cited on page 9.)

A. Venkatraman, W. Sun, M. Hebert, J. Bagnell, and B. Boots. Online instrumental variable regression with applications to online linear system identification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016. (Cited on page 4.)

M. Wang, E. X. Fang, and B. Liu. Stochastic compositional gradient descent: Algorithms for minimizing compositions of expected-value functions. *Mathematical Programming*, 161(1-2):419–449, 2017. (Cited on page 3.)

Y. Wang, S. Zou, and Y. Zhou. Non-asymptotic analysis for two time-scale tdc with general smooth function approximation. *Advances in Neural Information Processing Systems*, 34:9747–9758, 2021. (Cited on page 9.)

P. G. Wright. *The tariff on animal and vegetable oils*. Number 26. Macmillan, 1928. (Cited on page 3.)

L. Xu, Y. Chen, S. Srinivasan, N. de Freitas, A. Doucet, and A. Gretton. Learning deep features in instrumental variable regression. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=sy4Kg_ZQmS7. (Cited on page 3.)

T. Xu and Y. Liang. Sample complexity bounds for two timescale value-based reinforcement learning algorithms. In *International Conference on Artificial Intelligence and Statistics*, pages 811–819. PMLR, 2021. (Cited on page 9.)

J. Zhang and L. Xiao. Multilevel composite stochastic optimization via nested variance reduction. *SIAM Journal on Optimization*, 31(2):1131–1157, 2021. (Cited on page 3.)

W. Zhu, X. Chen, and W. B. Wu. Online covariance matrix estimation in stochastic gradient descent. *Journal of the American Statistical Association*, 118(541):393–404, 2023. (Cited on page 11.)

Y. Zhu, L. Gultchin, A. Gretton, M. J. Kusner, and R. Silva. Causal inference with treatment measurement error: a nonparametric instrumental variable approach. In *Uncertainty in Artificial Intelligence*, pages 2414–2424. PMLR, 2022. (Cited on page 3.)

# A Online updates of Della Vecchia and Basu (2024)

For the sake of clarity, we present the O2SLS algorithm proposed in (Della Vecchia and Basu, 2024, v3)[1] in the streaming format, without any explicit matrix inversions that we used in our experiments:

$$\theta_{t+1} = (I - U_t \gamma_t^\top Z_t Z_t^\top \gamma_t)\theta_t + U_t \gamma_t^\top Z_t Y_t$$

$$\gamma_{t+1} = (I - V_t Z_t Z_t^\top)\gamma_t + V_t Z_t X_t^\top$$

$$U_{t+1} = U_t - \frac{U_t \gamma_t^\top Z_t Z_t^\top \gamma_t U_t}{1 + Z_t^\top \gamma_t U_t \gamma_t^\top Z_t}$$

$$V_{t+1} = V_t - \frac{V_t Z_t Z_t^\top V_t}{1 + Z_t^\top V_t Z_t} \qquad V_0 = \lambda^{-1} I_{d_z},$$

where $U_t, V_t$ are two additional matrix sequences which tracks the matrix inverse of $\sum_{i=1}^t \gamma_i^\top Z_i Z_i^\top \gamma_i$, and $(\lambda I_{d_z} + \sum_{i=1}^t Z_t Z_t^\top)$ respectively for a user defined parameter $\lambda$. As mentioned in Della Vecchia and Basu (2024), we choose $\lambda = 0.1$. The major difference between O2SLS and Algorithm 2 is that O2SLS takes an online two-stage regression approach to minimize a suitably defined regret whereas we take a conditional stochastic optimization point of view which requires carefully chosen step-sizes. In our Algorithm 2, we do not need to explicitly or implicitly do matrix inverse which can potentially cause stability issues. Furthermore, unlike Della Vecchia and Basu (2024), we neither assume $\sum_{i=1}^t Z_i Z_i^\top$ is invertible for all $t$ nor do we assume that $Z$ is a bounded random variable for our analysis. Finally, the per-iteration computational complexity and memory requirement of Algorithm 2 is significantly better than O2SLS; see Section B.

# B Per-iteration Complexities

For the linear case, i.e., the underlying relationship between $Z$ and $X$ as well as $X$ and $Y$ are linear, Table 1 summarizes the per-iteration memory costs and number of arithmetic operations of the original O2SLS (Della Vecchia and Basu, 2023), the updated O2SLS (Della Vecchia and Basu, 2024) that we provide a matrix form update in Appendix A, TOSG-IVaR (Alg 1), and OTSG-IVaR (Alg 2) at the $t$-th iteration.

Notice that the original version of O2SLS (Della Vecchia and Basu, 2023) has a per-iteration and memory cost dependent on the iteration number $t$ as it needs to use all the samples accumulated till the iteration $t$ to conduct an offline 2SLS at each iteration. The updated O2SLS (Della Vecchia and Basu, 2024) (the algorithm that we compare to) uses samples obtained at iteration $t$ to perform the update. Although the updated O2SLS avoids explicit matrix inversion, it is obvious that its arithmetic operations and memory cost per iteration are larger than our TOSG-IVaR and OTSG-IVaR.

We highlight that the TOSG-IVaR, which uses two samples $X$ and $X'$ from the conditional distribution $\mathbb{P}(X \mid Z)$, requires only $\mathcal{O}(d_x)$ memory and arithmetic operations at each iteration.

For a fair comparison, we assume that two $n \times n$ matrices multiplication admits an $\mathcal{O}(n^3)$ complexity, i.e., using normal textbook matrix multiplication. We also assume computing the inversion of a $n \times n$ matrix admits an $\mathcal{O}(n^3)$ complexity. Interested readers may refer to (Papadimitriou, 2003) for more details about faster algorithms with better complexities for matrix operations.

---

[1]Note that the streaming algorithm was not present in version 1, i.e., (Della Vecchia and Basu, 2023, v1).

Table 1: Memory cost and the number of arithmetic operations at iteration $t$.

| Algorithm | Memory cost | Arithmetic Operations |
|---|---|---|
| O2SLS (Della Vecchia and Basu, 2023, v1) | $t(d_x + d_z) + d_z d_x + d_x$ | $\mathcal{O}(d_x^3 + t d_x^2 + t d_x d_z)$ |
| O2SLS (Della Vecchia and Basu, 2024, v3) (Sec. A) | $d_x^2 + d_z^2 + d_z d_x + d_x$ | $\mathcal{O}(d_x^2 + d_z^2 + d_z d_x)$ |
| TOSG-IVaR (our Alg 1) | $d_x$ | $\mathcal{O}(d_x)$ |
| OTSG-IVaR (our Alg 2) | $d_x d_z + d_x$ | $\mathcal{O}(\min(d_x^2, d_z^2) + d_z d_x)$ |

# C  Experimental Details

## C.1  Compute Resources

All experiments in Section 4 were conducted on a computer with an 11th Intel(R) Core(TM) i7-11370H CPU. The time and space required to run our experiments are negligible and we anticipate they can be conducted in almost all computers.

## C.2  Experimental Details for Figure 1

In Figure 1, we show an example where the updates (11) may diverge first before converging eventually and finite time performance can be much worse compared to Algorithm 2. For this experiment, we choose the model presented in (16) with $d_x = d_z = 1$, $\theta_* = 1$, $\gamma_* = -1$, $\rho = 4$, and $\sigma_\epsilon = 1$. When initialized at $\gamma_0 = 10$, and $\theta_0 = 0$, the updates in (11) keeps diverging rapidly at first whereas Algorithm 2 is much more stable. So, by the end of $100,000$ iterations, while Algorithm 2 achieves an error of $\approx 10^{-5}$, (11) achieves $\approx 10^4$ that is worse than it was at initialization because (11) has not recovered from the initial divergence phase yet. However, once (11) starts converging, the convergence rate of (11) is similar to Algorithm 2 as one can see from Figure 1 (also see our discussion on the convergence of (11) in Section E.3).

# D  Proofs for Section 2

## D.1  Proof of Theorem 1

*Proof.* We aim to find the optimal $\theta_*$. According to (2), we know

$$\mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\Big]\theta_* = \mathbb{E}_Z\Big[\mathbb{E}_{Y|Z}[Y] \cdot \mathbb{E}_{X|Z}[X]\Big] \tag{17}$$

The updates in Algorithm 1 can be written as

$$\theta_{t+1} = \theta_t - \alpha_{t+1}(X_t^\top \theta_t - Y_t)X_t^\top.$$

Hence we have

$$\begin{aligned}
&\theta_{t+1} - \theta_* \\
=&\theta_t - \alpha_{t+1}\mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\Big]\theta_t + \alpha_{t+1}\mathbb{E}_Z\Big[\mathbb{E}_{Y|Z}[Y] \cdot \mathbb{E}_{X|Z}[X]\Big] - \theta_* \\
&- \alpha_{t+1}\Big(X_t' X_t^\top - \mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\Big]\Big)\theta_t + \alpha_{t+1}\Big(Y_t X_t' - \mathbb{E}_Z\Big[\mathbb{E}_{Y|Z}[Y] \cdot \mathbb{E}_{X|Z}[X]\Big]\Big). \tag{18}
\end{aligned}$$

Now we analyze the convergence and variance separately. For the convergence part, we have

$$
\theta_t - \alpha_{t+1} \mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\Big]\theta_t + \alpha_{t+1}\mathbb{E}_Z\Big[\mathbb{E}_{Y|Z}[Y] \cdot \mathbb{E}_{X|Z}[X]\Big] - \theta_*
$$
$$
= \Big(I - \alpha_{t+1}\mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\Big]\Big)(\theta_t - \theta_*). \tag{19}
$$

For the variance part we have

$$
\mathbb{E}\Big[\|X_t' X_t^\top - \mathbb{E}_Z\big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\big]\|^2\Big]
$$
$$
= \mathbb{E}\Big[\|X_t' X_t^\top - \mathbb{E}_{X|Z_t}[X] \cdot \mathbb{E}_{X|Z_t}[X]^\top + \mathbb{E}_{X|Z_t}[X] \cdot \mathbb{E}_{X|Z_t}[X]^\top - \mathbb{E}_Z\big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\big]\|^2\Big]
$$
$$
\leq 2\mathbb{E}\Big[\|X_t' X_t^\top - \mathbb{E}_{X|Z_t}[X] \cdot \mathbb{E}_{X|Z_t}[X]^\top\|^2 + \|\mathbb{E}_{X|Z_t}[X] \cdot \mathbb{E}_{X|Z_t}[X]^\top - \mathbb{E}_Z\big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\big]\|^2\Big]
$$
$$
\leq 2C_x d_x^{\vartheta_1} + 2C_{xx} d_z^{\vartheta_3} =: \sigma_1^2 \tag{20}
$$

Similarly, we have

$$
\mathbb{E}\Big[\|Y_t X_t' - \mathbb{E}_Z\big[\mathbb{E}_{Y|Z}[Y] \cdot \mathbb{E}_{X|Z}[X]\big]\|^2\Big]
$$
$$
= \mathbb{E}\Big[\|Y_t X_t' - \mathbb{E}_{Y|Z_t}[Y] \cdot \mathbb{E}_{X|Z_t}[X]\|^2 + \|\mathbb{E}_{Y|Z_t}[Y] \cdot \mathbb{E}_{X|Z_t}[X] - \mathbb{E}_Z\big[\mathbb{E}_{Y|Z}[Y] \cdot \mathbb{E}_{X|Z}[X]\big]\|^2\Big]
$$
$$
\leq C_y d_x^{\vartheta_2} + C_{yx} d_z^{\vartheta_4} =: \sigma_2^2. \tag{21}
$$

Now we know from (18), (19), (20), and (21) that

$$
\|\theta_{t+1} - \theta_*\|^2 = \|A_t\|^2 + 2\alpha_{t+1}\langle A_t, B_t\rangle + \alpha_{t+1}^2\|B_t\|^2. \tag{22}
$$

where

$$
A_t = \Big(I - \alpha_{t+1}\mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\Big]\Big)(\theta_t - \theta_*)
$$
$$
B_t = -\Big(X_t' X_t^\top - \mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\Big]\Big)\theta_t + \Big(Y_t X_t' - \mathbb{E}_Z\Big[\mathbb{E}_{Y|Z}[Y] \cdot \mathbb{E}_{X|Z}[X]\Big]\Big).
$$

This implies

$$
\mathbb{E}_{\theta_{t+1}|\theta_t}\Big[\|\theta_{t+1} - \theta_*\|^2\Big]
$$
$$
= \|\Big(I - \alpha_{t+1}\mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\Big]\Big)(\theta_t - \theta_*)\|^2
$$
$$
+ \alpha_{t+1}^2\mathbb{E}_{X_t,X_t',Y_t,Z_t|\theta_t}\Big[\|\Big(X_t' X_t^\top - \mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\Big]\Big)\theta_t - \Big(Y_t X_t' - \mathbb{E}_Z\Big[\mathbb{E}_{Y|Z}[Y] \cdot \mathbb{E}_{X|Z}[X]\Big]\Big)\|^2\Big]
$$
$$
\leq (1 - \alpha_{t+1}\mu)^2\|\theta_t - \theta_*\|^2 + 3\alpha_{t+1}^2\Big(\sigma_1^2\|\theta_t - \theta_*\|^2 + \sigma_1^2\|\theta_*\|^2 + \sigma_2^2\|\theta_*\|^2\Big)
$$
$$
\leq ((1 - \alpha_{t+1}\mu)^2 + 3\alpha_{t+1}^2\sigma_1^2)\|\theta_t - \theta_*\|^2 + 3\alpha_{t+1}^2\sigma_1^2\|\theta_*\|^2 + 3\alpha_{t+1}^2\sigma_2^2\|\theta_*\|^2, \tag{23}
$$

where the first inequality uses Cauchy-Schwarz inequality, the definition of $\sigma_1, \sigma_2$ and Assumption 2.3. Choosing $\alpha_{t+1}$ such that

$$
((1 - \alpha_{t+1}\mu)^2 + 3\alpha_{t+1}^2\sigma_1^2) \leq 1 - \alpha_{t+1}\mu \Leftrightarrow \alpha \leq \frac{\mu}{\mu^2 + 3\sigma_1^2}
$$

and taking expectation on both sides of (23), we have

$$
\mathbb{E}\Big[\|\theta_{t+1} - \theta_*\|^2\Big] \leq (1 - \alpha_{t+1}\mu)\mathbb{E}\Big[\|\theta_t - \theta_*\|^2\Big] + 3\alpha_{t+1}^2\sigma_1^2\|\theta_*\|^2 + 3\alpha_{t+1}^2\sigma_2^2\|\theta_*\|^2.
$$

Now, we use the following result.

**Lemma 2.** *Suppose we have three sequences $\{a_t\}_{t=0}^{\infty}, \{b_t\}_{t=0}^{\infty}, \{r_t\}_{t=0}^{\infty}$ satisfying*

$$a_{t+1} \leq r_t a_t + b_t, r_t > 0 \tag{24}$$

*for any $t \geq 0$. Define $R_{t+1} = \prod_{i=0}^{t} r_i$, we have*

$$a_{t+1} \leq R_{t+1} a_0 + \sum_{i=0}^{t} \frac{R_{t+1} b_i}{R_{i+1}}.$$

By Lemma 2, we know

$$\mathbb{E}\Big[\|\theta_{t+1} - \theta_*\|^2\Big] \leq \prod_{i=0}^{t}(1 - \alpha_i \mu)\mathbb{E}\Big[\|\theta_0 - \theta_*\|^2\Big] + (3\sigma_1^2\|\theta_*\|^2 + 3\sigma_2^2\|\theta_*\|^2)\sum_{i=0}^{t}\alpha_i^2\prod_{j=i+1}^{t}(1 - \alpha_j \mu).$$

Now if we set $\alpha_i = \alpha$, we know

$$\mathbb{E}\Big[\|\theta_t - \theta_*\|^2\Big] \leq (1 - \alpha\mu)^t\mathbb{E}\Big[\|\theta_0 - \theta_*\|^2\Big] + \alpha^2\Big(\sum_{i=0}^{t}(1 - \alpha\mu)^i\Big)(3\sigma_1^2\|\theta_*\|^2 + 3\sigma_2^2\|\theta_*\|^2)$$

$$\leq e^{-t\alpha\mu}\mathbb{E}\Big[\|\theta_0 - \theta_*\|^2\Big] + \frac{\alpha}{\mu}(3\sigma_1^2\|\theta_*\|^2 + 3\sigma_2^2\|\theta_*\|^2)$$

Choosing $\alpha, T$ such that $\alpha = \frac{\log T}{\mu T} \leq \frac{\mu}{\mu^2 + 3\sigma_1^2}$, we know

$$\mathbb{E}\Big[\|\theta_T - \theta_*\|^2\Big] \leq \frac{\mathbb{E}\Big[\|\theta_0 - \theta_*\|^2\Big]}{T} + \frac{3\|\theta_*\|^2(\sigma_1^2 + \sigma_2^2)\log T}{\mu^2 T}.$$

$\square$

## D.2 Proof of Lemma 2

*Proof.* We notice from (24) that for any $0 \leq i \leq t$, we have

$$\frac{a_{i+1}}{R_{i+1}} \leq \frac{a_i}{R_i} + \frac{b_i}{R_{i+1}}.$$

Taking summation on both sides, we have

$$\frac{a_{t+1}}{R_{t+1}} \leq \frac{a_0}{R_0} + \sum_{i=0}^{t}\frac{b_i}{R_{i+1}}$$

which completes the proof by multiplying $R_{t+1}$ on both sides. $\square$

## D.3 Proof of Lemma 1

*Proof.* We first notice that Assumption 2.2 holds since

$$\mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X] \cdot \mathbb{E}_{X|Z}[X]^\top\Big] = \mathbb{E}_Z\Big[\phi(\gamma_*^\top Z) \cdot \phi(\gamma_*^\top Z)^\top\Big] \succeq \mu I.$$

20

For (5) and (6), we have

$$
\begin{aligned}
&\mathbb{E}\Big[\|X'X^\top - \mathbb{E}_{X|Z}[X]\mathbb{E}_{X|Z}[X]^\top\|^2\Big] \\
=&\mathbb{E}\Big[\|\epsilon_2'\phi(\gamma_*^\top Z)^\top + \phi(\gamma_*^\top Z)\epsilon_2^\top + \epsilon_2'\epsilon_2^\top\|^2\Big] \\
\leq&3\mathbb{E}\Big[\|\epsilon_2'\phi(\gamma_*^\top Z)^\top\|^2 + \|\phi(\gamma_*^\top Z)\epsilon_2^\top\|^2 + \|\epsilon_2'\epsilon_2^\top\|^2\Big] \\
=&3\mathbb{E}\Big[\|\phi(\gamma_*^\top Z)\epsilon_2'^\top \epsilon_2'\phi(\gamma_*^\top Z)^\top\| + \|\phi(\gamma_*^\top Z)\epsilon_2^\top \epsilon_2\phi(\gamma_*^\top Z)^\top\| + |\epsilon_2^\top \epsilon_2'|^2\Big] = \mathcal{O}(d_x^2),
\end{aligned}
\tag{25}
$$

and

$$
\begin{aligned}
&\mathbb{E}\Big[\|YX - \mathbb{E}_{Y|Z}[Y]\mathbb{E}_{X|Z}[X]\|^2\Big] \\
=&\mathbb{E}\Big[\|X'X^\top\theta_* + \epsilon_1 X - \mathbb{E}_{X|Z}[X]\mathbb{E}_{X|Z}[X]^\top\theta_*\|^2\Big] \\
\leq&2\mathbb{E}\Big[\|X'X^\top\theta_* - \mathbb{E}_{X|Z}[X]\mathbb{E}_{X|Z}[X]^\top\theta_*\|^2\Big] + 2\mathbb{E}\Big[\epsilon_1^2\|\phi(\gamma_*^\top Z) + \epsilon_2\|^2\Big] \\
=&\mathcal{O}(\|\theta_*\|^2\sigma_{\epsilon_2}^2 d_x^2 + \sigma_{\epsilon_1}^2 d_x + \sigma_{\epsilon_1,\epsilon_2}^2 d_x),
\end{aligned}
$$

where the first inequality uses Cauchy-Schwarz inequality, and the second equality uses (9), (10) and (25). For (7) we have

$$
\begin{aligned}
&\mathbb{E}\Big[\|\mathbb{E}_{X|Z}[X]\cdot\mathbb{E}_{X|Z}[X]^\top - \mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X]\cdot\mathbb{E}_{X|Z}[X]^\top\Big]\|^2\Big] \\
=&\mathbb{E}\Big[\|\phi(\gamma_*^\top Z)\phi(\gamma_*^\top Z)^\top - \mathbb{E}\Big[\phi(\gamma_*^\top Z)\phi(\gamma_*^\top Z)^\top\Big]\|^2\Big] = \mathcal{O}(d_z)
\end{aligned}
$$

where the last equality uses (10). Using the above conclusion in (8), we have

$$
\begin{aligned}
&\mathbb{E}\Big[\|\mathbb{E}_{Y|Z}[Y]\cdot\mathbb{E}_{X|Z}[X] - \mathbb{E}_Z\Big[\mathbb{E}_{Y|Z}[Y]\cdot\mathbb{E}_{X|Z}[X]\Big]\|^2\Big] \\
=&\mathbb{E}\Big[\|\mathbb{E}_{X|Z}[X]\cdot\mathbb{E}_{X|Z}[X]^\top\theta_* - \mathbb{E}_Z\Big[\mathbb{E}_{X|Z}[X]\cdot\mathbb{E}_{X|Z}[X]^\top\theta_*\Big]\|^2\Big] = \mathcal{O}(\|\theta_*\|^2 d_z).
\end{aligned}
$$

$\square$

# E   Proofs for Section 3

## E.1   Proof of Theorem 2

*Proof of Theorem 2 .*  Recall that $\xi_{Z_t}$, and $\xi_{Z_t Y_t}$ are the i.i.d. noise sequences

$$
\begin{aligned}
\xi_{Z_t} &= \Sigma_Z - Z_t Z_t^\top, \\
\xi_{Z_t Y_t} &= \Sigma_{ZY} - Z_t Y_t.
\end{aligned}
$$

Note $\gamma_*$, and $\theta_*$ can be written as $\gamma_* = \Sigma_Z^{-1}\Sigma_{ZX} \in \mathbb{R}^{d_z \times d_x}$, and $\theta_* = (\gamma_*^\top \Sigma_z \gamma_*)^{-1}\gamma_*^\top \Sigma_{ZY} \in \mathbb{R}^{d_x}$ which we are going to use throughout the proof.

To quantify the bias, we use the following bound on $\mathbb{E}\left[\|\gamma_t - \gamma_*\|_2^k\right]$, $k = 1, 2, 4$, proved in Lemma 3.2 of Chen et al. (2020).

**Lemma 3.** *Suppose Assumption 2.4, and Assumption 3.2 hold. Then we have*

$$
\mathbb{E}\left[\|\gamma_t - \gamma_*\|^k\right] = O\left(\sqrt{d_z^k \beta_t^k}\right) \quad \text{for} \quad k = 1, 2, 4.
\tag{26}
$$

We proceed by noting that if $\gamma_*$, $\Sigma_Z$, and $\Sigma_{ZY}$ were known beforehand, one could use the following deterministic gradient updates to obtain $\theta_*$.

$$\widetilde{\theta}_{t+1} = \widetilde{\theta}_t - \alpha_{t+1}\gamma_*^\top \left( \Sigma_Z \gamma_* \widetilde{\theta}_t - \Sigma_{ZY} \right). \tag{27}$$

**Lemma 4.** *Let Assumption 2.2 be true. Then, choosing $\eta_k = O(k^{-a})$ with $1/2 < a < 1$, we have $\|\widetilde{\theta}_t - \theta_*\| = O\left(\exp(-t^{1-a})\right)$.*

Define the sequence $\delta_t := \theta_t - \widetilde{\theta}_t$. We will establish the convergence rate of $\mathbb{E}\left[\|\delta_t\|_2^2\right]$. From (13), and (27), we have the following expansion of $\delta_{t+1}$.

$$\delta_{t+1} = Q_t \delta_t + \alpha_{t+1} D_t \theta_t + \alpha_{t+1}(\gamma_t - \gamma_*)^\top \Sigma_{ZY} - \alpha_{t+1}\gamma_t^\top \xi_{Z_t Y_t} + \alpha_{t+1}\gamma_t^\top \xi_{Z_t}\gamma_t\theta_t, \tag{28}$$

where

$$Q_t := (I - \alpha_{t+1}\gamma_*^\top \Sigma_Z \gamma_*),$$
$$D_t := \gamma_*^\top \Sigma_Z \gamma_* - \gamma_t^\top \Sigma_Z \gamma_t.$$

First we will establish an intermediate bound on $\mathbb{E}\left[\|\delta_t\|^2\right]$. To do so, we will need the following result which shows that $\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]$ is bounded for all $t$ which we prove in Section E.2.

**Lemma 5 (Boundedness of fourth moment of $\|\theta_t - \theta_*\|$).** *Let the conditions in Theorem 2 be true. Then, choosing $\alpha_t, \beta_t$ such that $\alpha_t \leq d_z^{-4\varkappa - \vartheta/2}$, and $\sum_{t=1}^\infty (\alpha_t^2 + \alpha_t\sqrt{\beta_t}) < \infty$, we have $\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]$ is bounded by some constant $M > 0$.*

**Lemma 6 (Intermediate bound on $\mathbb{E}[\|\delta_t\|_2^2]$).** *Let the conditions in Theorem 2 be true. We have the following intermediate bound on $\mathbb{E}\left[\|\delta_t\|^2\right]$:*

$$\mathbb{E}\left[\|\delta_t\|^2\right] = O\left(\beta_t d_z^{1+2\varkappa} + \alpha_{t+1} d_z^{4\varkappa + \vartheta/2} + \sqrt{d_z \beta_t}\right). \tag{29}$$

*Proof of Lemma 6.* Recall the update for $\delta_{t+1}$ obtained in (28).

$$\delta_{t+1} = Q_t \delta_t + \alpha_{t+1} D_t \theta_t + \alpha_{t+1}(\gamma_t - \gamma_*)^\top \Sigma_{ZY} - \alpha_{t+1}\gamma_t^\top \xi_{Z_t Y_t} + \alpha_{t+1}\gamma_t^\top \xi_{Z_t}\gamma_t\theta_t.$$

Then,

$$\begin{aligned}
\|\delta_{t+1}\|_2^2 =& \delta_t^\top Q_t^2 \delta_t + \alpha_{t+1}^2 \|D_t\theta_t + (\gamma_t - \gamma_*)^\top \Sigma_{ZY} - \gamma_t^\top \xi_{Z_t Y_t} + \gamma_t^\top \xi_{Z_t}\gamma_t\theta_t\|^2 \\
&+ 2\alpha_{t+1}\delta_t^\top Q_t \left(D_t\theta_t + (\gamma_t - \gamma_*)^\top \Sigma_{ZY}\right) \\
&+ 2\alpha_{t+1}\delta_t^\top Q_t \left(\gamma_t^\top \xi_{Z_t}\gamma_t\theta_t - \gamma_t^\top \xi_{Z_t Y_t}\right).
\end{aligned} \tag{30}$$

Then, choosing $\alpha_1(\|\gamma_*\|_2 \lambda_Z)^2 < 1$, using Young's inequality and Assumption 2.4, from (30) we get,

$$\begin{aligned}
\mathbb{E}_t\left[\|\delta_{t+1}\|_2^2\right] \leq& (1 - \alpha_{t+1}\mu)\|\delta_t\|^2 + 4\alpha_{t+1}^2 \left(\|D_t\theta_t\|^2 + \|(\gamma_t - \gamma_*)^\top \Sigma_{ZY}\|^2\right) \\
&+ 4\alpha_{t+1}^2 \left(\|\gamma_t\|_2^2 \mathbb{E}\left[\|\xi_{Z_t Y_t}\|_2^2\right] + \|\gamma_t\|_2^4 \mathbb{E}\left[\|\xi_{Z_t}\|^2\right]\|\theta_t\|^2\right) \\
&+ 2\alpha_{t+1}\delta_t^\top Q_t \left(D_t\theta_t + (\gamma_t - \gamma_*)^\top \Sigma_{ZY}\right) \\
\lesssim& (1 - \alpha_{t+1}\mu)\|\delta_t\|^2 + 4\alpha_{t+1}^2 \left(\|D_t\|^2\|\theta_t\|^2 + \|(\gamma_t - \gamma_*)^\top \Sigma_{ZY}\|^2\right) \\
&+ 4C\alpha_{t+1}^2 \left(d_z^{2\varkappa + \vartheta/2} + d_z^{4\varkappa + \vartheta/2}\|\theta_t\|^2\right) + \\
&2\alpha_{t+1}\delta_t^\top Q_t \left(D_t\theta_t + (\gamma_t - \gamma_*)^\top \Sigma_{ZY}\right),
\end{aligned}$$

where the last inequality follows by Assumption 3.1, and Assumption 3.3.

Now, taking expectation on both sides, we obtain

$$
\begin{aligned}
\mathbb{E}\left[\|\delta_{t+1}\|_2^2\right] \lesssim & (1 - \alpha_{t+1}\mu)\mathbb{E}\left[\|\delta_t\|^2\right] + 4\alpha_{t+1}^2\left(\mathbb{E}\left[\|D_t\|^2\|\theta_t\|^2\right] + \mathbb{E}\left[\|(\gamma_t - \gamma_*)^\top\Sigma_{ZY}\|^2\right]\right) \\
& + 4C\alpha_{t+1}^2\left(d_z^{2\varkappa+\vartheta/2} + d_z^{4\varkappa+\vartheta/2}\mathbb{E}\left[\|\theta_t\|^2\right]\right) \\
& + 2\alpha_{t+1}\left(\mathbb{E}\left[|\delta_t^\top Q_t D_t\theta_t|\right] + \mathbb{E}\left[|\delta_t^\top Q_t(\gamma_t - \gamma_*)^\top\Sigma_{ZY}|\right]\right).
\end{aligned}
\tag{31}
$$

Now, the following bounds are true:

1. We have that

$$
\alpha_{t+1}^2\mathbb{E}\left[\|D_t\|^2\|\theta_t\|^2\right] \leq \alpha_{t+1}^2\sqrt{\mathbb{E}\left[\|D_t\|_2^4\right]\mathbb{E}\left[\|\theta_t\|_2^4\right]} \lesssim d_z^{1+2\varkappa}\alpha_{t+1}^2\beta_t,
\tag{32}
$$

where the first inequality follows by Cauchy-Schwarz inequality, the second inequality follows by (42), and Lemma 5.

2. Using $\Sigma_{ZY} = O(1)$, and Lemma 3, we get

$$
\alpha_{t+1}^2\mathbb{E}\left[\|(\gamma_t - \gamma_*)^\top\Sigma_{ZY}\|^2\right] \lesssim d_z\beta_t\alpha_{t+1}^2.
\tag{33}
$$

3. We have that

$$
\begin{aligned}
\alpha_{t+1}\mathbb{E}\left[|\delta_t^\top Q_t D_t\theta_t|\right] \leq & \alpha_{t+1}\mathbb{E}\left[\|\delta_t\|_2\|Q_t\|_2\|D_t\|_2\|\theta_t\|_2\right] \\
\leq & \frac{\alpha_{t+1}\mu}{16}\mathbb{E}\left[\|\delta_t\|^2\right] + \frac{4\alpha_{t+1}}{\mu}\sqrt{\mathbb{E}\left[\|D_t\|_2^4\right]\mathbb{E}\left[\|\theta_t\|_2^4\right]} \\
\lesssim & \frac{\alpha_{t+1}\mu}{16}\mathbb{E}\left[\|\delta_t\|^2\right] + \frac{4d_z^{1+2\varkappa}\alpha_{t+1}\beta_t}{\mu},
\end{aligned}
\tag{34}
$$

where the first inequality follows by Hölder's inequality, the second inequality follows by Young's inequality, Cauchy-Schwarz inequality, and $\|Q_t\|_2 < 1$, and the third inequality follows by (42), and Lemma 5.

4. Using $\|Q_t\|_2 < 1$, $\|\Sigma_{ZY}\|_2 = O(1)$, Cauchy-Schwarz inequality, and Lemma 3, we get,

$$
\begin{aligned}
& \alpha_{t+1}\mathbb{E}\left[|\delta_t^\top Q_t(\gamma_t - \gamma_*)^\top\Sigma_{ZY}|\right] \\
\lesssim & \alpha_{t+1}\mathbb{E}\left[\|\delta_t\|_2\|\gamma_t - \gamma_*\|_2\right] \\
\leq & \alpha_{t+1}\sqrt{\mathbb{E}\left[\|\delta_t\|_2^2\right]\mathbb{E}\left[\|\gamma_t - \gamma_*\|_2^2\right]} \\
\leq & \frac{\sqrt{d_z\beta_t}\alpha_{t+1}}{2} + \frac{\sqrt{d_z\beta_t}\alpha_{t+1}\mathbb{E}\left[\|\delta_t\|_2^2\right]}{2}.
\end{aligned}
\tag{35}
$$

Combining (31), (32), (33), (34), (35), and Lemma 5, we have

$$
\begin{aligned}
& \mathbb{E}\left[\|\delta_{t+1}\|_2^2\right] \\
\lesssim & (1 - \alpha_{t+1}\mu)\mathbb{E}\left[\|\delta_t\|^2\right] + 4\alpha_{t+1}^2\beta_t d_z^{1+2\varkappa} + 4C\alpha_{t+1}^2 d_z^{4\varkappa+\vartheta/2} \\
& + 2\alpha_{t+1}\left(\mu\mathbb{E}\left[\|\delta_t\|^2\right]/16 + 4d_z^{1+2\varkappa}\beta_t/\mu + \sqrt{d_z\beta_t}/2 + \sqrt{d_z\beta_t}\mathbb{E}\left[\|\delta_t\|_2^2\right]/2\right) \\
\lesssim & (1 - 7\mu\alpha_{t+1}/8 + \alpha_{t+1}\sqrt{d_z\beta_t})\mathbb{E}\left[\|\delta_t\|^2\right] + (8\alpha_{t+1}\beta_t d_z^{1+2\varkappa}/\mu + 4C\alpha_{t+1}^2 d_z^{4\varkappa+\vartheta/2}) \\
& + \alpha_{t+1}\sqrt{d_z\beta_t} \\
\lesssim & (1 - 3\mu\alpha_{t+1}/4)\mathbb{E}\left[\|\delta_t\|^2\right] + (8\alpha_{t+1}\beta_t d_z^{1+2\varkappa}/\mu + 4C\alpha_{t+1}^2 d_z^{4\varkappa+\vartheta/2}) + \alpha_{t+1}\sqrt{d_z\beta_t}.
\end{aligned}
\tag{36, 37}
$$

23

In the above, the third inequality follows by choosing $\beta_t \le \mu^2/(64 d_z)$, and $\alpha_{t+1}\sqrt{d_z \beta_t} < 1$. Then, from (37), we have

$$\mathbb{E}\left[\|\delta_t\|_2^2\right] = O\left(\beta_t d_z^{1+2\varkappa} + \alpha_{t+1} d_z^{4\varkappa+\vartheta/2} + \sqrt{d_z \beta_t}\right).$$

$\square$

Coming back to the proof of Theorem 2, observe that, we can sharpen the bound in (35) using Lemma 6 which allows us to avoid the use of Young's inequality. This leads to the following improved version of the recursion in (37) using which we can improve the term $\sqrt{d_z \beta_t}$ in (29) as follows:

$$\mathbb{E}\left[\|\delta_{t+1}\|_2^2\right] \lesssim (1 - 7\mu\alpha_{t+1}/8)\mathbb{E}\left[\|\delta_t\|^2\right]$$
$$+ \alpha_{t+1}O\left(\beta_t d_z^{1+2\varkappa} + \alpha_{t+1} d_z^{4\varkappa+\vartheta/2} + \sqrt{\alpha_{t+1}\beta_t}d_z^{1/2+2\varkappa+\vartheta/4} + (\beta_t d_z)^{3/4}\right)$$
$$= O\left(\beta_t d_z^{1+2\varkappa} + \alpha_{t+1} d_z^{4\varkappa+\vartheta/2} + \sqrt{\alpha_{t+1}\beta_t}d_z^{1/2+2\varkappa+\vartheta/4} + (\beta_t d_z)^{3/4}\right).$$

In fact, this trick can be used repeatedly to sharpen the bound even further as shown in Lemma 7.

**Lemma 7** (**Final improved bound on** $\mathbb{E}[\|\delta_t\|_2^2]$). *Let the conditions in Theorem 2 be true. Then using Lemma 6, we have,*

$$\mathbb{E}\left[\|\delta_{t+1}\|_2^2\right]$$
$$\lesssim O\left((d_z \beta_t)^{1-2^{-r-1}} + \sum_{i=0}^{r}\left(\alpha_{t+1}^{2^{-i}}\beta_t^{1-2^{-i}}d_z^{1+(4\varkappa+\vartheta/2-1)2^{-i}} + \beta_t(1+\alpha_{t+1}^{2^{-i}})d_z^{1+2^{1-i}\varkappa}\right)\right),$$

*where $r$ is any non-negative integer.*

*Proof of Lemma 7.* If we have

$$\mathbb{E}\left[\|\delta_t\|^2\right] = O\left(\alpha_{t+1}d_z^{4\varkappa+\vartheta/2} + \beta_t d_z^{1+2\varkappa} + \sqrt{d_z \beta_t}\right),$$

then from (35), we have,

$$\mathbb{E}\left[|\delta_t^\top Q_t(\gamma_t - \gamma_*)^\top \Sigma_{ZY}|\right] \lesssim \sqrt{\mathbb{E}\left[\|\delta_t\|_2^2\right]\mathbb{E}\left[\|\gamma_t - \gamma_*\|_2^2\right]}$$
$$= O\left(\sqrt{\alpha_{t+1}\beta_t}d_z^{1/2+2\varkappa+\vartheta/4} + \beta_t d_z^{1+\varkappa} + (d_z \beta_t)^{3/4}\right). \qquad (38)$$

Then, similar to (36), we have,

$$\mathbb{E}\left[\|\delta_{t+1}\|_2^2\right]$$
$$\lesssim (1 - \alpha_{t+1}\mu)\mathbb{E}\left[\|\delta_t\|^2\right] + 4\alpha_{t+1}^2 \beta_t d_z^{1+2\varkappa} + 4C\alpha_{t+1}^2 d_z^{4\varkappa+\vartheta/2}$$
$$+ 2\alpha_{t+1}\left(\mu\mathbb{E}\left[\|\delta_t\|^2\right]/16 + 4d_z^{1+2\varkappa}\beta_t/\mu + \sqrt{\alpha_{t+1}\beta_t}d_z^{1/2+2\varkappa+\vartheta/4} + \beta_t d_z^{1+\varkappa} + (d_z \beta_t)^{3/4}\right)$$
$$\lesssim (1 - 7\mu\alpha_{t+1}/8)\mathbb{E}\left[\|\delta_t\|^2\right]$$
$$+ \alpha_{t+1}O\left((d_z \beta_t)^{3/4} + \sum_{i=0}^{1}\left(\alpha_{t+1}^{2^{-i}}\beta_t^{1-2^{-i}}d_z^{1+(4\varkappa+\vartheta/2-1)2^{-i}} + \beta_t(1+\alpha_{t+1}^{2^{-i}})d_z^{1+2^{1-i}\varkappa}\right)\right)$$
$$= O\left((d_z \beta_t)^{3/4} + \sum_{i=0}^{1}\left(\alpha_{t+1}^{2^{-i}}\beta_t^{1-2^{-i}}d_z^{1+(4\varkappa+\vartheta/2-1)2^{-i}} + \beta_t(1+\alpha_{t+1}^{2^{-i}})d_z^{1+2^{1-i}\varkappa}\right)\right).$$

24

Now if we repeat this step $r$ number of times (where $r$ is to be set later), by progressive sharpening we get the following bound.

$$\mathbb{E}\left[\|\delta_{t+1}\|_2^2\right]$$
$$\lesssim O\left((d_z\beta_t)^{1-2^{-r-1}} + \sum_{i=0}^r \left(\alpha_{t+1}^{2^{-i}}\beta_t^{1-2^{-i}}d_z^{1+(4\varkappa+\vartheta/2-1)2^{-i}} + \beta_t(1+\alpha_{t+1}^{2^{-i}})d_z^{1+2^{1-i}\varkappa}\right)\right).$$

$\square$

Coming back to the proof of Theorem 2, we have that by combining Lemma 4, and Lemma 7,

$$\mathbb{E}\left[\|\theta_t - \theta_*\|^2\right] \leq 2\mathbb{E}\left[\|\delta_t\|^2\right] + 2\mathbb{E}\left[\|\widetilde{\theta}_t - \theta_*\|^2\right]$$
$$= O\left((d_z\beta_t)^{1-2^{-r-1}} + \sum_{i=0}^r \left(\alpha_{t+1}^{2^{-i}}\beta_t^{1-2^{-i}}d_z^{1+(4\varkappa+\vartheta/2-1)2^{-i}} + \beta_t(1+\alpha_{t+1}^{2^{-i}})d_z^{1+2^{1-i}\varkappa}\right)\right). \quad (39)$$

Now, in (39), for some arbitrarily small number $\iota > 0$, choosing

$$\alpha_t = \min(0.5d_z^{-4\varkappa-\vartheta/2}\lambda_Z^{-1}C_\gamma^{-2}, 0.5(\|\gamma_*\|_2\lambda_Z)^{-2})t^{-1+\iota/2}, \qquad \beta_t = \mu^2 d_z^{-1-2\varkappa}t^{-1+\iota/2}/128,$$

and setting $r = \lceil\log_2\left((\iota/2)^{-1} - 1\right) - 1\rceil$ we get,

$$\mathbb{E}\left[\|\theta_t - \theta_*\|^2\right] = O\left(\max\left(t^{-1+\iota}, t^{-1+\iota/2}\log((\iota/2)^{-1} - 1)\right)\right).$$

$\square$

## E.2 Proof of Lemma 5

*Proof.* Using the form of $\theta_*$, from (13) we get,

$$\theta_{t+1} - \theta_* = \widehat{Q}_t(\theta_t - \theta_*) + \alpha_{t+1}(\gamma_t - \gamma_*)^\top\Sigma_{ZY} + \alpha_{t+1}D_t\theta_* + \alpha_{t+1}\gamma_t^\top\xi_{Z_t}\gamma_t(\theta_t - \theta_*)$$
$$+ \alpha_{t+1}\gamma_t^\top\xi_{Z_t}\gamma_t\theta_* + \alpha_{t+1}\gamma_t^\top\xi_{Z_tY_t}. \quad (40)$$

where $\widehat{Q}_t := \left(I - \alpha_{t+1}\gamma_t^\top\Sigma_Z\gamma_t\right) = Q_t + \alpha_{t+1}D_t$. Recall that $D_t = \gamma_*^\top\Sigma_Z\gamma_* - \gamma_t^\top\Sigma_Z\gamma_t$. By Assumption 3.3, we have the following bound on $\|D_t\|_2$.

$$\|D_t\|_2 = O(\lambda_Z C_\gamma^2 d_z^{2\varkappa}). \quad (41)$$

We have the following bound on $\mathbb{E}\left[\|D_t\|_2^4\right]$ by Lemma 3.

$$\mathbb{E}\left[\|D_t\|_2^4\right] = \mathbb{E}\left[\|(\gamma_* - \gamma_t)^\top\Sigma_Z\gamma_* + \gamma_t^\top\Sigma_Z(\gamma_* - \gamma_t)\|_2^4\right] = O(d_z^{2+4\varkappa}\beta_t^2). \quad (42)$$

From (40), we have

$$\|\theta_{t+1} - \theta_*\|_2^2 \leq (\theta_t - \theta_*)^\top\widehat{Q}_t^2(\theta_t - \theta_*) + 3\alpha_{t+1}^2\|\gamma_t^\top\xi_{Z_t}\gamma_t(\theta_t - \theta_*)\|_2^2$$
$$+ 2\alpha_{t+1}(\theta_t - \theta_*)^\top\widehat{Q}_t(\gamma_t - \gamma_*)^\top\Sigma_{ZY}$$
$$+ 2\alpha_{t+1}(\theta_t - \theta_*)^\top\widehat{Q}_tD_t\theta_* + A_{1,t} + A_{2,t}, \quad (43)$$

where

$$
\begin{aligned}
A_{1,t} =\alpha_{t+1}^2 \big( &\|(\gamma_t - \gamma_*)^\top \Sigma_{ZY}\|_2^2 + \|D_t \theta_*\|_2^2 \\
&+ 2\Sigma_{ZY}^\top(\gamma_t - \gamma_*)D_t\theta_* + 3\|\gamma_t{}^\top \xi_{Z_t}\gamma_t\theta_*\|_2^2 + 3\|\gamma_t{}^\top \xi_{Z_tY_t}\|_2^2 \big),
\end{aligned}
\tag{44}
$$

and

$$
\begin{aligned}
A_{2,t} =2\alpha_{t+1}(&\widehat{Q}_t(\theta_t - \theta_*) + \alpha_{t+1}(\gamma_t - \gamma_*)^\top \Sigma_{ZY} \\
&+ \alpha_{t+1}D_t\theta_*)^\top(\gamma_t{}^\top \xi_{Z_t}\gamma_t(\theta_t - \theta_*) + \gamma_t{}^\top \xi_{Z_t}\gamma_t\theta_* + \gamma_t{}^\top \xi_{Z_tY_t}).
\end{aligned}
$$

Define

$$
\begin{aligned}
A_{3,t} :=&3\alpha_{t+1}^2\|\gamma_t{}^\top \xi_{Z_t}\gamma_t(\theta_t - \theta_*)\|_2^2 + 2\alpha_{t+1}(\theta_t - \theta_*)^\top \widehat{Q}_t(\gamma_t - \gamma_*)^\top \Sigma_{ZY} \\
&+ 2\alpha_{t+1}(\theta_t - \theta_*)^\top \widehat{Q}_t D_t\theta_* + A_{1,t} + A_{2,t}.
\end{aligned}
\tag{45}
$$

Then, choosing $C_\gamma^2 d_z^{2\varkappa}\lambda_Z\alpha_{t+1} < 1$, which ensures $\|\widehat{Q}_t\| \le 1$, we have

$$
\|\theta_{t+1} - \theta_*\|_2^4 \le \|\theta_t - \theta_*\|_2^4 + 2(\theta_t - \theta_*)^\top \widehat{Q}_t^2(\theta_t - \theta_*)A_{3,t} + A_{3,t}^2.
\tag{46}
$$

We now have the following bounds:

1. Using Assumption 3.1, and Assumption 3.3,

$$
\alpha_{t+1}^4 \mathbb{E}\left[\|\gamma_t{}^\top \xi_{Z_t}\gamma_t(\theta_t - \theta_*)\|_2^4\right] \lesssim d_z^{8\varkappa+\vartheta}\alpha_{t+1}^4 \mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right].
\tag{47}
$$

2. We have that

$$
\begin{aligned}
&\mathbb{E}\left[((\theta_t - \theta_*)^\top \widehat{Q}_t(\gamma_t - \gamma_*)^\top \Sigma_{ZY})^2\right] \\
&\lesssim \mathbb{E}\left[\|\theta_t - \theta_*\|^2\|\gamma_t - \gamma_*\|^2\right] \\
&\le \sqrt{\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]\mathbb{E}\left[\|\gamma_t - \gamma_*\|_2^4\right]} \\
&\le d_z\beta_t\left(1 + \mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]\right)/2,
\end{aligned}
\tag{48}
$$

where, the first inequality follows by $\|\widehat{Q}_t\|_2 = O(1)$, and $\|\Sigma_{ZY}\|_2 = O(1)$. The second inequality follows by Cauchy-Schwarz inequality. The last inequality follows by $\sqrt{ab} \le (a + b)/2$, and Lemma 3.

3. We have that

$$
\begin{aligned}
&\mathbb{E}\left[((\theta_t - \theta_*)^\top \widehat{Q}_t D_t\theta_*)^2\right] \\
&\lesssim \mathbb{E}\left[\|\theta_t - \theta_*\|_2^2\|D_t\|_2^2\right] \\
&\le \sqrt{\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]\mathbb{E}\left[\|D_t\|_2^4\right]} \\
&\lesssim d_z^{1+2\varkappa}\beta_t\left(1 + \mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]\right)/2,
\end{aligned}
\tag{49}
$$

where, the first inequality follows by $\|\widehat{Q}_t\|_2 = O(1)$, and $\|\theta_*\|_2 = O(1)$. The second inequality follows by Cauchy-Schwarz inequality. The last inequality follows by $\sqrt{ab} \le (a + b)/2$, and (42).

4. Using Assumption 3.1, Assumption 3.3, (42), and Lemma 3, we have

$$\mathbb{E}\left[A_{1,t}^2\right] = O\left(d_z^{8\varkappa+\vartheta}\alpha_{t+1}^4\right). \tag{50}$$

.

5. Using Young's inequality, Assumption 3.1, Assumption 3.3, Lemma 3, $\|\Sigma_{ZY}\|_2 = O(1)$, $\|\theta_*\|_2 = O(1)$, and (42), we have

$$\begin{aligned}
\mathbb{E}\left[A_{2,t}^2\right] \leq &2\alpha_{t+1}^2\mathbb{E}\left[\|\widehat{Q}_t(\theta_t - \theta_*) + \alpha_{t+1}(\gamma_t - \gamma_*)^\top\Sigma_{ZY} + \alpha_{t+1}D_t\theta_*\|_2^4\right] \\
&+ 2\alpha_{t+1}^2\mathbb{E}\left[\|\gamma_t^\top\xi_{Z_t}\gamma_t(\theta_t - \theta_*) + \gamma_t^\top\xi_{Z_t}\gamma_t\theta_* + \gamma_t^\top\xi_{Z_tY_t}\|_2^4\right] \\
\lesssim &\alpha_{t+1}^2 d_z^{8\varkappa+\vartheta}(1 + \mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]).
\end{aligned} \tag{51}$$

6. Using $\|\widehat{Q}_t\|_2 = O(1)$, Assumption 3.1, and Assumption 3.3,

$$\alpha_{t+1}^2\mathbb{E}\left[(\theta_t - \theta_*)^\top\widehat{Q}_t^2(\theta_t - \theta_*)\|\gamma_t^\top\xi_{Z_t}\gamma_t(\theta_t - \theta_*)\|_2^2\right] \lesssim \alpha_{t+1}^2 d_z^{4\varkappa+\vartheta/2}\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]. \tag{52}$$

7. We have that

$$\begin{aligned}
&\alpha_{t+1}\mathbb{E}\left[|(\theta_t - \theta_*)^\top\widehat{Q}_t^2(\theta_t - \theta_*)(\theta_t - \theta_*)^\top\widehat{Q}_t(\gamma_t - \gamma_*)^\top\Sigma_{ZY}|\right] \\
\lesssim &\alpha_{t+1}\mathbb{E}\left[\|\theta_t - \theta_*\|_2^3\|\gamma_t - \gamma_*\|_2\right] \\
\leq &\alpha_{t+1}\left(\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]\right)^{3/4}\left(\mathbb{E}\left[\|\gamma_t - \gamma_*\|_2^4\right]\right)^{1/4} \\
\leq &\alpha_{t+1}\sqrt{d_z\beta_t}\left(\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]\right)^{3/4} \\
\leq &\frac{3\alpha_{t+1}\sqrt{d_z\beta_t}}{4}\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right] + \frac{\alpha_{t+1}\sqrt{d_z\beta_t}}{4},
\end{aligned} \tag{53}$$

where, the first inequality follows by $\|\widehat{Q}_t\|_2 = O(1)$, and $\|\Sigma_{ZY}\|_2 = O(1)$, the second inequality follows by Cauchy-Schwarz inequality, the third inequality follows by Lemma 3 and the fourth inequality follows by Young's inequality.

8. Similar to (53), we have,

$$\begin{aligned}
&\alpha_{t+1}\mathbb{E}\left[|(\theta_t - \theta_*)^\top\widehat{Q}_t^2(\theta_t - \theta_*)(\theta_t - \theta_*)^\top\widehat{Q}_tD_t\theta_*|\right] \\
\leq &\frac{3d_z^{1/2+\varkappa}\alpha_{t+1}\sqrt{\beta_t}}{4}\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right] + \frac{d_z^{1/2+\varkappa}\alpha_{t+1}\sqrt{\beta_t}}{4}.
\end{aligned} \tag{54}$$

9. Using $\|\widehat{Q}_t\|_2 = O(1)$, Cauchy-Schwarz inequality, (50), and Young's inequality,

$$\begin{aligned}
&\mathbb{E}\left[(\theta_t - \theta_*)^\top\widehat{Q}_t^2(\theta_t - \theta_*)A_{1,t}\right] \\
\leq &\mathbb{E}\left[\|\theta_t - \theta_*\|_2^2 A_{1,t}\right] \\
\leq &\sqrt{\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]\mathbb{E}\left[A_{1,t}^2\right]} \\
\lesssim &d_z^{4\varkappa+\vartheta/2}\alpha_{t+1}^2\left(1 + \mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]\right).
\end{aligned} \tag{55}$$

27

10. By Assumption 2.4, we have,

$$\mathbb{E}_t\left[(\theta_t - \theta_*)^\top \widehat{Q}_t^2 (\theta_t - \theta_*) A_{2,t}\right] = 0. \tag{56}$$

Now using Jensen's inequality, and combining (47), (48), (49), (50), and (51), we have,

$$
\begin{aligned}
\mathbb{E}\left[A_{3,t}^2\right] \leq & 45\alpha_{t+1}^4 \mathbb{E}\left[\|\gamma_t^\top \xi_{Z_t}\gamma_t(\theta_t - \theta_*)\|_2^4\right] + 20\alpha_{t+1}^2 \mathbb{E}\left[((\theta_t - \theta_*)^\top \widehat{Q}_t(\gamma_t - \gamma_*)^\top \Sigma_{ZY})^2\right] \\
& + 20\alpha_{t+1}^2 \mathbb{E}\left[((\theta_t - \theta_*)^\top \widehat{Q}_t D_t \theta_*)^2\right] + 5\mathbb{E}\left[A_{1,t}^2\right] + 5\mathbb{E}\left[A_{2,t}^2\right] \\
\lesssim & \alpha_{t+1}^4 d_z^{\vartheta 7 + 8\varkappa} \mathbb{E}\left[\|\theta_t - \theta^*\|_2^4\right] + d_z \alpha_{t+1}^2 \beta_t \left(1 + \mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]\right) \\
& + \alpha_{t+1}^2 d_z^{1+2\varkappa}\beta_t \left(1 + \mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]\right) + d_z^{8\varkappa + \vartheta}\alpha_{t+1}^4 + \alpha_{t+1}^2 d_z^{8\varkappa+\vartheta}(1 + \mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]) \\
\lesssim & \alpha_{t+1}^2 d_z^{8\varkappa+\vartheta}\left(1 + \mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]\right). \tag{57}
\end{aligned}
$$

Combining (52), (53), (54), (55), and (56), we get,

$$
\begin{aligned}
& \mathbb{E}\left[(\theta_t - \theta_*)^\top \widehat{Q}_t^2(\theta_t - \theta_*)A_{3,t}\right] \\
& \lesssim \alpha_{t+1}^2 d_z^{4\varkappa+\vartheta/2}\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right] + \frac{3\alpha_{t+1}\sqrt{d_z\beta_t}}{4}\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right] + \frac{\alpha_{t+1}\sqrt{d_z\beta_t}}{4} \\
& + \frac{3d_z^{1/2+\varkappa}\alpha_{t+1}\sqrt{\beta_t}}{4}\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right] + \frac{d_z^{1/2+\varkappa}\alpha_{t+1}\sqrt{\beta_t}}{4} + d_z^{4\varkappa+\vartheta/2}\alpha_{t+1}^2\left(1 + \mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]\right) \\
& \lesssim (\alpha_{t+1}^2 d_z^{4\varkappa+\vartheta/2} + \alpha_{t+1}\sqrt{\beta_{t+1}}d_z^{1/2+\varkappa})(1 + \|\theta_t - \theta^*\|_2^4). \tag{58}
\end{aligned}
$$

Combining (46), (57), and (58), we have,

$$\mathbb{E}\left[\|\theta_{t+1} - \theta_*\|_2^4\right] \lesssim (1 + \alpha_{t+1}^2 d_z^{8\varkappa+\vartheta} + \alpha_{t+1}\sqrt{\beta_{t+1}}d_z^{1/2+\varkappa})\left(1 + \mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right]\right). \tag{59}$$

Now choosing $\alpha_t, \beta_t$ such that $\alpha_t \leq d_z^{-4\varkappa-\vartheta/2}$, and $\sum_{t=1}^\infty (\alpha_{t+1}^2 + \alpha_{t+1}\sqrt{\beta_{t+1}}) < \infty$, we get

$$\mathbb{E}\left[\|\theta_t - \theta_*\|_2^4\right] \leq M, \tag{60}$$

for some constant $0 \leq M < \infty$. $\qquad\square$

## E.3 Comment on the convergence of (11)

We now discuss the convergence properties of the update sequence (11), which we refer to as the *conditional stochastic optimization* (CSO) based updates, which we restate below:

$$\theta_{t+1} = \theta_t - \alpha_{t+1}\gamma_t^\top Z_t(X_t^\top \theta_t - Y_t), \qquad \gamma_{t+1} = \gamma_t - \beta_{t+1}Z_t(Z_t^\top \gamma_t - X_t^\top).$$

Similar to (40), for the above updates, we have the following expansion:

$$
\begin{aligned}
\theta_{t+1} - \theta_* = & \widehat{Q}_t(\theta_t - \theta_*) + \alpha_{t+1}(\gamma_t - \gamma_*)^\top \Sigma_{ZY} + \alpha_{t+1}D_t\theta_* + \alpha_{t+1}\gamma_t^\top \xi_{Z_t}\gamma_*(\theta_t - \theta_*) \\
& + \alpha_{t+1}\gamma_t^\top \xi_{Z_t}\gamma_*\theta_* + \alpha_{t+1}\gamma_t^\top \xi_{Z_t Y_t} - \alpha_{t+1}\gamma_t^\top Z_t\epsilon_{2,t}^\top \theta_t,
\end{aligned}
$$

where $\xi_{Z_t} = \Sigma_Z - Z_t Z_t^\top$, $\xi_{Z_t Y_t} = \Sigma_{ZY} - Z_t Y_t$, $\widehat{Q}_t := \left(I - \alpha_{t+1}\gamma_t^\top \Sigma_Z \gamma_*\right) = Q_t + \alpha_{t+1}D_t$, and $D_t = (\gamma_* - \gamma_t)\Sigma_Z\gamma_*$.

Recall that the reason for the initial divergence of the updates in (11) are the potential negative eigenvalues of $\gamma_t^\top \Sigma_Z \gamma_*$. Here we will show that if $\gamma_t^\top \Sigma_Z \gamma_*$ is positive semi-definite or $\gamma_t$ is close enough to $\gamma_*$ such that the negative eigenvalues (if any) are not too large in absolute values, then the updates in (11) indeed exhibit the same convergence rate as Algorithm 2.

**Assumption E.1.** *Let either of the following two conditions be true. For all $t \geq t_0$,*

1. *$\gamma_t \Sigma_Z \gamma_*$ is positive semidefinite.*
2. *$\|\gamma_t - \gamma_*\|^2 \lesssim d_z \beta_t$.*

Note that Condition 1 of Assumption E.1 is an idealized condition which is difficult to ensure for all $t$ in reality. But of course if this is true, then $\gamma_t \Sigma_Z \gamma_*$ does not have a negative eigenvalue to cause divergence and the proof then follows exactly like Lemma 5.

Hence, we will focus on the more realistic Condition 2 of Assumption E.1 which holds true almost surely Polyak and Juditsky (1992). Since we are interested in the asymptotic rate of convergence of CSO updates (due to the requirement of Assumption E.1), we will only concentrate on the iterations $t \geq t_0$. In this case, the proof steps are similar to Theorem 2 except for two major differences, that we discuss below.

**Difference 1: Potential negative definiteness of $\gamma_t^\top \Sigma_Z \gamma_*$:**

Under Condition 2, $\gamma_t^\top \Sigma_Z \gamma_*$ can indeed be negative definite. In general, if $\gamma_t^\top \Sigma_Z \gamma_*$ is negative definite then that is undesirable as we explain Section 3. In terms of the proof, we can no longer write $(\theta_t - \theta^*)^\top \widehat{Q}_t^\top \widehat{Q}_t (\theta_t - \theta^*) \leq \|\theta_t - \theta_*\|^2$ (which was possible to do in (43) in the proof of Lemma 5). Subsequently, (46) breaks down. But we will show that under Condition 2 the negative eigenvalues are not too large in terms of absolute values. Specifically, we can write,

$$
\begin{aligned}
&(\theta_t - \theta^*)^\top \widehat{Q}_t^\top \widehat{Q}_t (\theta_t - \theta^*) \\
=&(\theta_t - \theta^*)^\top (Q_t^2 + \alpha_{t+1} Q_t^\top D_t + \alpha_{t+1} D_t^\top Q_t + \alpha_{t+1}^2 D_t^\top D_t)(\theta_t - \theta^*) \\
\leq&(1 + 2\alpha_{t+1}\|D_t\|)\|\theta_t - \theta_*\|^2 + \alpha_{t+1}^2\|D_t\|^2\|\theta_t - \theta_*\|^2 \\
\leq&(1 + 2\alpha_{t+1}\sqrt{d_z\beta_t})\|\theta_t - \theta_*\|^2 + \alpha_{t+1}^2\|D_t\|^2\|\theta_t - \theta_*\|^2.
\end{aligned} \tag{61}
$$

The term $\alpha_{t+1}^2\|D_t\|^2\|\theta_t - \theta_*\|^2$ is of the order of $A_{3,t}$ defined in (45). Now $\alpha_{t+1}\sqrt{d_z\beta_t}$ is small enough in the sense that we choose the stepsizes such that $\sum_{t=1}^\infty (\alpha_{t+1}^2 + \alpha_{t+1}\sqrt{\beta_t}) < \infty$. Using this one can now show a similar bound as (59) and consequently show $\mathbb{E}\left[\|\theta_t - \theta_*\|^4\right]$ is bounded.

Now let us see what happens in the absence of Condition 2. Here one could use the fact $(1+2\alpha_{t+1}\|D_t\|) \lesssim (1 + 2C_\gamma \alpha_{t+1} d_z^\varkappa)$ which is too big. Recall that we want something at least of the order of $\alpha_{t+1}\sqrt{\beta_t}$ to show that $\theta_t$ sequence is bounded. One could also try to use the fact that $\mathbb{E}\left[\|D_t\|\right]$ is small by Lemma 3. But since $D_t$ and $\theta_t$ are interdependent, one needs to decouple them. One way to do this would be to use Cauchy-Shwarz inequalityas shown below.

$$
\mathbb{E}\left[\|D_t\|\|\theta_t - \theta_*\|^2\right] \leq \sqrt{\mathbb{E}\left[\|D_t\|^2\right]\mathbb{E}\left[\|\theta_t - \theta_*\|^4\right]} \lesssim \sqrt{d_z\beta_t \mathbb{E}\left[\|\theta_t - \theta_*\|^4\right]}.
$$

But that leads to the presence of $\mathbb{E}\left[\|\theta_t - \theta_*\|^4\right]$ in (43) which is potentially problematic due to the fact that on the left-hand side we have $\mathbb{E}\left[\|\theta_{t+1} - \theta_*\|^2\right]$.

**Difference 2: Presence of additional error term $\alpha_{t+1}\gamma_t^\top Z_t \epsilon_{2,t}^\top \theta_t$:**

When comparing (12) with (40), yet another crucial difference is the presence of the term $\alpha_{t+1}\gamma_t^\top Z_t \epsilon_{2,t}^\top \theta_t$. We will show by the following observations that this error term gets absorbed by other terms already present in (40) without affecting the convergence rate. Specifically, the following holds.

1. Using the independence between $Z$, and $\epsilon_{2,t}$, and by Assumption 2.4, we have,

$$
\begin{aligned}
&\mathbb{E}_t[(\widehat{Q}_t(\theta_t - \theta_*) + \alpha_{t+1}(\gamma_t - \gamma_*)^\top \Sigma_{ZY} + \alpha_{t+1}D_t\theta_* + \alpha_{t+1}\gamma_t^\top \xi_{Z_t}\gamma_*(\theta_t - \theta_*) \\
&+ \alpha_{t+1}\gamma_t^\top \xi_{Z_t}\gamma_*\theta_*)^\top \gamma_t^\top Z_t \epsilon_{2,t}^\top \theta_t] = 0.
\end{aligned}
$$

2. We also have that

$$\alpha_{t+1}^2 \mathbb{E}_t \left[ (\gamma_t^\top \xi_{Z_t Y_t})^\top \gamma_t^\top Z_t \epsilon_{2,t}^\top \theta_t \right]$$
$$= \alpha_{t+1}^2 (\gamma_t^\top \Sigma_Z \gamma_t \|\theta_*\|^2 + \gamma_t^\top \Sigma_Z \gamma_t \theta_*^\top (\theta_t - \theta_*))$$
$$\leq \alpha_{t+1}^2 (\gamma_t^\top \Sigma_Z \gamma_t \|\theta_*\|^2 + \|\gamma_t^\top \Sigma_Z \gamma_t (\theta_t - \theta_*)\|^2 + \|\theta_*\|^2)$$

This shows that the above term is of the same order as $A_{1,t}$ and $A_{3,t}$ defined in (44), and (45).

3. Finally, we have

$$\alpha_{t+1}^2 \mathbb{E}_t \left[ \|\gamma_t^\top Z_t \epsilon_{2,t}^\top \theta_t\|^2 \right] \lesssim \alpha_{t+1}^2 (\|\gamma_t\|^2 \|\theta_t - \theta_*\|^2 + \|\gamma_t\|^2 \|\theta_*\|^2).$$

So this term is of the order of $A_{3,t}$ as well.

Combining the above facts and following similar procedure as the proof of Theorem 2, one can show that the CSO updates achieve a similar rate under additional Assumption E.1.