DATA VALUE ESTIMATION ON PRIVATE GRADIENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

For gradient-based machine learning (ML) methods commonly adopted in practice such as stochastic gradient descent, the *de facto* differential privacy (DP) technique is perturbing the gradients with random Gaussian noise. Data valuation attributes the ML performance to the training data and is widely used in privacy-aware applications that require enforcing DP such as data pricing, collaborative ML, and federated learning (FL). *Can existing data valuation methods still be used when DP is enforced via gradient perturbations?* We show that the answer is no with the default approach of injecting i.i.d. random noise to the gradients because the *estimation uncertainty* of the data value estimates almost like random guesses. To address this issue, we propose to instead inject carefully correlated noise to provably remove the linear scaling of estimation uncertainty w.r.t. the budget with some assumptions on the gradient distribution. We also empirically demonstrate that our method gives better data value estimates on various ML tasks and is applicable to use cases including dataset valuation and FL.

024

026

000

001 002 003

004

006 007

008 009

010

011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

With growing data privacy regulations (Bukaty, 2019; Council of European Union, 2014) and
machine learning (ML) model attacks (Shokri et al., 2017), privacy has become a primary concern
in many scenarios such as collaborative ML (Sim et al., 2020; Xu et al., 2021) and federated
learning (FL) (McMahan et al., 2017; Yang et al., 2019). Differential privacy (DP) (Dwork & Roth, 2014) is commonly adopted as the *de facto* framework to provide privacy protection for training data
with theoretical guarantees. In deep learning, DP is typically achieved by perturbing the gradients
w.r.t. the model's loss on the training data (Abadi et al., 2016).

Data valuation (Ghorbani & Zou, 2019) has received increased interest from providing attribution for 034 parties in collaborative ML (Sim et al., 2022) and identifying data quality for data curation (Ghorbani & Zou, 2019) and data marketplace (Agarwal et al., 2019). Data values are often estimated with sampling-based methods such as Monte Carlo (Castro et al., 2009) on user statistics, e.g. gradients 037 of user data (Ghorbani & Zou, 2019). However, the sensitive nature of the user statistics requires privacy during data valuation to protect the data of participants in collaboration (Sim et al., 2023) or to facilitate interaction between data buyers and sellers in a marketplace (Chen et al., 2023). 040 Specifically, enforcing DP on the data is desirable in these scenarios for two reasons: 1) DP enjoys 041 the post-processing immunity (Dwork & Roth, 2014), allowing further access to data without risking 042 privacy leakage; 2) DP offers a natural trade-off between privacy protection and data quality: Data 043 contributors can determine at their discretion the level of information protection at the expense of 044 degraded data quality. While some previous works considered DP in data valuation (Sim et al., 2023; Wang et al., 2023; Watson et al., 2022), they either focused on a limited class of ML models or required a trusted central server, both of which are difficult to satisfy in real-world scenarios (Sim 046 et al., 2022). A natural question arises: Can we circumvent these two challenges in data valuation 047 while enforcing DP? 048

One might think of perturbing gradients on user data before using them for updating a gradienttrained parametric model to overcome the challenges. Unfortunately, we demonstrate that the naive
approach of adding i.i.d. Gaussian noise (Dwork & Roth, 2014) to user gradients leads to practically
useless data value estimates: The perturbation on the gradients erodes the information they carry,
leading to a lowered data value and increased data value estimation uncertainty because the noise
introduced by such perturbation accumulates with repeated evaluations (i.e., sampling methods)

054 methods (Castro et al., 2009; Maleki et al., 2014)). While a lowered (mean of) data value estimate 055 is intuitive and expected as a result of information loss (Dwork et al., 2015; Kairouz et al., 2015), 056 increased estimation uncertainty can render existing data valuation methods useless because the 057 magnitude of the noise scales with the evaluation budget for a fixed DP guarantee. As an illustration, 058 the 3rd figure of Fig. 2 shows that as the number of evaluations k increases, removing data with high data value estimates produces a curve closer to that with random removal, implying that more evaluation samples, paradoxically, have worsened the quality of data value estimates. We theoretically 060 account for this counter-intuitive phenomenon by showing that perturbation introduced by DP can 061 cause the estimation uncertainty of data value estimates, measured in terms of the variance brought 062 by the perturbation, to scale linearly in the number of evaluations. 063

064 Fortunately, through our developed theoretical analysis, we derive an insight that leads to a method for controlling the estimation uncertainty. We focus the analysis on the family of semivalues widely 065 adopted as data valuation metrics such as data Shapley (Ghorbani & Zou, 2019), Beta Shapley (Kwon 066 & Zou, 2022), and data Banzhaf (Wang & Jia, 2023). We revisit the paradigm of data valuation in 067 the context of gradient-based DP ML and thus design a technique that perturbs the gradients with 068 correlated noise in repeated evaluations to mitigate the above issue. Contrary to the linear scaling of 069 estimation uncertainty of the data value estimates with the evaluation budget in the naive approach, our proposed method is shown to control the estimation uncertainty to a *constant*. Additionally, we 071 empirically demonstrate that, on various ML tasks, our proposed method produces (i) greater model 072 degradation from removing high-value data; and (ii) higher AUC scores in identifying label-corrupted 073 data. In comparison, the naive approach performs similarly to random selection. We also apply our 074 approach to other scenarios including dataset valuation (Wu et al., 2022) and FL (McMahan et al., 075 2017). Our specific contributions are summarized as follows:

- We formalize a notion of *estimation uncertainty* (Eq. (3)) to specifically target the uncertainty due to DP. We then theoretically identify that under the naive approach of injecting i.i.d. noise for DP, the estimation uncertainty grows in $\Omega(k)$ with the evaluation budget k (Prop. 5.1), resulting in low-quality data value estimates.
 - As mitigation, we propose a simple yet effective approach (in Algorithm 1) via injecting correlated noise to control the estimation uncertainty to $\mathcal{O}(1)$ (Prop. 5.4) as opposed to $\Omega(k)$.
- We empirically demonstrate the implications of the escalating estimation uncertainty shown by the near-random performance on the data removal task (Sec. 6.1) and noisy label detection task (Sec. 6.2). Our approach outperforms the baseline approach on these tasks (Sec. 6.2) and on noisy label detection task in dataset valuation and FL (Sec. 6.3).

2 RELATED WORK

076

077

079

080

081

082

083

084

085

087

090

Prior works (Sim et al., 2023; Wang et al., 2023; Watson et al., 2022; Usynin et al., 2024) that consider 091 privacy-aware data valuation have limited applicability due to limited settings. (Sim et al., 2023) 092 considered perturbing user statistics to ensure DP but is restricted on the class of Bayesian models whereas we consider a wider family of models trained with gradient-based methods including neural 094 networks. (Wang et al., 2023) proposed a private variant of KNN-Shapley but did not generalize 095 to other semivalues, whereas our method applies to all semivalues. Watson et al. (2022) directly 096 perturbed the semivalue estimates. However, both (Wang et al., 2023; Watson et al., 2022) require a 097 trusted server to centralize the original gradients which may not reflect real-world scenarios where 098 untrusted central servers pose added privacy risks, whereas our approach does not require a trusted central server. (Usynin et al., 2024) assessed using variance of gradients and privacy loss-input 099 susceptibility score to select useful data points for DP training. The authors further propose methods 100 to compute the DP version of these scores. (Bani-Harouni et al., 2023) considered improving the 101 performance of DP-SGD by utilizing cosine similarity between privatized per-sample gradient and 102 original gradient to decide whether to include the gradient in averaged gradient. 103

Li & Yu (2023); Wang & Jia (2023) identified that stochastic utility functions can lead to noisy data
values which deteriorated the rank preservation and proposed to use (weighted) Banzhaf values as the
semivalue metric. Although the noise introduced by DP also renders the utility functions stochastic,
their methods do not consider the DP setting where noise scales with the number of evaluation
budgets. We specifically consider mitigating the issue of the scaling noise.

Dwork et al. (2010); Li et al. (2015); Nasr et al. (2020) introduced correlated noise to mitigate the effect of noise brought by DP, which has since been widely adopted in online learning with DP requirement to improve learning performance (Choquette-Choo et al., 2023a;b; Denisov et al., 2023; Kairouz et al., 2021). In contrast, we apply correlated noise to the setting of data valuation with privacy needs, to improve the quality of data value estimates. (Koloskova et al., 2024) analyzed the use of correlated noise under a DP follow-the-regularized-leader setting.

3 PRELIMINARIES

114 115

116 117

118

119 120

121

125

134 135 136

137

We recall the definition of semivalue for data valuation (Ghorbani & Zou, 2019) and the necessary preliminaries on DP.

3.1 DATA VALUATION

122 Semivalues. Denote $[n] := \{1, 2, ..., n\}$. The *semivalue* of i in a set [n] of parties w.r.t. a utility 123 function $V : 2^{[n]} \to \mathbb{R}$ and a weight function $w : [n] \to \mathbb{R}$ s.t. $\sum_{r=1}^{n} {n-1 \choose r-1} w(r) = n$ is (Dubey 124 et al., 1981)

$$\phi_i \coloneqq \sum_{r=1}^n n^{-1} w(r) \sum_{S \subseteq N \setminus \{i\}, |S| = r-1} [V(S \cup \{i\}) - V(S)] .$$
(1)

Leave-one-out (Cook, 1977), Shapley value (Shapley, 1953), and Banzhaf value (Wang & Jia, 2023) are examples of semivalues. In data valuation, a party can be represented by a data point, a dataset, or (the data of) an agent in FL setting. In ML, semivalues are often treated as a random variable and estimated using Monte Carlo methods (Castro et al., 2009; Maleki et al., 2014) since [n] is usually large and V is stochastic (more details in App. B.1). Denote P_j^{π} as the set of predecessors of party j in a permutation π uniformly randomly drawn from the set of all permutations Π , and let $p_j(\pi) \coloneqq 2^{n-1}n^{-1}w(|P_j^{\pi} \cup \{j\}|)$, then $\phi_j = \mathbb{E}[\psi_j]$ with ψ_j an average over k random draws:

$$\psi_j = (1/k) \sum_{t=1}^k p_j(\pi^t) [V(P_j^{\pi^t} \cup \{j\}) - V(P_j^{\pi^t})].$$
⁽²⁾

3.2 DIFFERENTIALLY PRIVATE MACHINE LEARNING (DP ML)

Definition 3.1 ((ϵ, δ)-Differential Privacy (Dwork & Roth, 2014, Def. 2.4)). A randomized algorithm *M* with domain \mathcal{D} and range \mathcal{R} is said to be (ϵ, δ)-differentially private if for any two neighboring¹ datasets $d, d' \in \mathcal{D}$, and for all event $S \subseteq \mathcal{R}$, $\Pr(\mathcal{M}(d) \in S) \leq \exp(\epsilon)\Pr(\mathcal{M}(d') \in S) + \delta$.

141 Importantly, the DP guarantee of \mathcal{M} is immune against post-processing (Dwork & Roth, 2014, 142 Proposition 2.1): The composition $f \circ \mathcal{M}$ with an arbitrary randomized mapping f have the same DP 143 guarantee as \mathcal{M} . We adopt this definition of DP to show the linearly scaling effect of perturbation 144 (Sec. 5.1). We elaborate in App. B.2 that our analysis can be extended to other DP frameworks.

145 146 147

4 SETTINGS AND PROBLEM STATEMENT

148 Settings. Our analysis is based on the G-Shapley framework (Ghorbani & Zou, 2019) for gradient-149 based ML methods where a parametric ML model learns from the data of each party via the perturbed 150 gradients of the data against a deterministic loss function $\mathcal{L}: [n] \times \mathbb{R}^d \to \mathbb{R}$ which maps the (data 151 of a) party and model parameters ($\in \mathbb{R}^d$) to a score ($\in \mathbb{R}$). The utility improvement reflects the 152 data value after the model updates its parameters with the gradient. For an evaluation budget k, k153 uniformly random permutations $\pi^1, \ldots, \pi^k \in \Pi$ are sampled, and for each sampled permutation 154 π (superscript omitted), a model is randomly initialized with θ_{π} and updated by the parties in sequence according to π . Then, denote $\theta_{\pi,i}^p$ the model parameters immediately before party j 155 updates the model in permutation π . For each subsequent party j in π , the Gaussian mechanism 156 is adopted to obtain a perturbed gradient $\tilde{g}_{\pi,j}$ to update the model $\theta_{\pi,j} \coloneqq \theta_{\pi,j}^p - \alpha \tilde{g}_{\pi,j}$ where 157 $\tilde{g}_{\pi,j} \coloneqq \hat{g}_{\pi,j} + z$ from a Gaussian noise $z \sim \mathcal{N}(\mathbf{0}, (C\sigma)^2 \mathbf{I})$ (with $C, \sigma > 0$) and the norm clipped gradient $\hat{g}_{\pi,j} \coloneqq g_{\pi,j} / \max(1, \|g_{\pi,j}\|_2 / C)$ based on the gradient $g_{\pi,j} \coloneqq \nabla_{\boldsymbol{\theta}} \mathcal{L}(j, \boldsymbol{\theta}_{\pi,j}^p)$ derived from 158 159 j's data. For a fixed test dataset, a utility $V(P_i^{\pi} \cup \{j\})$ representing the test performance depends on 160

¹⁶¹

¹Two inputs x, x' are neighboring if they differ by one training example (Abadi et al., 2016).

162 the model parameter $\theta_{\pi,j}$ (e.g. V is the negated test loss), so we replace $V(P_j^{\pi} \cup \{j\})$ with $V(\theta_{\pi,j})$ hereafter to highlight the interaction between model parameters and utility. While it is possible 163 164 $P_i^{\pi} \cup \{j\}$ results in different $\theta_{\pi,j}$ due to the random model initialization and varying orders of parties 165 in π , our subsequent definition of estimation uncertainty (Eq. (3)) carefully excludes their effects and focuses on the uncertainty due to DP if V is *deterministic* w.r.t. model parameters θ as we fix $\theta_{\pi,i}^p$. 166 167

Algorithm 1 outlines this i.i.d. noise approach, our method and a variant, explained in Sec. 5.

Algorithm 1 i.i.d. Corr. Noise (X) Corr. Noise (Y)

168 169

170 171

172

173

197

199

200 201

211

1: **Input:** number of parties n, utility function V, clipping norm C, loss function \mathcal{L} , noise multiplier σ , number of evaluations k, weight coefficient p_i , burn-in ratio q. 2: for $t \leftarrow 1$ to k do

174 Draw $\pi^t \stackrel{\text{unif.}}{\sim} \Pi$ and initialize the model with θ_{π^t} ; 3: 175 $i \leftarrow \pi^t[0]; \boldsymbol{\theta}_{\pi^t,i} \leftarrow \boldsymbol{\theta}_{\pi^t}$ 4: 176 for $j \in \{\pi^t[1], \pi^t[2], \dots, \pi^t[n]\}$ do $g_{\pi^t, j} \leftarrow \nabla_{\boldsymbol{\theta}} \mathcal{L}(j, \boldsymbol{\theta}_{\pi^t, j}^p)$ 5: 177 6: 178 $\tilde{g}_{\pi^t,j} \leftarrow g_{\pi^t,j} / \max(1.0, \|g_{\pi^t,j}\|_2/C) + \mathcal{N}(\mathbf{0}, k(C\sigma)^2 \mathbf{I})$ if Correlated Noise (\mathbf{X}) or Correlated Noise (\mathbf{Y}) then 7: 179 8: $\tilde{g}_{\pi^t,j} \leftarrow (1 - \boldsymbol{X}_{t,t}) \times \tilde{g}_j^{\text{roll}} + \boldsymbol{X}_{t,t} \times \tilde{g}_{\pi^t,j}$ 9: 181 $\tilde{g}_{j}^{\text{roll}} \leftarrow (t-1)/t \times \tilde{g}_{j}^{\text{roll}} + 1/t \times \tilde{g}_{\pi^{t},j}$ 182 10: 183 11: end if $\boldsymbol{\theta}_{\pi^t,j} \leftarrow \boldsymbol{\theta}_{\pi^t,j}^p - \alpha \tilde{g}_{\pi^t,j}$ 12: 185 if Correlated Noise (Y) then 13: 186 if t > kq then 14: 187 $\psi_j \leftarrow \frac{t-kq-1}{t-kq}\psi_j + \frac{p_j(\pi^t)}{t-kq}(V(\boldsymbol{\theta}_{\pi^t,j}) - V(\boldsymbol{\theta}_{\pi^t,i}))$ 188 15: 189 end if 16: 190 $i \leftarrow j$ and **Continue** 17: 191 18: end if 192 $\psi_j \leftarrow \frac{t-1}{t}\psi_j + p_j(\pi^t)(V(\boldsymbol{\theta}_{\pi^t,j}) - V(\boldsymbol{\theta}_{\pi^t,i}))/t ; i \leftarrow j$ 19: 193 20: end for 194 21: end for 22: **Return** ψ 196

Problem Statement. We study how noise introduced by DP impacts the estimation uncertainty of semivalue estimates. Formally, the estimation uncertainty for each party j is defined as

$$\operatorname{Var}[\psi_j | \boldsymbol{\theta}_{\pi^1,j}^p, \boldsymbol{\theta}_{\pi^2,j}^p, \dots, \boldsymbol{\theta}_{\pi^k,j}^p], \qquad (3)$$

202 and we aim to precisely understand its (asymptotic) relationship with k under i.i.d. noise and using our 203 proposed method (in Sec. 5) respectively. While Eq. (3) may seem more complicated than $Var[\psi_i]$, 204 the conditioning effectively removes the stochasticity in permutation sampling, gradient descent, and 205 model initialization, thus focusing on the impact of noise introduced by DP. Indeed, when there is no 206 perturbation due to DP and V is deterministic (w.r.t. $\boldsymbol{\theta}$), $\operatorname{Var}[\psi_j | \boldsymbol{\theta}_{\pi^1,j}^p, \boldsymbol{\theta}_{\pi^2,j}^p, \dots, \boldsymbol{\theta}_{\pi^k,j}^p] = 0$, meaning 207 Eq. (3) solely depends on the randomness from the perturbations of DP. Importantly, the definition of 208 estimation uncertainty in Eq. (3) admits precise and intuitive results about the randomness due to DP, 209 enabling us to pinpoint an issue with the i.i.d. noise approach; then, we propose the correlated noise 210 approach to asymptotically reduce this uncertainty and mitigate the issue.

212 We note that minimizing the estimation uncertainty *does not* necessarily recover the Remark. 213 original (i.e., without DP) data value estimate since DP exerts an upper bound on the marginal contribution (more discussed in App. B.4). Nevertheless, we highlight that this limitation does not 214 diminish the merit of our analysis as reducing estimation uncertainty helps improve the preservation 215 of the ranking of data value estimates of different parties compared to the naive approach of injecting

216 i.i.d. noise, which is an important axiomatic characterization and a common application of such data 217 value metrics (Ghorbani & Zou, 2019; Zhou et al., 2023). 218

Our subsequent analysis is w.r.t. a particular party j and for notational brevity, omits the subscript j(e.g., $\theta_{\pi,j} = \theta_{\pi}, \tilde{g}_{\pi,j} = \tilde{g}_{\pi}$) where the context is clear. Table 4 in App. A consolidates the important notations used throughout this work.

APPROACH AND THEORETICAL RESULTS 5

We summarize the results on estimation uncertainty for the naive approach (i.i.d. noise), our method, and a variant in Table 1. All proofs and derivations are deferred to App. C.

0.70

0.65

40

-60

Table 1: Summary of theoretical results. σ_a^2 refers to the average variance of unperturbed

219

220

221 222

223 224

225

226 227

228

229

230

231

232

233

234

235

236

237 238 239

240

241

242

243 244

246

267



Choice of V. For the purpose of mathematical analysis, we consider, in our theoretical analysis, the case where V is deterministic, non-positive, and Lipschitz smooth. An example of such V is the (average) negated loss on a fixed test dataset (Ghorbani & Zou, 2019), elaborated in Sec. 5.1. Nevertheless, Sec. 6.1 additionally empirically investigates test accuracy as V, which is discrete, to show our method works even if these properties are relaxed.

245 5.1 I.I.D. NOISE CAUSES SCALING ESTIMATION UNCERTAINTY

We first reveal that the propagation of noise from each z_t to V can be catastrophic to the estimator ψ : 247 a higher evaluation budget k, surprisingly, leads to a higher estimation uncertainty. 248

249 The estimator ψ in Eq. (2) aggregates various marginal contributions $m(\pi) = V(\theta_{\pi}) - V(\theta_{\pi}^{p})$ over 250 different $\pi \in \Pi$. Thus, if ψ requires k evaluations of m, party j needs to reveal its gradients k times. The repeated release of gradients increases the privacy risk, requiring greater perturbation to maintain 251 the same DP guarantee. In the Gaussian mechanism, the best known variance of the Gaussian noise 252 grows in $\Theta(k)$ (Abadi et al., 2016, Theorem 1) (more discussed in App. B.3), so each Gaussian noise 253 is expressed as $z_t \sim \mathcal{N}(\mathbf{0}, k(C\sigma)^2 \mathbf{I})$ where both C and σ^2 are constants satisfying a fixed (ϵ, δ) -DP 254 guarantee hereafter. 255

256 The linear scaling of $Var[z_t]$ creates a vicious cycle: More marginal contributions are needed to obtain a more certain semivalue estimation, incurring a larger k, which in turn necessitates greater perturba-257 tions for DP and increasing the estimation uncertainty. Formally, we show that the perturbations used 258 to guarantee DP are amplified by a strongly concave utility function V (achievable with a strongly 259 convex loss w.r.t. model parameters, such as K-Means, Lasso, and logistic regression with weight 260 decay), and causes the estimation uncertainty of semivalues to grow in the order of $\Omega(k)$, via two 261 specific globally strongly convex loss functions: 262

Proposition 5.1. (I.I.D. Noise) $\forall t \in [k]$, denote $\boldsymbol{\theta}_{\pi^t} \coloneqq \boldsymbol{\theta}_{\pi^t}^p - \alpha \tilde{g}_{\pi^t} = \boldsymbol{\theta}_{\pi^t}^p - \alpha (\hat{g}_{\pi^t} + z_t)$ where 263 $\forall t \in [k], z_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, k(C\sigma)^2 \mathbf{I})$ under the Gaussian mechanism. Denote a test dataset $\mathcal{D}_{\text{test}} =$ 264 $\{(x_1, y_1), \dots, (x_l, y_l)\}$ of l data points. If V is the negated mean-squared error loss on a linear 265 regression model 266

$$V(\boldsymbol{\theta}) \coloneqq -l^{-1} \sum_{i=1}^{l} (\boldsymbol{\theta}^{\top} \boldsymbol{x}_i - y_i)^2$$

268 or the negated ℓ_2 -regularized cross-entropy loss on a logistic regression model 269

$$V(\boldsymbol{\theta}) \coloneqq l^{-1} \sum_{i=1}^{l} (1 - y_i) \log(1 - \operatorname{Sig}(\boldsymbol{\theta}, \boldsymbol{x}_i)) + l^{-1} \sum_{i=1}^{l} y_i \log(\operatorname{Sig}(\boldsymbol{\theta}, \boldsymbol{x}_i)) - \lambda \|\boldsymbol{\theta}\|_2^2$$

where Sig $(\theta, x_i) = (1 + e^{-\theta^\top x_i})^{-1}$ is the sigmoid function and $\lambda > 0$ is the regularization hyperparameter. Denote $m(\pi^t) \coloneqq V(\theta_{\pi^t}) - V(\theta_{\pi^t}^p)$. Further denote a regular semivalue estimator $\psi \coloneqq k^{-1} \sum_{t=1}^k p(\pi^t) m(\pi^t)$.² Then, Var $[\psi | \theta_{\pi^1}^p, \theta_{\pi^2}^p, \dots, \theta_{\pi^k}^p] = \Omega(k)$.

273

288

307

313

318

321 322

323

As an intuition, consider that in each iteration, a Gaussian noise of variance $k(C\sigma)^2$ is injected to each gradient and amplified when propagated to the loss (i.e., -V by definition) due to strong convexity to become $\Omega(k^2)$, and the averaging over multiple evaluations shaves off a factor of k, yielding the final $\Omega(k)$ estimation uncertainty. Worse, the estimation uncertainty grows asymptotically with k: increasing the evaluation budget increases the estimation uncertainty due to the noise of DP.

Convex loss. As the loss surface of neural networks can exhibit strong convexity around local minima (Kleinberg et al., 2018; Milne, 2019), we derive, with a few additional assumptions including a notion of local strong convexity, a counter-part result to Prop. 5.1 w.r.t. deep neural networks where the estimation uncertainty is also $\Omega(k)$ via Prop. C.4 in App. C. In short, using i.i.d. noise can lead to scaling estimation uncertainty in many cases, raising a significant issue for data valuation, which we mitigate next.

287 5.2 CORRELATED NOISE TO REDUCE ESTIMATION UNCERTAINTY

The key to mitigating this "linearly scaling" estimation uncertainty lies in reducing the variance of the noise z_t in each iteration to render it less significant after amplification when propagated to the loss (equivalently V). Taking advantage of the way data values are estimated where private gradients are continuously released in each iteration, we can add a carefully correlated noise z_t^* to the gradients \hat{g}_{π^t} instead of independent noise $z_t \sim \mathcal{N}(\mathbf{0}, k(C\sigma)^2 \mathbf{I})$ while achieving the same DP guarantee.

Constructing z_t^* exploits the post-processing property of DP: reusing previously released private statistics does not affect the DP guarantee level. At each iteration t, instead of directly injecting 295 i.i.d. noise to \hat{g}_{π^t} to become $\tilde{g}_{\pi^t} := \hat{g}_{\pi^t} + z_t$, estimation uncertainty can be reduced by reusing 296 the i.i.d.-perturbed gradients $\tilde{g}_{\pi^1}, \ldots, \tilde{g}_{\pi^t}$. To ease the understanding of the core idea, we begin by 297 assuming that the original gradients at each iteration are identical, i.e. $\hat{g}_{\pi^1} = \hat{g}_{\pi^2} = \ldots = \hat{g}_{\pi^k}$, 298 and relax it later. Instead of \tilde{g}_{π^t} , a weighted sum (e.g., $\tilde{g}^*_{\pi^t} \coloneqq t^{-1} \sum_{l=1}^t \tilde{g}_{\pi^l}$) of previously released 299 private gradients can be utilized. This implicitly constructs the correlated noise $z_t^* \coloneqq t^{-1} \sum_{l=1}^t z_l \sim z_l$ 300 $\mathcal{N}(\mathbf{0}, (k/t)(C\sigma)^2 \mathbf{I})$, resulting a reduction in variance by a factor of t. To see how this variance 301 reduction is propagated to V, compare the variance on the following strongly convex function (to mimic V): $\operatorname{Var}[\|\tilde{g}_{\pi^t}\|_2^2] \to t^{-2}\operatorname{Var}[\|\tilde{g}_{\pi^t}\|_2^2]$ as $k \to \infty$ (see Obs. C.1). Notice that the variance 302 303 reduction is amplified from t^{-1} to t^{-2} . Indeed, such variance reduction can lead to an $\mathcal{O}(\log^2 k)$ 304 bound for the estimation uncertainty (see Prop. C.9). However, two questions remain: (i) Is the 305 "identical gradient" assumption satisfied for semivalue estimation, and what if it is not? (ii) How to 306 cleverly reuse previously privatized gradients to obtain a lower estimation uncertainty?

For question (i), unfortunately, the assumption is *not* satisfied: The gradients at each permutation \hat{g}_{π^t} , though obtained from the same underlying data, are not identical in general since $\theta_{\pi^t}^p$ are different for different π^t . Specifically, consider a party j whose unperturbed gradients in k evaluations are $\hat{g}_{\pi^1}, \hat{g}_{\pi^2}, \ldots, \hat{g}_{\pi^k}$. The unperturbed gradients are first injected with a Gaussian noise to produce the perturbed gradients $\tilde{g}_{\pi^t} = \hat{g}_{\pi^t} + z_t$ for $t \in [k]$. Then, party j will release the gradient

$$\tilde{g}_{\pi^t}^* \coloneqq t^{-1} \sum_{l=1}^t \tilde{g}_{\pi^l} = t^{-1} \sum_{l=1}^t \hat{g}_{\pi^l} + t^{-1} \sum_{l=1}^t z_l$$

which, however, is not an unbiased estimator for \hat{g}_{π^t} since $\mathbb{E}[\tilde{g}_{\pi^t}] = t^{-1} \sum_{l=1}^t \hat{g}_{\pi^l} \neq \hat{g}_{\pi^t}$ unless $\hat{g}_{\pi^1} = \hat{g}_{\pi^2} = \ldots = \hat{g}_{\pi^t}$ (i.e., identical gradients). Nevertheless, empirical observations suggest that as compared to \tilde{g}_{π^t} , \hat{g}_{π^t} is much "closer" to $\tilde{g}_{\pi^t}^*$ in terms of the mean difference in cosine similarity

$$\Delta_{\cos} \coloneqq n^{-1} \sum_{j \in [n]} k^{-1} \sum_{t=1}^{k} [\cos(\hat{g}_{\pi^{t},j}, \tilde{g}_{\pi^{t},j}^{*}) - \cos(\hat{g}_{\pi^{t},j}, \tilde{g}_{\pi^{t},j})]$$

where $\cos(a, b) \coloneqq |a \cdot b| / (||a||_2 ||b||_2)$, and the difference in ℓ_2 distance

$$\Delta_{\ell_2} \coloneqq n^{-1} \sum_{j \in [n]} k^{-1} \sum_{t=1}^k \|\hat{g}_{\pi^t,j} - \tilde{g}_{\pi^t,j}^*\|_2 - \|\hat{g}_{\pi^t,j} - \tilde{g}_{\pi^t,j}\|_2,$$

²A regular semivalue has $p(\pi) > 0$ for all $\pi \in \Pi$ (Carreras & Freixas, 2002). Examples include Shapley value and Banzhaf value.

as shown in Fig. 1 where we empirically investigate the similarity with 400 randomly selected data points from the diabetes dataset (Efron et al., 2004) trained with logistic regression. The difference is even more pronounced with a higher budget k, suggesting that $\tilde{g}_{\pi t}^*$ is better than $\tilde{g}_{\pi t}$ in approximating $\hat{g}_{\pi t}$. Therefore, relaxing the assumption of identical gradients, we assume that they are i.i.d. samples of a distribution (also assumed in (Faghri et al., 2020; Zhang et al., 2013)). Additionally, we further assume an diagonal sub-Gaussian distribution, defined as follows.

330 Definition 5.2 (Diagonal Multivariate Sub-Gaussian Distribution). Let $X \in \mathbb{R}^n$ be a random **331** vector with $\mathbb{E}[X] = \mathbf{0}$. X is said to follow a diagonal multivariate sub-Gaussian distribution if **332** $\forall i \in [n], X_i$ follows a sub-Gaussian distribution and $\forall i, j \in [n]$ s.t. $i \neq j$, $\text{Cov}[X_i, X_j] = 0$.

For question (ii), we specifically consider linear combinations of the private gradients to use linearity of expectation to ensure unbiasedness (see P2). Formally, we express the finally released gradients in all iterations as a matrix product W = XA, where $A = (\tilde{g}_{\pi^1}, \ldots, \tilde{g}_{\pi^k})^{\top}$ and X is a square matrix (which will be relaxed in Sec. 5.3) mapping A to the matrix of released gradients W = $(\tilde{g}_{\pi^1}^*, \ldots, \tilde{g}_{\pi^k}^*)^{\top}$. Since there are many possible matrices W, we describe two key principles (specific to data valuation) in selecting W, for a square matrix X.

P1. Lower traingularity: X should be lower triangular.

P2. Unbiasedness: the combined gradient $\tilde{g}_{\pi t}^*$ should be a *weighted sum* of the previous gradients, i.e., $\tilde{g}_{\pi t}^* := \sum_{l=1}^t X_{t,l} \tilde{g}_{\pi l}$ where $\sum_{l=1}^t X_{t,l} = 1$.

P1 ensures that each revealed gradient is calculated as a weighted prefix sum of the preceding perturbed gradients as future gradients are unknown. P2 ensures unbiased gradient estimate $\mathbb{E}[\tilde{g}_{\pi^t}] = \mathbb{E}[\hat{g}_{\pi^t}]$ (note that we treat \hat{g}_{π^t} as a random variable here), and thus an unbiased estimate for the utility. The following square matrix X^* satisfies P1 and P2:

$$\mathbf{X}_{i,j}^* \coloneqq i^{-1} \cdot \mathbf{1}_{j \le i} \,. \tag{4}$$

349 Despite its simplicity, X^* is surprisingly a "one-size-fits-all" matrix as it delivers a lower asymptotic 350 bound than $\Omega(k)$ (with i.i.d. noise) when k is large, even if the unperturbed gradients are not identical: 351 **Proposition 5.3** (Correlated Noise with X, informal). Let $\tilde{g}_{\pi^{l}}, l \in [t]$ be perturbed gradients using 352 the Gaussian mechanism that satisfies (ϵ, δ) -DP. $\forall t \in [k]$, denote $\boldsymbol{\theta}_{\pi^t}^* \coloneqq \boldsymbol{\theta}_{\pi^t}^p - \alpha \sum_{l=1}^t \boldsymbol{X}_{t,l} \tilde{g}_{\pi^l}$, 353 $m^*(\pi^t) \coloneqq V(\boldsymbol{\theta}_{\pi^t}^*) - V(\boldsymbol{\theta}_{\pi^t}^p) \text{ and } \psi^* \coloneqq k^{-1} \sum_{t=1}^k p(\pi^t) m^*(\pi^t). \text{ Assume that } \forall t \in [k], \hat{g}_{\pi^t} - \mathbb{E}[\hat{g}_{\pi^t}]$ 354 i.i.d. follow an diagonal multivariate sub-Gaussian distribution with covariance $\Sigma \in \mathbb{R}^{(d \times d)}$ and let 355 $\sigma_g^2 \coloneqq d^{-1} \sum_{r=1}^d \Sigma_{r,r}$. Then, using a suitable matrix \boldsymbol{X} can produce $\operatorname{Var}[\psi^* | \boldsymbol{\theta}_{\pi^1}^p, \boldsymbol{\theta}_{\pi^2}^p, \dots, \boldsymbol{\theta}_{\pi^k}^p] =$ 356 $\mathcal{O}(\log^2 k + \sigma_a^4)$ and $\mathbb{E}[\psi^* - \psi] = \mathcal{O}(\log k + \sigma_a^2)$ while satisfying (ϵ, δ) -DP. Moreover, as $k \to \infty$, 357 358 $X o X^*$. 359

A formal statement of Prop. 5.3 and its proof is found in Appendix (Prop. C.10). Prop. 5.3 offers a better bound than $\Omega(k)$ in Prop. 5.1 as σ_g^2 is a constant. Notice that as $k \to \infty$, $X \to X^*$. An exciting implication of Prop. 5.3 is that one can use X^* to approximate X by setting a large kregardless of σ_g^2 and achieve the given asymptotic bound. Another implication is that implementation of Prop. 5.3 is simple since all $X_{t,l}$ are identical except for $X_{t,t}$. Thus X is fully specified by just k parameters: $X_{t,t}$ for $t \in [k]$. As shown in Algorithm 1 in red, on top of i.i.d., the model is updated using a weighted sum between $\tilde{g}_{\pi^{t-1}}^*$ and \tilde{g}_{π^t} at each iteration t.

367 368

333

340

341

342 343

344

345

346

347 348

5.3 MORE VARIANCE REDUCTION WITH NON-SQUARE MATRIX

369 The bound in Prop. 5.3, though asymptotically better than $\Omega(k)$, still grows in terms of k: estimates 370 still become worse even with more samples. Inspired by the "burn-in" technique (Neiswanger et al., 371 2014) widely employed in MCMC, we show that it is possible to achieve a constant bound via a 372 non-square matrix (i.e., relaxing P1). For an arbitrary square matrix X, consider its counterpart Y373 defined with a hyperparameter $q \in (0,1)$, $\boldsymbol{Y}_{(k-kq) \times k} \coloneqq \boldsymbol{X}(kq+1:k;1:k)$ where the bracket means taking a sub-matrix with the selected rows and columns. In other words, given a X, Y is a 374 375 sub-matrix of X with the (kq + 1)th row to the last row, and thus a counterpart to X and relaxes **P1.** The intuition is that in the first few iterations, t is small, causing $\tilde{g}_{\pi t}^{*}$ to still incur relatively 376 large variances even with correlated noise via X. As a remedy, Y effectively discards these highly 377 fluctuating marginal contributions to yield an asymptotically even lower total variance than X.

378 **Proposition 5.4** (Correlated Noise with Y, *informal*). Let $\tilde{g}_{\pi^l}, l \in [t]$ be perturbed gradients 379 using the Gaussian mechanism that satisfies (ϵ, δ) -DP. $\forall t \in \{kq + 1, \dots, k\}$, denote $\theta_{\pi^t}^* \coloneqq \theta_{\pi^t}^p - \theta_{\pi^t}^{p}$ $\alpha \sum_{l=1}^{t} \boldsymbol{Y}_{t-kq,l} \tilde{g}_{\pi^{l}}, m^{*}(\pi^{t}) \coloneqq V(\boldsymbol{\theta}_{\pi^{t}}^{*}) - V(\boldsymbol{\theta}_{\pi^{t}}^{p}), \text{ and } \psi^{*} \coloneqq (k-kq)^{-1} \sum_{t=kq+1}^{k} p(\pi^{t}) m^{*}(\pi^{t}) \text{ for } q \in (0, 1). \text{ With the same assumption and the suitable matrix } \boldsymbol{X} \text{ discussed in Prop. 5.3, setting } \forall t \in [0, 1]$ 380 381 382 $[k], \forall l \in [t], \boldsymbol{Y}_{t,l} = \boldsymbol{X}_{t,l} \text{ produces } \operatorname{Var}[\psi^* | \boldsymbol{\theta}_{\pi^1}^p, \boldsymbol{\theta}_{\pi^2}^p, \dots, \boldsymbol{\theta}_{\pi^k}^p] = \mathcal{O}\left((1-q)^{-2} \log^2 (1/q) + \sigma_g^4\right)$ 383 and $\mathbb{E}[\psi^* - \psi] = \mathcal{O}((1 - q)\log 1/q + \sigma_q^2)$ while satisfying (ϵ, δ) -DP. 384

A formal statement of Prop. 5.4 is in Appendix (Prop. C.11). Prop. 5.4 shows that we can control 386 the estimation uncertainty by a *constant* (i.e., entirely removing the effect of k) with a combination 387 of correlated noise and burn-in. Intuitively, the first term represents the injected noise controlled by q. In particular, since $kq \in \mathbb{N}^+$, we have $q \ge 1/k$ which implies $\log^2(1/q) \le \log^2 k$. As 389 such, when $k \to \infty$ and $q \to 1/k$, the bound is reduced to that in Prop. 5.3. As $q \to 1$, too many evaluation samples are "burnt", leaving insufficient samples to average out the noise as reflected 390 by the exploding first term where $(1-q)^{-2}\log^2(1/q) \to \infty$. We show in detail how q affects the 391 392 estimation uncertainty in Sec. 6.2. Pseudo-code is in Algorithm 1 (in green). On top of red and 393 blue, green lines exclude the first kq marginal contributions from being included in ψ .

EXPERIMENTS 6

394

396 397

410

411

We fix C = 1.0 and $(\epsilon = 1, \delta = 5 \times 10^{-5})$ -DP guarantee unless otherwise specified. All experiments 398 are repeated over 5 independent trials. We focus on classification tasks as they are more susceptible 399 to noise and defer regression task to App. D.3. We consider data selection and noisy label detection 400 tasks as standard evaluations of the effectiveness of a data value estimate (Ghorbani & Zou, 2019; 401 Kwon & Zou, 2022; Wang & Jia, 2023; Zhou et al., 2023). Exploiting Prop. 5.3, we set a large 402 $k \geq 200$ (except FL) and use X^* and Y^* for adding correlated noise. Additional experimental 403 settings and results are in App. D. 404

While our theoretical results have provided a (ϵ, δ) -DP guarantee level, in App. D.1, we verify the 405 privacy protection of our method by constructing a membership inference attack (MIA) following the 406 setting in (Wang et al., 2023), and demonstrate that our method can successfully defend against the 407 constructed MIA. For experiments in the main text, we focus on how our method improves the data 408 value estimation with privacy protection. 409

INCREASING k UNDER I.I.D. NOISE *Does Not* REDUCE THE ESTIMATION UNCERTAINTY 6.1

412 We empirically show the scaling estimation uncertainty as the evaluation budget k increases and its 413 implication. As a setup, we randomly choose 400 training examples from the diabetes dataset (Efron et al., 2004) with the remaining data points as the test dataset. We train a logistic regression 414 using the negated cross-entropy loss on the test dataset as the utility function V. To evaluate 415 the estimation uncertainty, we first examine the quality of data value estimates through mean-416 adjusted variance of ψ_j (Zhou et al., 2023), which is the ratio between the empirical variance $s_j^2 := k^{-1}(k-1)^{-1} \sum_{t=1}^k [m_j(\pi^t) - \mu_j]^2$ and the empirical mean $\mu_j := k^{-1} \sum_{t=1}^k m_j(\pi^t)$. We also examine how test accuracy changes when removing training examples with the highest ψ 's. The 417 418 419 leftmost figure of Fig. 2 shows that the mean-adjusted variance $s_{i.i.d.}^2/|\mu_{i.i.d.}|$ with i.i.d. noise increases 420 with k, indicating greater estimation uncertainty, whereas using correlated noise, $s_{\text{corr.}}^2/|\mu_{\text{corr.}}|$ not 421 only decreases but is also smaller by several magnitudes (in 10^5). Moreover, $\mu_{i.i.d.}$'s computed with 422 i.i.d. noise are increasingly negative as k increases, whereas $\mu_{\text{corr.}}$'s computed with correlated noise 423 stay positive, suggesting that the estimated ψ 's are less affected by noise. The 3rd figure of Fig. 2 424 shows that ψ 's computed with greater k produces higher test accuracy during removal (i.e., closer to random removal), verifying our identified paradox that ψ 's computed with more budget are poorer 426 estimates of data values.

427

428 6.2 CORRELATED NOISE IMPROVES THE QUALITY OF THE ESTIMATES

429

Following the same setup as Sec. 6.1, we compute the ψ 's using correlated noise (with q = 0.8) and 430 no injected noise due to DP respectively. The results are shown in the rightmost figure of Fig. 2. 431 Removing data points with high ψ 's computed with no injected noise produces low test accuracy



Figure 2: (a) $n^{-1} \sum_{j \in [n]} s_j^2 / |\mu_j|$ and (b) μ_j vs. k using i.i.d. noise and correlated noise. (c) error bars of test accuracy vs. ratio of data removed with *the highest* ψ 's using different k with i.i.d. noise and (d) also with correlated noise and no DP. "Random" means random removal.

close to that with correlated noise. In contrast, the curve of ψ 's computed with i.i.d. noise lies far above (i.e., close to random removal), suggesting that ψ 's computed with correlated noise are much more reflective of the true worth of data as compared to ψ 's computed with i.i.d. noise.

Ablation study on the influence of q. We study how much q affects the data value estimation in a noisy label detection setting on two datasets. We randomly perturb 30% of labels on a selection of 800 training examples from Covertype dataset (Blackard, 1998) (and MNIST dataset (LeCun et al., 1990) in App. D.4) respectively. The datasets are trained with logistic regression (LR) and a more complex convolutional network (CNN). Ideally, a good data value estimate should assign the lowest ψ 's to the perturbed training examples. To measure this, we plot the AUC-ROC curve (AUC) in Fig. 3 with q = 0 equivalent to X^* . With increased k, the AUC with our method *increases* especially in the large-q region while the AUC with i.i.d. noise *decreases*. We also observe that AUC generally increases with q when $q \leq 0.9$.

Adopting other utility functions. While our theoretical analysis is w.r.t. negated loss as the utility function V, we demonstrate similar results with test accuracy as V, shown in the right column of Fig. 3. Our method still outperforms i.i.d. noise, although the AUC is lower than when V is negated loss. We think this is because accuracy is a less fine-grained metric. Hence subtle changes in the model performance (often due to perturbation) in this task cannot be well reflected. Moreover, we observe that as $q \rightarrow 1$, AUC degrades, consistent with the result in Prop. 5.4.



Figure 3: Plots of AUC v.s. burn-in ratio $q \in [0, 1)$ (with q = 0 equivalent to X^*). V is (left and middle) negated test loss and (right) test accuracy. Lines represent mean and shades represent 1 standard deviation. Higher is better.

Experiments on other semivalues. We compare the effect of utilizing correlated noise with data Banzhaf (Wang & Jia, 2023) and Beta Shapley (Kwon & Zou, 2022) in Table 2. We consider the same setup as Sec. 6.2 with k = 1000. We compare the performance using X^* and Y^* with

Table 2: Mean (std. errors) of AUC on Covertype trained with LR (top) and MNIST trained
with CNN (bottom). The best score is highlighted.
Higher is better.

Table 3: Mean (std. errors) of $\Delta_{i.i.d.}$ and $\Delta_{corr.}$ for dataset valuation (top) and FL (bottom). Best scores are highlighted (lower is better). *B* is the dataset size.

semivalues	no DP	i.i.d. noise	X^*	$Y^{*}(q = 0.5)$	$Y^{*}(q = 0.9)$				
Shapley	0.905 (1.00e-03)	0.675 (6.00e-03)	0.735 (2.00e-02)	0.774 (8.00e-03)	0.788 (4.00e-03)	ML Task	$\{k, n, B\}$	$\Delta_{i.i.d.}$	$\Delta_{\text{corr.}}$ (Ours)
Banzhaf	0.896 (2.00e-03)	0.533 (1.70e-02)	0.725 (2.00e-02)	0.770 (8.00e-03)	0.777 (4.00e-03)	MNIST + CNN	{200, 800, 8}	0.290 (1.63e-02)	0.170 (3.73e-02)
Beta(4, 1)	0.882 (1.00e-03)	0.612 (1.40e-02)	0.721 (2.10e-02)	0.766 (8.00e-03)	0.777 (4.00e-03)	CIFAR10 + CNN	$\{1000, 100, 32\}$	0.178 (5.85e-02)	0.0303 (1.53e-02)
Beta(16, 1)	0.875 (0.00e+00)	0.557 (2.40e-02)	0.707 (2.00e-02)	0.757 (8.00e-03)	0.767 (5.00e-03)	CIFAR10 + ResNet18	$\{1000, 100, 32\}$	0.204 (5.33e-02)	0.119 (1.73e-02)
Shanley	0.988 (0.00)	0.600 (1.50e=02)	0.627 (2.30e=02)	0.803 (4.00e=03)	0.828 (6.00e=03)	CIFAR10 + ResNet34	$\{1000, 100, 32\}$	0.141 (8.51e-03)	0.0433 (7.10e-03)
Banzhaf	0.985 (1.00e-03)	0.532 (2.00e-02)	0.616 (2.30e-02)	0.808 (4.00e-03)	0.827 (7.00e-03)	MNIST + CNN	$\{50, 50, 32\}$	0.195 (9.40e-03)	0.0810 (8.06e-03)
Beta(4, 1)	0.991 (0.00)	0.569 (1.50e-02)	0.618 (2.10e-02)	0.789 (4.00e-03)	0.810 (7.00e-03)	CIFAR10 + CNN	$\{50, 50, 32\}$	0.167 (3.96e-02)	0.0415 (9.63e-03)
Beta(16, 1)	0.992 (0.00)	0.539 (1.00e-02)	0.615 (1.10e-02)	0.773 (5.00e-03)	0.793 (9.00e-03)				

495 496 497

498

499

500

501 502

504

 $q \in \{0.5, 0.9\}$. For all variants, the AUC is higher than that with i.i.d. noise. Particularly, q = 0.9 works the best for all 4 semivalues tested. In contrast, the AUCs with i.i.d. noise using data Banzhaf are ≈ 0.53 on both datasets, *close to randomness* (0.5). The results show that our method generalizes to various data valuation metrics. We also note a similar improvement for LOO (Cook, 1977) in App. D.4 despite it not being a regular semivalue.

6.3 APPLICATION TO OTHER USE CASES

We verify the effectiveness of using correlated noise for dataset valuation (Wu et al., 2022) and 505 collaborator attribution in federated learning (Wang et al., 2020) (refer to App. D.2 for setup). We 506 consider MNIST dataset and CIFAR10 (Krizhevsky et al., 2012) dataset, trained on CNN and fine-507 tuned on pretrained ResNet18/34 (He et al., 2015). We tabulate the difference in AUC between 508 using correlated noise with q = 0.9, denoted as $\Delta_{corr.} \coloneqq AUC_{no DP} - AUC_{corr.}$, and using i.i.d. noise, 509 denoted as $\Delta_{i.i.d.} \coloneqq AUC_{no DP} - AUC_{i.i.d.}$. For ResNet34, we use $\epsilon = 10$ as the model is more 510 complex, causing both i.i.d. and our method to have degraded performance with strict privacy. Our 511 methods outperform i.i.d. noise as shown in Table 3 top. For FL, we notice that the continually 512 updating characteristic of the global model poses two challenges: (i) the overall scale of the loss 513 values decreases in each round as the model gradually converges such that marginal contributions 514 computed in later rounds are less significant, and (ii) the variance of gradients σ_a^2 is larger than in 515 common data valuation scenarios. To tackle these challenges, we adopt test accuracy as V so that Vare in the same scale in each round, use a matrix with $X_{t,t} > 1/t$ to control σ_q^2 , and choose a small 516 burn-in ratio q = 0.2 to keep more evaluations in the first few rounds (detailed in App. D.2). The 517 results with these modifications are shown in Table 3 bottom. Our method outperforms i.i.d. noise by 518 a large margin, showing the applicability of our approach beyond the default data valuation setting 519 under which our theoretical results are developed. 520

521 522

523

7 CONCLUSION, LIMITATIONS, AND FUTURE WORKS

In this work, we identify a problem in data valuation where DP is enforced via perturbing gradients 524 with i.i.d. noise: The estimation uncertainty scales linearly (i.e., $\Omega(k)$) with more budget k and 525 renders the data value estimates almost useless (i.e., close to random guesses in some investigated 526 cases). As a solution, we propose to use correlated noise and theoretically show that using a weighted 527 sum via matrix form using X provably reduces the estimation uncertainty of semivalues from $\Omega(k)$ 528 to $\mathcal{O}(1)$ and empirically demonstrate the implications of our method on various ML tasks and data 529 valuation metrics. One limitation is the need to store the gradients. However, this limitation is 530 alleviated when the number of parties is small (e.g. dataset valuation) or when the memory load 531 can be distributed across parties (e.g. FL). Another limitation is that our theoretical result assumes 532 an diagonal multivariate sub-Gaussian distribution of the gradients. Nevertheless, we empirically 533 demonstrate that our method works for neural networks where the assumption is not explicitly satisfied. A future direction is to explore other possible X to reduce estimation uncertainty further. 534

535

536 REFERENCES

 Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and
 Li Zhang. Deep learning with differential privacy. In *Proceedings of the ACM SIGSAC Conference* on Computer and Communications Security, 2016.

540 541	A. Agarwal, M. Dahleh, and T. Sarkar. A marketplace for data: An algorithmic solution. In <i>Proceedings of the ACM Conference on Economics and Computation</i> , 2019.
543 544 545	Shahab Asoodeh, Jiachun Liao, Flavio P. Calmon, Oliver Kosut, and Lalitha Sankar. Three variants of differential privacy: Lossless conversion and applications. <i>IEEE Journal on Selected Areas in Information Theory</i> , 2021.
546 547 548	Siva Balakrishnan. Lecture 7: Moments revisited and the central limit theorem. https://www.stat.cmu.edu/~siva/teaching/700/lec7.pdf, 2016. Accessed: 2024-11-27.
549 550 551 552	David Bani-Harouni, Tamara T. Mueller, Daniel Rueckert, and Georgios Kaissis. Gradient self- alignment private deep learning. In <i>Medical Image Computing and Computer Assisted Intervention</i> – <i>MICCAI 2023 Workshops: ISIC 2023, Care-AI 2023, MedAGI 2023, DeCaF 2023, Held in</i> <i>Conjunction with MICCAI 2023, Vancouver, BC, Canada, October 8–12, 2023, Proceedings, 2023.</i>
553 554	Jock Blackard. Covertype. UCI Machine Learning Repository, 1998.
555	Preston Bukaty. The california consumer privacy act, 2019.
556 557 558	Mark Bun and Thomas Steinke. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In <i>Proceedings of the International Conference on Theory of Cryptography</i> , 2016.
559 560	Francesc Carreras and Josep Freixas. Semivalue versatility and applications. Annals of Operations Research 109, 343–358 (2002), 2002.
561 562 563	Javier Castro, Daniel Gómez, and Juan Tejada. Polynomial calculation of the shapley value based on sampling. <i>Computers & Operations Research</i> , 36(5):1726–1730, 2009.
564 565 566	Lingjiao Chen, Bilge Acun, Newsha Ardalani, Yifan Sun, Feiyang Kang, Hanrui Lyu, Yongchan Kwon, Ruoxi Jia, Carole-Jean Wu, Matei Zaharia, and James Zou. Data acquisition: A new frontier in data-centric ai, 2023.
567 568 569 570	Christopher A. Choquette-Choo, Arun Ganesh, Ryan McKenna, H. Brendan McMahan, Keith Rush, Abhradeep Thakurta, and Zheng Xu. (amplified) banded matrix factorization: A unified approach to private training, 2023a.
571 572 573	Christopher A. Choquette-Choo, Hugh Brendan Mcmahan, J Keith Rush, and Abhradeep Guha Thakurta. Multi-epoch matrix factorization mechanisms for private machine learning. In <i>Proceedings of the International Conference on Machine Learning</i> , 2023b.
574 575	R. Dennis Cook. Detection of influential observation in linear regression. <i>Technometrics</i> , 19(1): 15–18, 1977.
576 577 578	Paulo Cortez, A. Cerdeira, F. Almeida, T. Matos, and J. Reis. Wine Quality. UCI Machine Learning Repository, 2009.
579 580	Council of European Union. Council regulation (EU) no 269/2014, 2014.
581 582 583	Sergey Denisov, Brendan McMahan, Keith Rush, Adam Smith, and Abhradeep Guha Thakurta. Improved differential privacy for sgd via optimal private linear operators on adaptive streams. In <i>Advances in Neural Information Processing Systems</i> , 2023.
584 585 586	Pradeep Dubey, Abraham Neyman, and Robert James Weber. Value theory without efficiency. <i>Mathematics of Operations Research</i> , 6(1):122–128, 1981.
587 588	Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. <i>Foundations and Trends in Theoretical Computer Science</i> , aug 2014.
589 590 591	Cynthia Dwork, Moni Naor, Toniann Pitassi, and Guy N. Rothblum. Differential privacy under continual observation. In <i>Proceedings of the ACM Symposium on Theory of Computing</i> , 2010.
592 593	Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Leon Roth. Preserving statistical validity in adaptive data analysis. In <i>Proceedings of the Annual ACM</i> <i>Symposium on Theory of Computing</i> , 2015.

594 595 596	Bradley Efron, Trevor Hastie, Iain Johnstone, and Robert Tibshirani. Least angle regression. <i>The Annals of Statistics</i> , 32(2), 2004.
597 598	Fartash Faghri, David Duvenaud, David J. Fleet, and Jimmy Ba. A study of gradient variance in deep learning, 2020.
599 600 601	Amirata Ghorbani and James Zou. Data Shapley: Equitable valuation of data for machine learning. In <i>Proceedings of the International Conference on Machine Learning</i> , 2019.
602 603	Sivakanth Gopi, Yin Tat Lee, and Lukas Wutschitz. Numerical composition of differential privacy. In <i>Advances in Neural Information Processing Systems</i> , 2021.
604 605 606	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. <i>arXiv preprint arXiv:1512.03385</i> , 2015.
607 608	Peter Kairouz, Sewoong Oh, and Pramod Viswanath. The composition theorem for differential privacy, 2015.
609 610 611 612	Peter Kairouz, Brendan Mcmahan, Shuang Song, Om Thakkar, Abhradeep Thakurta, and Zheng Xu. Practical and private (deep) learning without sampling or shuffling. In <i>Proceedings of the International Conference on Machine Learning</i> , 2021.
613 614	Bobby Kleinberg, Yuanzhi Li, and Yang Yuan. An alternative view: When does SGD escape local minima? In <i>Proceedings of the International Conference on Machine Learning</i> , 2018.
615 616 617 618	Anastasia Koloskova, Ryan McKenna, Zachary Charles, Keith Rush, and Brendan McMahan. Gra- dient descent with linearly correlated noise: theory and applications to differential privacy. In <i>Advances in Neural Information Processing Systems</i> , 2024.
619 620	Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convolu- tional neural networks. In <i>Advances in Neural Information Processing Systems</i> , 2012.
621 622 623	Yongchan Kwon and James Zou. Beta shapley: a unified and noise-reduced data valuation framework for machine learning. In <i>Proceedings of the International Conference on Artificial Intelligence and Statistics</i> , 2022.
625 626 627	Yann LeCun, Bernhard Boser, John Denker, Donnie Henderson, R. Howard, Wayne Hubbard, and Lawrence Jackel. Handwritten digit recognition with a back-propagation network. In <i>Advances in Neural Information Processing Systems</i> , 1990.
628 629 630	Chao Li, Gerome Miklau, Michael Hay, Andrew Mcgregor, and Vibhor Rastogi. The matrix mechanism: Optimizing linear counting queries under differential privacy. <i>The VLDB Journal</i> , dec 2015.
631 632 633	Weida Li and Yaoliang Yu. Robust data valuation with weighted banzhaf values. In Advances in Neural Information Processing Systems, 2023.
634 635	Sasan Maleki, Long Tran-Thanh, Greg Hines, Talal Rahwan, and Alex Rogers. Bounding the estimation error of sampling-based shapley value approximation, 2014.
636 637 638 639	H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In <i>Proceedings of the</i> <i>International Conference on Artificial Intelligence and Statistics</i> , 2017.
640 641	Tristan Milne. Piecewise strong convexity of neural networks. In Advances in Neural Information Processing Systems, 2019.
643	Ilya Mironov. Rényi differential privacy. In IEEE Computer Security Foundations Symposium, 2017.
644 645	Milad Nasr, Reza Shokri, and Amir houmansadr. Improving deep learning with differential privacy using gradient encoding and denoising, 2020.
647	Willie Neiswanger, Chong Wang, and Eric Xing. Asymptotically exact, embarrassingly parallel mcmc, 2014.

Opacus. Opacus PyTorch library. Available from opacus.ai, 2021.

648

- 649 650 Ahmed El Ouadrhiri and Ahmed Abdelhadi. Differential privacy for deep and federated learning: A survey. IEEE Access, 2022. 651 652 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor 653 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward 654 Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, 655 Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning 656 library, 2019. 657 Philippe Rigollet. Chapter 1: Sub-gaussian random variables. https://ocw.mit. 658 edu/courses/18-s997-high-dimensional-statistics-spring-2015/ 659 a69e2f53bb2eeb9464520f3027fc61e6_MIT18_S997S15_Chapter1.pdf, 2015. 660 Accessed: 2024-11-27. 661 L. S. Shapley. A value for n-person games. In H. W. Kuhn and A. W. Tucker (eds.), Contributions to 662 the Theory of Games, volume 2, pp. 307-317. Princeton Univ. Press, 1953. 663 Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks 665 against machine learning models. In Proceedings of the IEEE Symposium on Security and Privacy, 666 2017. 667 R. H. L. Sim, X. Xu, and K. H. Low. Data valuation in machine learning: "ingredients", strategies, 668 and open challenges. In Proceedings of the International Joint Conference on Artificial Intelligence, 669 2022. Survey Track. 670 671 Rachael Hwee Ling Sim, Yehong Zhang, Mun Choon Chan, and Kian Hsiang Low. Collaborative 672 machine learning with incentive-aware model rewards. In Proceedings of the International 673 Conference on Machine Learning, 2020. 674 Rachael Hwee Ling Sim, Yehong Zhang, Nghia Hoang, Xinyi Xu, Bryan Kian Hsiang Low, and 675 Patrick Jaillet. Incentives in private collaborative machine learning. In Advances in Neural 676 Information Processing Systems, 2023. 677 678 Dmitrii Usynin, Daniel Rueckert, and Georgios Kaissis. Incentivising the federation: gradient-based 679 metrics for data selection and valuation in private decentralised training. In Proceedings of the 2024 European Interdisciplinary Cybersecurity Conference, 2024. 680 681 Jiachen T. Wang and Ruoxi Jia. Data banzhaf: A robust data valuation framework for machine 682 learning. In Proceedings of the International Conference on Artificial Intelligence and Statistics, 683 2023. 684 Jiachen T. Wang, Yuqing Zhu, Yu-Xiang Wang, Ruoxi Jia, and Prateek Mittal. A privacy-friendly 685 approach to data valuation. In Advances in Neural Information Processing Systems, 2023. 686 687 Tianhao Wang, Johannes Rausch, Ce Zhang, Ruoxi Jia, and Dawn Song. A principled approach to 688 data valuation for federated learning. Lecture Notes in Computer Science, 12500:153–167, 2020. 689 Lauren Watson, Rayna Andreeva, Hao-Tsung Yang, and Rik Sarkar. Differentially private shapley 690 values for data evaluation, 2022. 691 692 Z. Wu, Y. Shu, and K. H. Low. DAVINZ: Data valuation using deep neural networks at initialization. 693 In Proceedings of the International Conference on Machine Learning, 2022. 694 Xinyi Xu, Lingjuan Lyu, Xingjun Ma, Chenglin Miao, Chuan Sheng Foo, and Kian Hsiang Low. Gradient driven rewards to guarantee fairness in collaborative machine learning. In Advances in 696 Neural Information Processing Systems, 2021. 697 Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept and 699 applications. ACM Transactions on Intelligent Systems and Technology, 2019. 700
- 701 Yuchen Zhang, John C. Duchi, and Martin J. Wainwright. Communication-efficient algorithms for statistical optimization. *Journal of Machine Learning Research*, 2013.

702	Zijian Zhou, Xinyi Xu, Rachael Hwee Ling Sim, Chuan Sheng Eoo, and Kian Hsiang Low, Probably
703	approximate shapley fairness with applications in machine learning. In <i>Proceedings of the AAAI</i>
704	Conference on Artificial Intelligence, 2023.
705	
706	
707	
708	
709	
710	
711	
712	
713	
714	
715	
716	
717	
718	
719	
720	
721	
722	
723	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
730	
737	
730	
739	
740	
7/12	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	

APPENDIX

756

758

759 760 761

763

765

766

767

768

769

770

771

774

775

777 778

780 781

782 783

784

785

786

787

788

789

Α TABLE OF NOTATIONS

Table 4: Notations. The subscript *j* is omitted where the context is clear.



В ADDITIONAL DISCUSSIONS

B.1 SEMIVALUE AS A RANDOM VARIABLE.

The exact computation of semivalues is often intractable in practice due to the need to compute an exponential number of marginal contributions $V(S \cup \{i\}) - V(S)$. Moreover, in data valuation, the utility function V is not deterministic w.r.t. S: V is commonly defined as the test accuracy or test loss (Ghorbani & Zou, 2019), which is stochastic due to (stochastic) gradient descent and random model initialization. Such stochasticity becomes more pronounced when gradients are injected with artificial noise to ensure privacy such as in DP-SGD (Abadi et al., 2016) (recalled later). Hence, the semivalue estimator ψ_i is treated as a *random variable* whose randomness comes from the marginal contributions.

790 791

793

B.2 DIFFERENTIAL PRIVACY FRAMEWORK.

794 We emphasize that our analysis adopts this definition as we rely on the Gaussian mechanism and leverage its composition property (Dwork & Roth, 2014). Our analysis holds for other DP frame-796 works capturing privacy guarantee with Gaussian mechanism and possessing composition and 797 post-processing properties, such as Rényi DP (Mironov, 2017) and z-CDP (Bun & Steinke, 2016). 798 Typically, both Rényi DP and z-CDP satisfy post-processing immunity and have composition proper-799 ties. Moreover, both provide a DP guarantee for the Gaussian mechanism. Notably, while both DP 800 frameworks provide a modestly tighter privacy analysis than compared to (Dwork & Roth, 2014), to the best of our knowledge, they still require the variance of the Gaussian noise to linearly scale 801 with k to satisfy a given level of DP guarantee. In fact, some existing works have identified that 802 the moment accountant method (Abadi et al., 2016) which we adopt is an instantiation of Rényi 803 DP (Ouadrhiri & Abdelhadi, 2022) translated back to (ϵ, δ)-DP. Lastly, we note that our theoretical 804 analysis is independent of the choice of DP framework except for requiring $z_t \sim \mathcal{N}(\mathbf{0}, k\sigma^2 \mathbf{I})$. 805

806

807 Gaussian mechanism. We note that other mechanisms exist for achieving privacy with DP guarantees. We choose the Gaussian mechanism for the convenience of mathematical analysis compared to 808 other mechanisms such as the Laplace mechanism as well as a manifold of existing works discussing the theoretical properties of the Gaussian mechanism.

B.3 COMPOSITION OF DP MECHANISMS.

811 812

While it is possible to find a theoretical privacy budget lower than the result in (Abadi et al., 2016, Theorem 1), the current asymptotic bound $\sigma^2 = \Omega(k)$ is the best we know. We also note that there have been efforts at improving the bound (Asoodeh et al., 2021). However, the results are not asymptotically better in terms of k, hence do not affect our theoretical results. Moreover, we emphasize that our theoretical analysis and proof strategy do not depend on the exact form of σ^2 other than assuming it depends on k. If a lower bound is found in the future, our analysis can be readily adapted.

819

826 827

828

829

830

831

832 833

834 835

B.4 DP GUARANTEE LEVEL AND DATA VALUE ESTIMATES

Noise introduced by DP causes a lower mean of data value estimates. Mathematically, DP imposes a limit on the leave-one-out property of the privatized mechanism \mathcal{M} . In particular, if $V(\cdot) \in [0,1]$ (e.g., when test accuracy is used as V), then we have, for any party i and subsets $S \subseteq [n] \setminus \{i\}$, the following (Dwork et al., 2015, Lemma 6):

$$|\psi_i| \coloneqq |\mathbb{E}_S[V(S \cup \{i\}) - V(S)]| \le e^{\epsilon} - 1 + \delta$$

where (ϵ, δ) are the parameters satisfying (ϵ, δ) -DP. Assuming a fixed δ , stronger DP implies a lower ϵ , hence decreased right-hand side value, i.e., the upper bound for the marginal contribution. This inequality suggests that stronger DP results in lower *absolute value* of the data value estimates. We find this a reasonable behavior as the decreased value reflects the erosion of information carried by the data. Instead, having a lower absolute data value does not forbid us to still preserve the relative order of the estimates by minimizing the estimation uncertainty.

842 843

844

845

846

847

848

849

Impact of different (ϵ, δ) -**DP guarantee levels on our method.** Through our theoretical development, we have established that the estimation uncertainty can be controlled to be *independent* of the evaluation budget k, i.e., the noise due to DP is only affected by the final DP guarantee level. That said, an overly large final DP guarantee level shrinks the gap between the performance of i.i.d. noise and without DP noise since perturbation is small, and *vice versa*. Therefore, to highlight the effectiveness of our method in controlling the estimation uncertainty, we fix a moderate final DP guarantee level $(\epsilon, \delta) = (1.0, 5 \times 10^{-5})$ throughout the main text s.t. a clear contrast in performance on various ML tasks can be observed. We also provide additional experiments on different ϵ values in App. D.5.

850 851 852

Computing Var $[\psi_j]$. One can derive Var $[\psi_j]$ given Var $[\psi_j|\theta_{\pi^1,j}^p,\ldots,\theta_{\pi^k,j}^p]$ via the law of total variance and assumptions on the inter-dependence between $\theta_{\pi^t,j}^p$'s and k, which is not the focus of this work – how DP impacts data valuation, particularly semivalue estimation – and thus left for future work.

858

B.5 SOCIETAL IMPACT

859 860

As discussed in the introduction in the main text, we believe our contribution has a huge potential
 societal impact in improving privacy, especially with the rising awareness of protecting personal
 data (Bukaty, 2019; Council of European Union, 2014). We do not find a direct path to any negative
 societal impact with our contribution.

С **PROOFS AND ADDITIONAL RESULTS**

C.1 PROOF OF THE EXAMPLE IN SEC. 5.2.

Observation C.1. For a particular party j, assume, for $t \in [k]$, the norm-clipped gradients $\hat{g}_{\pi^1,j} =$ $\hat{g}_{\pi^{2},j} = \dots = \hat{g}_{\pi^{k},j}. \text{ Denote } \forall t \in [k], \\ \tilde{g}_{\pi^{t},j} \coloneqq \hat{g}_{\pi^{t},j} + z_{t} \text{ where } z_{t} \overset{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, k(C\sigma)^{2}\boldsymbol{I}) \text{ are } i.i.d. \text{ drawn. Let } \\ \tilde{g}_{\pi^{t},j}^{*} \coloneqq t^{-1} \sum_{l=1}^{t} \tilde{g}_{\pi^{l},j}. \text{ Then, } \operatorname{Var}[\|\tilde{g}_{\pi^{t},j}^{*}\|_{2}^{2}] \to t^{-2}\operatorname{Var}[\|\tilde{g}_{\pi^{t},j}\|_{2}^{2}] \text{ as } k \to \infty.$

Proof. We make the following notations to help ease understanding. Denote $z_t^* \coloneqq t^{-1} \sum_{l=1}^t z_l \sim t^{-1} \sum_{l=1}^t z_l = t^{-1$ $\mathcal{N}(\mathbf{0}, ((C\sigma)^2/t)\mathbf{I})$. Denote $\hat{g} \coloneqq \hat{g}_{\pi^1,j} = \hat{g}_{\pi^2,j} = \ldots = \hat{g}_{\pi^k,j}$. Further denote $(\cdot)_r$ the *r*th element of a vector. Then, we have

$$\operatorname{Var}[\|\tilde{g}_{\pi^{t},j}^{*}\|_{2}^{2}] = \operatorname{Var}[\|\hat{g} + z_{t}^{*}\|_{2}^{2}]$$

$$\begin{split} &= \operatorname{Var}[\sum_{r=1}^{d} (\hat{g} + z_{t}^{*})_{r}^{2}] \\ &= \sum_{r=1}^{d} \operatorname{Var}[(\hat{g} + z_{t}^{*})_{r}^{2}] \\ &= \sum_{r=1}^{d} (2\hat{g})_{r}^{2} \operatorname{Var}[(z_{t}^{*})_{r}] + \sum_{r=1}^{d} \operatorname{Var}[(z_{t}^{*})_{r}^{2}] + 2 \sum_{r=1}^{d} \operatorname{Cov}[2(\hat{g})_{r}(z_{t}^{*})_{r}, (z_{t}^{*})_{r}^{2}] \\ &= 4 \|\hat{g}\|_{2}^{2} \operatorname{Var}[(z_{t}^{*})_{1}] + d \operatorname{Var}[(z_{t}^{*})_{1}^{2}] + 2 \sum_{r=1}^{d} \left(\mathbb{E}[2(\hat{g})_{r}(z_{t}^{*})_{r}^{3}] - \mathbb{E}[2(\hat{g})_{r}(z_{t}^{*})_{r}] \mathbb{E}[(z_{t}^{*})_{r}^{2}] \right) \\ &= 4 \|\hat{g}\|_{2}^{2} k(C\sigma)^{2} / t + 2dk^{2}(C\sigma)^{4} / t^{2} \end{split}$$

where in the last step we use the fact $\mathbb{E}[z] = \mathbb{E}[z^3] = 0$ if z follows a Normal distribution with mean 0 as well as the fact that $(t/\sigma^2)(z_t^*)_1^2 \sim \chi_1^2$ where χ_1^2 is a chi-squared distribution with degree of freedom 1. On the other hand,

$$\begin{aligned} \operatorname{Var}[\|\tilde{g}_{\pi^{t},j}\|_{2}^{2}] &= \operatorname{Var}[\|\hat{g} + z_{t}\|_{2}^{2}] \\ &= \sum_{r=1}^{d} \operatorname{Var}[(\hat{g} + z_{t})_{r}^{2}] \end{aligned}$$

d

_

=

$$2\sum_{r=1}^{d} (2\hat{g})_r \operatorname{Var}[(z_t)_r] + \sum_{r=1}^{d} \operatorname{Var}[(z_t)_r^2] + 2\sum_{r=1}^{d} \operatorname{Cov}[2(\hat{g})_r(z_t)_r, (z_t)_r^2]$$

$$4\|\hat{g}\|_2^2 \operatorname{Var}[(z_t)_1] + d\operatorname{Var}[(z_t)_1^2] + 2\sum_{r=1}^{d} \left(\mathbb{E}[2(\hat{g})_r(z_t)_r^3] - \mathbb{E}[2(\hat{g})_r(z_t)_r] \mathbb{E}[(z_t)_r^2] \right)$$

$$= 4 \|\hat{g}\|_2^2 k (C\sigma)^2 + 2dk^2 (C\sigma)^4$$

where the last step follows the same logic as the last equation. We can easily verify that

d

$$\lim_{k \to \infty} \frac{\operatorname{Var}[\|\tilde{g}_{\pi^t,j}\|_2^2]}{\operatorname{Var}[\|\tilde{g}_{\pi^t,j}\|_2^2]} = \frac{2d(C\sigma)^4}{2d(C\sigma)^4/t^2} = t^2 \ .$$

C.2 PROOF OF PROP. 5.1.

We first establish the following lemma to facilitate the proof.

Lemma C.2. If random variables X and Y are independent, then

$$\operatorname{Var}[XY] \ge \mathbb{E}[Y]^2 \operatorname{Var}[X]$$

Proof. Since X and Y are independent, X|Y = X and Y|X = Y. Then, by the law of total variance, 919 we have

920
920
920
921
921
922
923
924
924
925
926
Var[XY] =
$$\mathbb{E}[Var[XY|X]] + Var[\mathbb{X}\mathbb{E}[Y|X]]$$

 $= \mathbb{E}[X^2Var[Y]] + Var[\mathbb{X}\mathbb{E}[Y]]$
 $= \mathbb{E}[Y]^2Var[X] + Var[Y]\mathbb{E}[X^2]$
 $\geq \mathbb{E}[Y]^2Var[X]$.

Proposition C.3 (Reproduced from Prop. 5.1.). $\forall t \in [k]$, denote $\theta_{\pi^t} \coloneqq \theta_{\pi^t}^p - \alpha \tilde{g}_{\pi^t} = \theta_{\pi^t}^p - \alpha (\hat{g}_{\pi^t} + z_t)$ where $\forall t \in [k], z_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, k(C\sigma)^2 \mathbf{I})$ with $\sigma, C > 0$. Given a test dataset consisting of l data points $\mathcal{D}_{\text{test}} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_l, y_l)\}$. If V is the negated mean-squared error loss on a linear regression model

$$V(\boldsymbol{\theta}) \coloneqq -l^{-1} \sum_{i=1}^{l} (\boldsymbol{\theta}^{\top} \boldsymbol{x}_i - y_i)^2$$

or the negated ℓ_2 -regularized cross-entropy loss on a logistic regression model

V

$$\begin{split} T(\boldsymbol{\theta}) &\coloneqq l^{-1} \sum_{i=1}^{l} (1-y_i) \log(1 - \operatorname{Sig}(\boldsymbol{\theta}, \boldsymbol{x}_i)) \\ &+ l^{-1} \sum_{i=1}^{l} y_i \log(\operatorname{Sig}(\boldsymbol{\theta}, \boldsymbol{x}_i)) - \lambda \|\boldsymbol{\theta}\|_2^2 \end{split}$$

where $\operatorname{Sig}(\boldsymbol{\theta}, \boldsymbol{x}_i) = (1 + e^{-\boldsymbol{\theta}^\top \boldsymbol{x}_i})^{-1}$ is the sigmoid function and $\lambda > 0$ is the regularization hyperparameter. Denote $m(\pi^t) \coloneqq V(\boldsymbol{\theta}_{\pi^t}) - V(\boldsymbol{\theta}_{\pi^t}^p)$. Further denote a regular semivalue estimator $\psi \coloneqq k^{-1} \sum_{t=1}^k p(\pi^t) m(\pi^t)$.³ The estimation uncertainty satisfies $\operatorname{Var}[\psi|\boldsymbol{\theta}_{\pi^1}^p, \boldsymbol{\theta}_{\pi^2}^p, \dots, \boldsymbol{\theta}_{\pi^k}^p] = \Omega(k)$.

Proof of Prop. 5.1. We first analyze the variance of z_t on each data point in $\mathcal{D}_{\text{test}}$. Note that $\theta_{\pi^t,j} = \theta_{\pi^t,j}^p - \alpha(\hat{g}_{\pi^t,j} + z_t) = (\theta_{\pi^t,j}^p - \alpha \hat{g}_{\pi^t,j}) - \alpha z_t$. Denote $\bar{\theta}_{\pi^t,j} \coloneqq \theta_{\pi^t,j}^p - \alpha \hat{g}_{\pi^t,j}$. Notice that, conditional on $\theta_{\pi^t,j}^p, \bar{\theta}_{\pi^t,j}$ can be *deterministically calculated* via gradient computation since the underlying data is fixed. Then we have $\theta_{\pi^t,j} = \bar{\theta}_{\pi^t,j} - \alpha z_t$ where the randomness of $\theta_{\pi^t,j}$ arises from the randomness of z_t . Denote Denote $\underline{p} \coloneqq \min_{\pi \in \Pi} p_j(\pi) > 0$ the minimal weight, which is positive since ψ_j is a regular semivalue. Due to the independence of all z_t 's, we first consider the variance of utility

$$\operatorname{Var}\left[\frac{\sum_{t=1}^{k} p_{j}(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j})}{k} \middle| \boldsymbol{\theta}_{\pi^{1},j}^{p}, \dots, \boldsymbol{\theta}_{\pi^{k},j}^{p} \right] = \frac{\sum_{t=1}^{k} \operatorname{Var}[p_{j}(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}) | \boldsymbol{\theta}_{\pi^{t},j}^{p}]}{k^{2}}$$
$$\geq \frac{\sum_{t=1}^{k} \mathbb{E}[p_{j}(\pi^{t}) | \boldsymbol{\theta}_{\pi^{t},j}^{p}]^{2} \operatorname{Var}[V(\boldsymbol{\theta}_{\pi^{t},j}) | \boldsymbol{\theta}_{\pi^{t},j}^{p}]}{k^{2}}$$
$$\geq \frac{\sum_{t=1}^{k} \underline{p}^{2} \operatorname{Var}[V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}) | \boldsymbol{\theta}_{\pi^{t},j}^{p}]}{k^{2}}$$
$$= \frac{\sum_{t=1}^{k} \underline{p}^{2} \operatorname{Var}[V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}) | \bar{\boldsymbol{\theta}}_{\pi^{t},j}^{p}]}{k^{2}}$$

where the 2nd step is derived from Lemma C.2 since p_j and V are independent. We analyze the conditional variance w.r.t. the two loss functions respectively in the following.

Linear regression. *V* here is the negated MSE on $\mathcal{D}_{\text{test}}$. The above variance is further analyzed w.r.t. each single test data point $(\boldsymbol{x}_i, y_i) \in \mathcal{D}_{\text{test}}$ as

$$\operatorname{Var}[V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})|\bar{\boldsymbol{\theta}}_{\pi^{t},j}] = \operatorname{Var}[-\frac{1}{l}\sum_{i=1}^{l}((\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})^{\top}\boldsymbol{x}_{i} - y_{i})^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}]$$
$$= \frac{1}{l^{2}}\sum_{i=1}^{l}\operatorname{Var}[(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})^{\top}\boldsymbol{x}_{i} - y_{i})^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}].$$

³A regular semivalue has w(s) > 0 for all $s \in [n]$ (i.e., $p(\pi) > 0$ for all $\pi \in \Pi$) (Carreras & Freixas, 2002). Examples of regular semivalues include Shapley value, Banzhaf value, and Beta Shapley. 973 Note that $\forall t \in [k]$ and $\forall i \in [l]$, the parameters $\bar{\theta}_{\pi^t,j}$, test data point feature x_i , and test data point 1 label y_i are fixed conditional on $\bar{\theta}_{\pi^t,j}$. Hence, $(\bar{\theta}_{\pi^t,j} - \alpha z_t)^\top x_i - y_i$ is an affine transformation of a 974 normally distributed random variable z_t :

$$(ar{m{ heta}}_{\pi^t,j} - lpha z_t)^{ op} m{x}_i - y_i |ar{m{ heta}}_{\pi^t,j} \sim \mathcal{N}(ar{m{ heta}}_{\pi^t,j}^{ op} m{x}_i - y_i, lpha^2 k \sigma^2 m{x}_i^{ op} m{x}_i)$$

the square of which produces a noncentral chi-squared random variable with 1 degree of freedom

$$((\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t)^\top \boldsymbol{x}_i - y_i)^2 |\bar{\boldsymbol{\theta}}_{\pi^t,j} = \lambda_i \chi_1^2 (\mu_{i,t}^2 / \lambda_i)$$

where $\lambda_i = \alpha^2 k \sigma^2 \boldsymbol{x}_i^\top \boldsymbol{x}_i$ and $\mu_{i,t} = \boldsymbol{\bar{\theta}}_{\pi^t,j}^\top \boldsymbol{x}_i - y_i$.

Then, use the closed-form expression for the variance of a noncentral chi-squared random variable, Var[$(\bar{\theta}_{\pi^t,j} - \alpha z_t)^{\top} x_i - y_i)^2 |\bar{\theta}_{\pi^t,j}] = \operatorname{Var}[\lambda_i \chi_1^2(\mu_{i,t}^2) |\bar{\theta}_{\pi^t,j}] = 2\lambda_i^2(1 + 2\mu_{i,t}^2/\lambda_i) = 2\lambda_i^2 + 4\mu_{i,t}^2\lambda_i$. Denote $\beta \coloneqq 2\sum_{i=1}^l \alpha^4 \sigma^4 (x_i^{\top} x_i)^2$ and $\gamma_t \coloneqq 4\alpha^2 \sigma^2 \sum_{i=1}^l (x_i^{\top} x_i) (\bar{\theta}_{\pi^t,j}^{\top} x_i - y_i)^2$. Note that β and γ_t are constants conditional on $\bar{\theta}_{\pi^t,j}$. We may rewrite them as $\beta = \frac{1}{k^2} \sum_{i=1}^l 2\lambda_i^2$ and $\gamma_t = \frac{1}{k} \sum_{i=1}^l 4\mu_{i,t}^2\lambda_i$. With these, the variance of V is

$$\begin{aligned} \operatorname{Var}[V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})|\bar{\boldsymbol{\theta}}_{\pi^{t},j}] &= \operatorname{Var}[-\frac{1}{l}\sum_{i=1}^{l}((\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})^{\top}\boldsymbol{x}_{i} - y_{i})^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}] \\ &= \frac{1}{l^{2}}\sum_{i=1}^{l}\operatorname{Var}[(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})^{\top}\boldsymbol{x}_{i} - y_{i})^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}] \end{aligned}$$

$$= \frac{1}{l^2} \sum_{i=1} \operatorname{Var}[(\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t)^\top \boldsymbol{x}_i - y_i)^2 | \bar{\boldsymbol{\theta}}_i$$

= $\frac{\beta}{l^2} k^2 + \frac{\gamma_t}{l^2} k$.

With this, we have the lower bound of the total variance of V as

$$\begin{aligned} \operatorname{Var}\left[\frac{\sum_{t=1}^{k} p_{j}(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j})}{k} \middle| \boldsymbol{\theta}_{\pi^{1},j}^{p}, \dots, \boldsymbol{\theta}_{\pi^{k},j}^{p} \right] &\geq \frac{\sum_{t=1}^{k} \underline{p}^{2} \operatorname{Var}[V(\boldsymbol{\theta}_{\pi^{t},j}) | \boldsymbol{\theta}_{\pi^{t},j}^{p}]}{k^{2}} \\ &= \frac{\sum_{t=1}^{k} \underline{p}^{2} \operatorname{Var}[V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}) | \bar{\boldsymbol{\theta}}_{\pi^{t},j}]}{k^{2}} \\ &\geq \frac{\underline{p}^{2} \sum_{t=1}^{k} \frac{\beta}{l^{2}} k^{2}}{k^{2}} \\ &= \frac{\underline{p}^{2} k \beta}{l^{2}} \\ &= \Omega(k) \;. \end{aligned}$$

Logistic regression. V here is the negated ℓ_2 -regularized logistic loss. Similarly, we can break the variance of the utility into the variance of the loss on each test data point

1014
1015
$$\operatorname{Var}[V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})|\bar{\boldsymbol{\theta}}_{\pi^{t},j}] = \operatorname{Var}[\frac{1}{l}\sum_{i=1}^{l} y_{i} \log(1/(1 + e^{-(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})^{\top}\boldsymbol{x}_{i}})) + (1 - y) \log(1 - 1/(1 + e^{-(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})^{\top}\boldsymbol{x}_{i}}))) - \|\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}\|_{2}^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}]$$
1016
$$- \|\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}\|_{2}^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}]$$

1019
1020
$$= \frac{1}{l^2} \sum_{i=1}^{l} \operatorname{Var}[y_i \log(1/(1 + e^{-(\bar{\boldsymbol{\theta}}_{\pi^t, j} - \alpha z_t)^\top \boldsymbol{x}_i})) + (1 - y_i) \log(1 - 1/(1 + e^{-(\bar{\boldsymbol{\theta}}_{\pi^t, j} - \alpha z_t)^\top \boldsymbol{x}_i}))$$

 $- \|\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t\|_2^2 |\bar{\boldsymbol{\theta}}_{\pi^t,j}| .$

1023 Denote the logistic loss $p_i(\boldsymbol{\theta}) \coloneqq y_i \log(1/(1 + e^{-\boldsymbol{\theta}^\top \boldsymbol{x}_i})) + (1 - y_i) \log(1 - 1/(1 + e^{-\boldsymbol{\theta}^\top \boldsymbol{x}_i})) \ge 0.$ 1024 Further denote the ℓ_2 regularizer $g(\boldsymbol{\theta}) = \lambda \|\boldsymbol{\theta}\|_2^2 \ge 0$. To ease notation, we let

$$\xi_t \coloneqq (\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t)^\top \boldsymbol{x}_i | \bar{\boldsymbol{\theta}}_{\pi^t,j} \sim \mathcal{N}(\bar{\boldsymbol{\theta}}_{\pi^t,j}^\top \boldsymbol{x}_i, \alpha^2 k \sigma^2 \boldsymbol{x}_i^\top \boldsymbol{x}_i) \,.$$

1026 Let $u \coloneqq \bar{\boldsymbol{\theta}}_{\pi^t,j}^\top \boldsymbol{x}_i$ and $s^2 \coloneqq \alpha^2 k \sigma^2 \boldsymbol{x}_i^\top \boldsymbol{x}_i$. Now we show that $\mathbb{E}[-p_i(\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t)|\bar{\boldsymbol{\theta}}_{\pi^t,j}] = \mathcal{O}(\sqrt{k})$. First, consider when the label $y_i = 1$, 1027 1028 $0 \leq \mathbb{E}[-p_i(\bar{\boldsymbol{\theta}}_{\pi^t, j} - \alpha z_t) | \bar{\boldsymbol{\theta}}_{\pi^t, j}] = -\mathbb{E}[\log(1/(1 + e^{-\xi_t}))]$ 1029 $= \mathbb{E}[\log(1 + e^{-\xi_t})]$ 1030 1031 $= \int_{0}^{\infty} p(\xi_t) \log(1 + e^{-\xi_t}) \mathrm{d}\xi_t$ 1032 1033 $= \int_{-\infty}^{0} p(\xi_t) \log(1 + e^{-\xi_t}) d\xi_t + \int_{0}^{\infty} p(\xi_t) \log(1 + e^{-\xi_t}) d\xi_t$ 1034 1035 $\leq \int_{0}^{0} p(\xi_t)(1+\log(e^{-\xi_t}))\mathrm{d}\xi_t + \int_{0}^{\infty} p(\xi_t)\log 2\mathrm{d}\xi_t$ $\leq \int_{-\infty}^{0} p(\xi_t)(1-\xi_t) \mathrm{d}\xi_t + \int_{-\infty}^{\infty} p(\xi_t) \log 2\mathrm{d}\xi_t$ 1039 1040 1041 $= \int_{0}^{0} p(\xi_t)(1-\xi_t) \mathrm{d}\xi_t + \log 2$ 1043 $\leq \log 2 + 1 + \int_0^0 p(\xi_t)(-\xi_t) \mathrm{d}\xi_t$ 1045 $= \log 2 + 1 + \frac{1}{s\sqrt{2\pi}} \int_{-\infty}^{0} (-\xi_t) \exp\left(-\frac{(\xi_t - u)^2}{2s^2}\right) d\xi_t$ 1046 1047 1048 $= \log 2 + 1 + \left(-\frac{1}{2}\mu \operatorname{erfc}\left(\frac{\sqrt{2}\mu}{2s}\right) + \sqrt{\frac{1}{2\pi}}\operatorname{sexp}\left(-\frac{u^2}{2s^2}\right)\right)$ 1049 1050 1051 $\leq \log 2 + 1 + |u| + s \sqrt{\frac{1}{2\pi}} \exp\left(-\frac{u^2}{2s^2}\right)$ 1052 1053 $\leq \log 2 + 1 + |u| + s$ 1054 $= \mathcal{O}(s)$ 1055 1056 $= \mathcal{O}(\sqrt{k})$ 1057 where erfc represents the complementary error function and the integral $\frac{1}{s\sqrt{2\pi}}\int_{-\infty}^{0}(-\xi_t)\exp\left(-\frac{(\xi_t-u)^2}{2s^2}\right)d\xi_t = \left(-\frac{1}{2}\mu\operatorname{erfc}\left(\frac{\sqrt{2}\mu}{2s}\right) + \sqrt{\frac{1}{2\pi}}\operatorname{sexp}\left(-\frac{u^2}{2s^2}\right)\right)$ can be de-1058 1059 rived with a math solver (in our case, the "sympy" package of Python, and a notebook involving the code snippet for reproducing the result is included in the supplementary materials). Similarly, when 1061 the label $y_i = 0$, 1062 $\mathbb{E}[-p_i(\bar{\boldsymbol{\theta}}_{\pi^t,i} - \alpha z_t)|\bar{\boldsymbol{\theta}}_{\pi^t,i}] = \mathbb{E}[-\log(e^{-\xi_t}/(1 + e^{-\xi_t}))]$ 1064 $= -\mathbb{E}[-\mathcal{E}_t - \log(1 + e^{-\xi_t})]$ $= \mathbb{E}[\xi_t] + \mathbb{E}[\log(1 + e^{-\xi_t}]]$ 1067 $\leq \mathbb{E}[\xi_t] + \log 2 + 1 + \int_0^0 p(\xi_t)(-\xi_t) d\xi_t$ 1068 1069 $= \log 2 + 1 + \int_{0}^{\infty} p(\xi_t) \xi_t \mathrm{d}\xi_t$ 1070 1071 $=\log 2 + 1 + \frac{1}{e\sqrt{2\pi}} \int_0^\infty \xi_t \exp\left(-\frac{(\xi_t - u)^2}{2s^2}\right) \mathrm{d}\xi_t$ 1072 1074 $\leq \log 2 + 1 + \frac{1}{2s} \left(2|us| + \sqrt{\frac{2}{\pi}} s^2 \exp\left(-\frac{u^2}{2s^2}\right) \right)$ 1075 1076 1077 $= \mathcal{O}(\sqrt{k})$. 1078 In conclusion, $\mathbb{E}[-p(\bar{\theta}_{\pi^t,j} - \alpha z_t)|\bar{\theta}_{\pi^t,j}] = \mathcal{O}(\sqrt{k})$. Next, we work out the expectation of $\|\bar{\theta}_{\pi^t,j} - \alpha z_t\|$ 1079 $\alpha z_t \|_2^2 |\bar{\theta}_{\pi^t,j}$. Note that $\bar{\theta}_{\pi^t,j} - \alpha z_t \in \mathbb{R}^d$ where d is the dimension of the model parameters. By

linearity of expectation, we have

$$\mathbb{E}[\|\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}\|_{2}^{2} |\bar{\boldsymbol{\theta}}_{\pi^{t},j}] = \sum_{r=1}^{d} \mathbb{E}[(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})_{r}^{2} |\bar{\boldsymbol{\theta}}_{\pi^{t},j}]$$
$$= \sum_{r=1}^{d} \left(\mathbb{E}[\bar{\boldsymbol{\theta}}_{\pi^{t},j}^{2} |\bar{\boldsymbol{\theta}}_{\pi^{t},j}] - 2\alpha \bar{\boldsymbol{\theta}}_{\pi^{t},j} \mathbb{E}[z_{t}] + \alpha^{2} \mathbb{E}[z_{t}^{2}] \right)$$
$$= \bar{\boldsymbol{\theta}}_{\pi^{t},j}^{2} + \alpha^{2} dk \sigma^{2}$$
$$= \mathcal{O}(k) .$$

With this, we can expand the following square of expectation as

$$\begin{split} & \mathbb{E}[-p(\bar{\theta}_{\pi^{t},j} - \alpha z_{t}) + \|\bar{\theta}_{\pi^{t},j} - \alpha z_{t}\|_{2}^{2} |\bar{\theta}_{\pi^{t},j}|^{2} \\ & = \mathbb{E}[-p(\bar{\theta}_{\pi^{t},j} - \alpha z_{t})|\bar{\theta}_{\pi^{t},j}]^{2} - 2\mathbb{E}[-p(\bar{\theta}_{\pi^{t},j} - \alpha z_{t})|\bar{\theta}_{\pi^{t},j}]\mathbb{E}[\|\bar{\theta}_{\pi^{t},j} - \alpha z_{t}\|_{2}^{2} |\bar{\theta}_{\pi^{t},j}] + \mathbb{E}[\|\bar{\theta}_{\pi^{t},j} - \alpha z_{t}\|_{2}^{2} |\bar{\theta}_{\pi^{t},j}]^{2} \\ & = \mathcal{O}(k) + \mathcal{O}(k\sqrt{k}) + \mathbb{E}[\|\bar{\theta}_{\pi^{t},j} - \alpha z_{t}\|_{2}^{2} |\bar{\theta}_{\pi^{t},j}]^{2} \\ & = \mathcal{O}(k\sqrt{k}) + \mathbb{E}[\|\bar{\theta}_{\pi^{t},j} - \alpha z_{t}\|_{2}^{2} |\bar{\theta}_{\pi^{t},j}]^{2} \\ & = \mathcal{O}(k\sqrt{k}) + \mathbb{E}[\|\bar{\theta}_{\pi^{t},j} - \alpha z_{t}\|_{2}^{2} |\bar{\theta}_{\pi^{t},j}]^{2} . \end{split}$$
Next, consider that

ext, consider that

$$\mathbb{E}[(-p(\bar{\boldsymbol{\theta}}_{\pi^{t},j}-\alpha z_{t})+\|\bar{\boldsymbol{\theta}}_{\pi^{t},j}-\alpha z_{t}\|_{2}^{2})^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}] \geq \mathbb{E}[\|\bar{\boldsymbol{\theta}}_{\pi^{t},j}-\alpha z_{t}\|_{2}^{4}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}]$$

So, we can derive a bound of the variance from below by

$$\begin{aligned} & \text{Var}[-p(\bar{\theta}_{\pi^{t},j} - \alpha z_{t}) + \|\bar{\theta}_{\pi^{t},j} - \alpha z_{t}\|_{2}^{2} |\bar{\theta}_{\pi^{t},j}] \geq \mathbb{E}[\|\bar{\theta}_{\pi^{t},j} - \alpha z_{t}\|_{2}^{4} |\bar{\theta}_{\pi^{t},j}] - (\mathcal{O}(k\sqrt{k}) + \mathbb{E}[\|\bar{\theta}_{\pi^{t},j} - \alpha z_{t}\|_{2}^{2} |\bar{\theta}_{\pi^{t},j}]^{2}) \\ & = \text{Var}[\|\bar{\theta}_{\pi^{t},j} - \alpha z_{t}\|_{2}^{2} |\bar{\theta}_{\pi^{t},j}] - \mathcal{O}(k\sqrt{k}) \\ & \text{1105} \\ & = \sum_{r=1}^{d} \text{Var}[(\bar{\theta}_{\pi^{t},j} - \alpha z_{t})_{r}^{2} |\bar{\theta}_{\pi^{t},j}] - \mathcal{O}(k\sqrt{k}) . \end{aligned}$$

Consider that for each $r \in [d]$ where d is the dimension of the model parameters,

$$(\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t)_r | \bar{\boldsymbol{\theta}}_{\pi^t,j} \sim \mathcal{N}((\bar{\boldsymbol{\theta}}_{\pi^t,j})_r, \alpha^2 k \sigma^2)$$

and squaring it produces a noncentral chi-squared random variable with 1 degree of freedom

$$(\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t)_r^2 | \bar{\boldsymbol{\theta}}_{\pi^t,j} \sim \lambda \chi_1^2(\mu^2/\lambda)$$

where $\lambda = \alpha^2 k \sigma^2$ and $\mu = (\theta_{\pi^t,j})_r$. By the closed-form expression for the variance of a noncentral chi-squared random variable,

$$\begin{aligned} \operatorname{Var}[(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})_{r}^{2} | \bar{\boldsymbol{\theta}}_{\pi^{t},j}] &= \operatorname{Var}[\lambda \chi_{1}^{2}(\mu^{2}/\lambda)] \\ &= 2\lambda^{2}(1 + 2\mu^{2}/\lambda) \\ &= 2\alpha^{4}k^{2}\sigma^{2} + 4(\bar{\boldsymbol{\theta}}_{\pi^{t},j})_{r}^{2}\alpha^{2}k\sigma^{2} \,. \end{aligned}$$

Plug this result back into the previous inequality,

$$\begin{aligned} &\operatorname{Var}[-p(\bar{\boldsymbol{\theta}}_{\pi^{t},j}-\alpha z_{t})+\|\bar{\boldsymbol{\theta}}_{\pi^{t},j}-\alpha z_{t}\|_{2}^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}] \geq 2d\alpha^{4}k^{2}\sigma^{2}+4d(\bar{\boldsymbol{\theta}}_{\pi^{t},j})_{r}^{2}\alpha^{2}k\sigma^{2}-\mathcal{O}(k\sqrt{k}) \\ &\operatorname{Hence}, \end{aligned}$$

$$\begin{aligned} \operatorname{Var}\left[\frac{\sum_{t=1}^{k} p_{j}(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j})}{k} \middle| \boldsymbol{\theta}_{\pi^{1},j}^{p}, \dots, \boldsymbol{\theta}_{\pi^{k},j}^{p} \right] &\geq \frac{\sum_{t=1}^{k} \underline{p}^{2} \operatorname{Var}[V(\boldsymbol{\theta}_{\pi^{t},j}) | \boldsymbol{\theta}_{\pi^{t},j}^{p}]}{k^{2}} \\ &= \frac{\underline{p}^{2} \sum_{t=1}^{k} \operatorname{Var}[V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}) | \bar{\boldsymbol{\theta}}_{\pi^{t},j}]}{k^{2}} \\ &\geq \frac{\underline{p}^{2} \sum_{t=1}^{k} 2d\alpha^{4}k^{2}\sigma^{2} - \mathcal{O}(k\sqrt{k})}{k^{2}} \\ &= 2\underline{p}^{2}d\alpha^{4}\sigma^{2}k - \mathcal{O}(\sqrt{k}) \\ &= \Omega(k) \;. \end{aligned}$$

Lastly, we derive the conditional variance of a semivalue estimator with independent noise as

1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146

$$Var[\psi_{j}|\theta_{\pi^{1},j}^{p}, \theta_{\pi^{2},j}^{p}, \dots, \theta_{\pi^{k},j}^{p}] = Var\left[\frac{\sum_{t=1}^{k} p_{j}(\pi^{t})V(\theta_{\pi^{t},j})}{k} \middle| \bar{\theta}_{\pi^{1},j}, \dots, \bar{\theta}_{\pi^{k},j} \right]$$

$$\geq Var\left[\frac{\sum_{t=1}^{k} p_{j}(\pi^{t})V(\theta_{\pi^{t},j})}{k} \middle| \bar{\theta}_{\pi^{1},j}, \dots, \bar{\theta}_{\pi^{k},j} \right]$$

$$\geq \frac{\sum_{t=1}^{k} p^{2} Var[V(\theta_{\pi^{t},j})|\bar{\theta}_{\pi^{t},j}]}{k^{2}}$$

$$= \Omega(k) .$$

1147 C.3 MORE GENERAL PROPOSITION FOR ESTIMATION UNCERTAINTY WITH I.I.D. NOISE

Proposition C.4. (I.I.D. Noise More General) Denote $\bar{\theta}_{\pi^t,j} \coloneqq \theta_{\pi^t,j}^p - \alpha \hat{g}_{\pi^t,j}$ and $\theta_{\pi^t,j} \coloneqq$ $\theta_{\pi^t,j}^p - \alpha \tilde{g}_{\pi^t,j}$. Suppose V is "locally" m-strongly concave in a domain \mathcal{D} which contains the model parameters in all k evaluations up to a margin αK , i.e. the union of $\ell_{\infty}(\alpha K)$ balls $\bigcup_{t \in [k]} \{ \theta \in$ \mathbb{R}^d : $\|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}_{\pi^t,j}\|_{\infty} \leq \alpha K \} \subseteq \mathcal{D}$. Let $\tilde{g}_{\pi^t,j} \coloneqq \hat{g}_{\pi^t,j} + z_t$ where each z_t follows a truncated Gaussian distribution bounded in the range $[-K\mathbf{1}, K\mathbf{1}]$, i.e., $z_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_{\text{trunc}}(\mathbf{0}, k(C\sigma)^2 \mathbf{I}, -K\mathbf{1}, K\mathbf{1})$.⁴ Assume that $\forall t \in [k], \ \bar{\boldsymbol{\theta}}_{\pi^t, j}$ is close to a local optimum in the sense that $\|\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^t, j})\|_2 < 1$ $\min(\frac{m}{2}, -V(\bar{\theta}_{\pi^{t}, j})). \text{ Let } \psi_{j} \coloneqq k^{-1} \sum_{t=1}^{k} p_{j}(\pi^{t}) [V(\theta_{\pi^{t}, j}) - V(\theta_{\pi^{t}, j}^{p})]. \text{ Then, for } k \leq K^{2}/(C\sigma)^{2},$ $\operatorname{Var}[\psi_j | \boldsymbol{\theta}_{\pi^1, j}^p, \dots, \boldsymbol{\theta}_{\pi^k, j}^p] = \Omega(k).$

Proof. First, denote φ and Φ the pdf and cdf of a standard normal distribution. Let a truncated normal 1160 distribution $X \sim \mathcal{N}_{\text{trunc}}(\mu, \sigma^2, a, b)$. Further denote $\alpha \coloneqq \frac{a-\mu}{\sigma}, \beta \coloneqq \frac{b-\mu}{\sigma}$. Then, the probability 1161 density function (pdf) of X is

$$f(x;\mu,\sigma,a,b) \coloneqq \frac{1}{\sigma} \frac{\varphi(\xi)}{\Phi(\beta) - \Phi(\alpha)}$$

X admits a closed-form formula for the variance as

$$\operatorname{Var}[X] = \sigma^2 \left[1 - \frac{\beta \varphi(\beta) - \alpha \varphi(\alpha)}{Z} - \left(\frac{\varphi(\alpha) - \varphi(\beta)}{Z}\right)^2 \right] \le \mathbb{E}[X^2]$$

$$(x\varphi(x))' = \frac{1}{\sqrt{2\pi}}(xe^{-x^2/2})' = \frac{1}{\sqrt{2\pi}}(1-2x^2)e^{-x^2/2} < 0$$

 $\sqrt{2\pi}$ $\sqrt{2\pi}$ 1176 Therefore, $\beta\varphi(\beta) - \alpha\varphi(\alpha) \le 2\varphi(1) = \frac{2}{\sqrt{2\pi}}e^{-1/2}$. With these, we have that

$$1 - \frac{\beta\varphi(\beta) - \alpha\varphi(\alpha)}{Z} - \left(\frac{\varphi(\alpha) - \varphi(\beta)}{Z}\right)^2 \ge 1 - \frac{2}{0.6 \times \sqrt{2\pi}} - \left(\frac{e^{-1/2}}{0.6 \times \sqrt{2\pi}}\right)^2 \ge 1 - 0.81 - 0.17$$
$$\ge 0.02.$$

1182 Consider that

$$\mathbb{E}[\|\alpha z_t\|_2^2] = \alpha^2 \sum_{r=1}^d \mathbb{E}[\xi_r^2] \ge \alpha^2 \sum_{r=1}^d \operatorname{Var}[\xi] \ge 0.02\alpha^2 dk (C\sigma)^2 .$$
1186

1187 ⁴Let $X \sim \mathcal{N}(\mu, \sigma^2)$. Then, conditional on a < X < b, X follows a truncated normal distribution $\mathcal{N}_{\text{trunc}}(\mu, \sigma^2, a, b)$.

Next, notice that $\theta_{\pi^t,j} - \bar{\theta}_{\pi^t,j} = \alpha(\hat{g}_{\pi^t,j} - \tilde{g}_{\pi^t,j}) = -\alpha z_t$, i.e., $\theta_{\pi^t,j} = \bar{\theta}_{\pi^t,j} - \alpha z_t$. Since z_t is clipped by K1, we have that $\forall t \in k, |\theta_{\pi^t,j} - \bar{\theta}_{\pi^t,j}| = \alpha |z_t| \leq \alpha K1$. This suggests that, $\forall t \in [k], \theta_{\pi^t,j}$ lies in the union of $\ell_{\infty}(\alpha K)$ balls, and therefore, $\forall t \in [k], \theta_{\pi^t,j} \in \mathcal{D}$. Let $M \coloneqq \max_{t \in [k]} \|\nabla_{\theta} V(\bar{\theta}_{\pi^t,j})\|_2 < \frac{m}{2}$. Recall that if a function f is m-strongly convex, there is

$$f(y) \ge f(x) + \nabla f(x)^{\top} (y - x) + \frac{m}{2} ||y - x||^2$$
.

1195 Since V is m-strongly concave, -V is m-strongly convex. Therefore, we have for -V

$$-V(\boldsymbol{\theta}_{\pi^{t},j}) = -V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}) \geq -V(\bar{\boldsymbol{\theta}}_{\pi^{t},j}) + \nabla_{\boldsymbol{\theta}}V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{\top}\alpha z_{t} + \frac{m}{2}\|\alpha z_{t}\|_{2}^{2} \geq 0.$$

By strong convexity, we also have

$$-V(\bar{\boldsymbol{\theta}}_{\pi^{t},j}) \geq -V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}) - \nabla_{\boldsymbol{\theta}}V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})^{\top}\alpha z_{t} + \frac{m}{2}\|\alpha z_{t}\|_{2}^{2}$$
$$-V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}) \leq -V(\bar{\boldsymbol{\theta}}_{\pi^{t},j}) + \nabla_{\boldsymbol{\theta}}V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})^{\top}\alpha z_{t} - \frac{m}{2}\|\alpha z_{t}\|_{2}^{2}.$$

1204 Consider that

 $\mathbb{E}[(-V(\boldsymbol{\theta}_{\pi^{t},j}))^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}]$ $= \mathbb{E}[(-V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}))^{2} | \bar{\boldsymbol{\theta}}_{\pi^{t},j}]$ $\geq \mathbb{E}[(-V(\bar{\boldsymbol{\theta}}_{\pi^{t},j}) + \nabla_{\boldsymbol{\theta}}V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{\top}\alpha z_{t} + \frac{m}{2}\|\alpha z_{t}\|_{2}^{2})^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}]$ $\geq V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{2} + \frac{m^{2}}{4} \mathbb{E}[\|\alpha z_{t}\|_{2}^{4}] - mV(\bar{\boldsymbol{\theta}}_{\pi^{t},j})\mathbb{E}[\|\alpha z_{t}\|_{2}^{2}]$ $-2\alpha V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})\mathbb{E}[\nabla_{\boldsymbol{\theta}}V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{\top}z_{t}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}]+m\alpha\mathbb{E}[(\nabla_{\boldsymbol{\theta}}V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{\top}z_{t})||z_{t}||_{2}^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}]$ $\geq V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{2} + \frac{m^{2}}{{}^{\scriptscriptstyle{\mathcal{I}}}} \mathbb{E}[\|\alpha z_{t}\|_{2}^{4}] - 2\alpha V(\bar{\boldsymbol{\theta}}_{\pi^{t},j}) \mathbb{E}[\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{\top} z_{t} | \bar{\boldsymbol{\theta}}_{\pi^{t},j}] + m\alpha \mathbb{E}[(\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{\top} z_{t}) \|z_{t}\|_{2}^{2} | \bar{\boldsymbol{\theta}}_{\pi^{t},j}] \,.$

¹²¹⁶ Consider that, similar to a (non-truncated) normal distribution, $\mathbb{E}[\xi_r] = \mathbb{E}[\xi_r^3] = 0$ due to symmetry ¹²¹⁷ of ξ_r , we have

$$\mathbb{E}[\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{\top} z_{t} | \bar{\boldsymbol{\theta}}_{\pi^{t},j}] = \sum_{r=1}^{d} \mathbb{E}[(\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^{t},j}))_{r} \xi_{r} | \bar{\boldsymbol{\theta}}_{\pi^{t},j}] = \sum_{r=1}^{d} (\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^{t},j}))_{r} \mathbb{E}[\xi_{r}] = 0$$

1222 and

1223
1224
$$\mathbb{E}[(\nabla_{\theta}V(\bar{\theta}_{\pi^{t},j})^{\top}z_{t})\|z_{t}\|_{2}^{2}|\bar{\theta}_{\pi^{t},j}] = \sum_{r=1}^{d}\mathbb{E}[(\nabla_{\theta}V(\bar{\theta}_{\pi^{t},j}))_{r}\xi_{r}^{3}|\bar{\theta}_{\pi^{t},j}] + \sum_{r\neq r'}\mathbb{E}[(\nabla_{\theta}V(\bar{\theta}_{\pi^{t},j}))_{r}\xi_{r}\xi_{r}^{2}|\bar{\theta}_{\pi^{t},j}]$$
1226
1227
$$= \sum_{r=1}^{d}(\nabla_{\theta}V(\bar{\theta}_{\pi^{t},j}))_{r}\mathbb{E}[\xi_{r}^{3}] + \sum_{r\neq r'}\mathbb{E}[(\nabla_{\theta}V(\bar{\theta}_{\pi^{t},j}))_{r}\xi_{r}|\bar{\theta}_{\pi^{t},j}]\mathbb{E}[\xi_{r'}^{2}]$$
1228
$$= 0 + 0$$
1230
$$= 0.$$

1232 Therefore,

$$\mathbb{E}[(-V(\boldsymbol{\theta}_{\pi^{t},j}))^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}] \geq V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{2} + \frac{m^{2}}{4}\mathbb{E}[\|\alpha z_{t}\|_{2}^{4}].$$

Since $|z_t| \leq K\mathbf{1}$, $(\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t)$ is closed and bounded. By the extreme value theorem, there exists a maximum of $\|\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t)\|_2^2$. We may let $\|\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t)\|_2^2 \leq D$. Then, $\|\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t)\|_2^2 \leq D$. Then, $\|\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t)\|_2^2 \leq D$. With this, we have $\mathbb{E}[\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t)^\top \alpha z_t |\bar{\boldsymbol{\theta}}_{\pi^t,j}] \leq \mathbb{E}[\|\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t)\|_2 \|\alpha z_t\|_2 |\bar{\boldsymbol{\theta}}_{\pi^t,j}]$

1239
$$\mathbb{E}[\nabla_{\boldsymbol{\theta}} V(\boldsymbol{\theta}_{\pi^{t},j} - \alpha z_{t}) \ \alpha z_{t} | \boldsymbol{\theta}_{\pi^{t},j}] \leq \mathbb{E}[\|\nabla_{\boldsymbol{\theta}} V(\boldsymbol{\theta}_{\pi^{t},j} - \alpha z_{t})\|_{2} \|\alpha z_{t}\|_{2} | \boldsymbol{\theta}_{\pi^{t},j}]$$
1240
$$= \|\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})\|_{2} \mathbb{E}[\|\alpha z_{t}\|_{2}]$$

$$\leq \sqrt{D}\mathbb{E}[\|\alpha z_t\|_2]$$

1242 and 1243 $\mathbb{E}[\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})^{\top} \alpha z_{t} | \bar{\boldsymbol{\theta}}_{\pi^{t},j}] \geq \mathbb{E}[-\|\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t})\|_{2} \|\alpha z_{t}\|_{2} | \bar{\boldsymbol{\theta}}_{\pi^{t},j}]$ 1244 $= -\|\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^{t}, i} - \alpha z_{t})\|_{2} \mathbb{E}[\|\alpha z_{t}\|_{2}]$ 1245 $\geq -\sqrt{D}\mathbb{E}[\|\alpha z_t\|_2],$ 1246 where the 1st step in both inequalities above is by Cauchy-Scharwz inequality. We can further bound 1247 1248 $\mathbb{E}[\|\alpha z_t\|_2]$ by $\mathbb{E}[\|\alpha z_t\|_2] = \sqrt{\mathbb{E}[\|\alpha z_t\|_2]^2}$ 1249 1250 $\leq \sqrt{\mathbb{E}[\|\alpha z_t\|_2^2]}$ 1251 1252 $= \sqrt{\alpha^2 \sum_{i=1}^{d} \mathbb{E}[\xi_r^2]}$ 1253 1255 $= \alpha \sqrt{\sum_{i=1}^{d} \operatorname{Var}[\xi]}$ 1256 1257 1259 $= \alpha \sqrt{\sum_{r=1}^{d} k(C\sigma)^2}$ 1260 1261 1262 $= \alpha C \sigma \sqrt{dk}$ 1263 where in the 2nd and 4th steps we use the formula $\operatorname{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$. With this, we have 1264 1265 $\mathbb{E}[-V(\bar{\boldsymbol{\theta}}_{\pi^{t},j}-\alpha z_{t})|\bar{\boldsymbol{\theta}}_{\pi^{t},j}]^{2} \leq \mathbb{E}[-V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})+\nabla_{\boldsymbol{\theta}}V(\bar{\boldsymbol{\theta}}_{\pi^{t},j}-\alpha z_{t})^{\top}\alpha z_{t}-\frac{m}{2}\|\alpha z_{t}\|_{2}^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}|^{2}$ 1266 $\leq V(\bar{\theta}_{\pi^{t},j})^{2} + D\mathbb{E}[\|\alpha z_{t}\|_{2}]^{2} + \frac{m^{2}}{4}\mathbb{E}[\|\alpha z_{t}\|_{2}^{2}]^{2}$ 1267 1268 $-2\alpha V(\bar{\boldsymbol{\theta}}_{\pi^{t},i})\mathbb{E}[\nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^{t},i}-\alpha z_{t})^{\top}\alpha z_{t}|\bar{\boldsymbol{\theta}}_{\pi^{t},i}] + mV(\bar{\boldsymbol{\theta}}_{\pi^{t},i})\mathbb{E}[\|\alpha z_{t}\|_{2}^{2}]$ 1270 $-m\mathbb{E}[\nabla_{\boldsymbol{\theta}}V(\bar{\boldsymbol{\theta}}_{\pi^{t}}|_{i}-\alpha z_{t})^{\top}\alpha z_{t}|\bar{\boldsymbol{\theta}}_{\pi^{t}}|_{i}]\mathbb{E}[\|\alpha z_{t}\|_{2}^{2}]$ $\leq V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{2} + D\mathbb{E}[\|\alpha z_{t}\|_{2}^{2}] + \frac{m^{2}}{4}\mathbb{E}[\|\alpha z_{t}\|_{2}^{2}]^{2} - 2\alpha V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})\sqrt{D}\mathbb{E}[\|\alpha z_{t}\|_{2}]$ 1272 $+ m\sqrt{D}\mathbb{E}[\|\alpha z_t\|_2]\mathbb{E}[\|\alpha z_t\|_2^2]$ 1274 $\leq V(\bar{\boldsymbol{\theta}}_{\pi^t,j})^2 + \frac{m^2}{{}^{\mathcal{A}}} \mathbb{E}[\|\alpha z_t\|_2^2]^2 + (D + m\sqrt{D}\alpha\sigma\sqrt{dk})\mathbb{E}[\|\alpha z_t\|_2^2] - 2\alpha^2 V(\bar{\boldsymbol{\theta}}_{\pi^t,j})\sqrt{D}\sigma\sqrt{dk}]$ 1276 $= V(\bar{\boldsymbol{\theta}}_{\pi^t,j})^2 + \frac{m^2}{4} \mathbb{E}[\|\alpha z_t\|_2^2]^2 + (D + m\sqrt{D}\alpha\sigma\sqrt{dk})\alpha^2 dk\sigma^2 - 2\alpha^2 V(\bar{\boldsymbol{\theta}}_{\pi^t,j})\sqrt{D}\sigma\sqrt{dk}]^2 + (D + m\sqrt{D}\alpha\sigma\sqrt{dk})\alpha^2 dk\sigma^2 - 2\alpha^2 V(\bar{\boldsymbol{\theta}}_{\pi^t,j})\sqrt{D}\sigma\sqrt{dk}$ 1278 1279 $= \frac{m^2}{4} \mathbb{E}[\|\alpha z_t\|_2^2]^2 + \mathcal{O}(k\sqrt{k})$ 1280 1281 where in the last step we make use of the fact $\mathbb{E}[\|\alpha z_t\|_2^2] = \alpha^2 \sum_{r=1}^d \mathbb{E}[\xi_r^2] = \alpha^2 dk (C\sigma)^2$. Therefore, 1282 1283 $\operatorname{Var}[V(\boldsymbol{\theta}_{\pi^{t},j})|\bar{\boldsymbol{\theta}}_{\pi^{t},j}] = \mathbb{E}[V(\boldsymbol{\theta}_{\pi^{t},j})^{2}|\bar{\boldsymbol{\theta}}_{\pi^{t},j}] - \mathbb{E}[V(\boldsymbol{\theta}_{\pi^{t},j})|\bar{\boldsymbol{\theta}}_{\pi^{t},j}]^{2}$ 1284 $\geq \frac{m^2}{4} \mathbb{E}[\|\alpha z_t\|_2^4] - \frac{m^2}{4} \mathbb{E}[\|\alpha z_t\|_2^2]^2 + \Omega(k\sqrt{k})$ 1285 1286 $= \operatorname{Var}[\|\alpha z_t\|_2^2] + \Omega(k\sqrt{k}) .$ 1287 Note that $\operatorname{Var}[\|\alpha z_t\|_2^2] = \operatorname{Var}[\sum^d \xi_r^2]$ 1290

1290 1290 1291 1292 1293 1294 $Var[\|\alpha z_t\|_2^2] = Var[\sum_{r=1}^{d} \xi_r^2]$ $= \sum_{r=1}^{d} Var[\xi_r^2] .$

Intuitively, $Var[\xi_r^2]$ should be asymptotically the same as the variance of a chi-squared random variable, i.e., $\Omega(k^2)$. However, because of truncation, the expression of $Var[\xi_r^2]$ is much more complicated than that of a chi-squared random variable. We use a program built with Python Sympy to analyze the exact form of $Var[\xi_r^2]$. The relevant code (written in jupyter notebook) has been included in the supplementary materials. Particularly, its variance has the following expression:

$$\begin{array}{c} 1300\\ 1301\\ 1302\\ 1302\\ 1303\\ 1304\\ 1305 \end{array} - \frac{k^2 \sigma^4 \left(-\frac{\sqrt{2K}e^{-\frac{K^2}{2k\sigma^2}}}{\sqrt{\pi}\sqrt{k\sigma}} + \operatorname{erf}\left(\frac{\sqrt{2K}}{2\sqrt{k\sigma}}\right)\right)^2}{\operatorname{erf}^2\left(\frac{\sqrt{2K}}{2\sqrt{k\sigma}}\right)} + \frac{k^2 \sigma^4 \left(-\frac{6\sqrt{2K}e^{-\frac{K^2}{2k\sigma^2}}}{\sqrt{\pi}\sqrt{k\sigma}} - \frac{\sqrt{2}\left(\frac{K^3e^{-\frac{K^2}{2k\sigma^2}}}{k^{3/2}\sigma^3} - \frac{3Ke^{-\frac{K^2}{2k\sigma^2}}}{\sqrt{k\sigma}}\right)}{\sqrt{\pi}} + 3\operatorname{erf}\left(\frac{\sqrt{2K}}{2\sqrt{k\sigma}}\right) \right)} \\ = \frac{erf^2\left(\frac{\sqrt{2K}}{2\sqrt{k\sigma}}\right)}{\operatorname{erf}\left(\frac{\sqrt{2K}}{2\sqrt{k\sigma}}\right)} + \frac{erf\left(\frac{\sqrt{2K}}{2\sqrt{k\sigma}}\right)}{\operatorname{erf}\left(\frac{\sqrt{2K}}{2\sqrt{k\sigma}}\right)} + \frac{erf\left(\frac{\sqrt{2K}}{2\sqrt{k\sigma}}\right)}{\operatorname{erf}\left(\frac{\sqrt{2K}$$

Note that we have $K^2 \ge k\sigma^2$. As such, we may denote $v := K^2/(k\sigma^2) \ge 1$. Then, the above expression can be reduced to

$$k^{2}\sigma^{4}\left(-\frac{\left(-\frac{\sqrt{2}v^{0.5}e^{-\frac{v}{2}}}{\sqrt{\pi}}+\operatorname{erf}\left(\frac{\sqrt{2}v^{0.5}}{2}\right)\right)^{2}}{\operatorname{erf}^{2}\left(\frac{\sqrt{2}v^{0.5}}{2}\right)}+\frac{-\frac{6\sqrt{2}v^{0.5}e^{-\frac{v}{2}}}{\sqrt{\pi}}-\frac{\sqrt{2}\left(-3v^{0.5}e^{-\frac{v}{2}}+v^{1.5}e^{-\frac{v}{2}}\right)}{\sqrt{\pi}}+3\operatorname{erf}\left(\frac{\sqrt{2}v^{0.5}}{2}\right)}{\operatorname{erf}\left(\frac{\sqrt{2}v^{0.5}}{2}\right)}\right)$$

Let the part in the large bracket be f(v). A plot of f is in Fig. 4. Note that as $v \to \infty$, $\operatorname{erf}\left(\frac{\sqrt{2}v^{0.5}}{2}\right)$ becomes the dominant term, thus $f \rightarrow 2$ as $v \rightarrow 1$, f decreases, but stays positive. With these, we conclude that $\operatorname{Var}[\xi_r^2] \propto k^2$. As such, $\operatorname{Var}[V(\boldsymbol{\theta}_{\pi^t,j})] = \Omega(k^2)$.



Figure 4: Plot of f vs. $v \in [1, 50]$.

Note that, when $\theta_{\pi^t, j}^p$ is determined, so is $\bar{\theta}_{\pi^t, j}$. Denote $p \coloneqq \min_{\pi \in \Pi} p_j(\pi)$ and apply Lemma C.2. Since each z_t is i.i.d., we have

$$\begin{aligned} \operatorname{Var}[\psi_{j}|\boldsymbol{\theta}_{\pi^{1},j}^{p},\ldots,\boldsymbol{\theta}_{\pi^{k},j}^{p}] &= \operatorname{Var}\left[\frac{\sum_{t=1}^{k}p_{j}(\pi^{t})[V(\boldsymbol{\theta}_{\pi^{t},j})-V(\boldsymbol{\theta}_{\pi^{t},j}^{p})]}{k}\middle|\boldsymbol{\theta}_{\pi^{1},j}^{p},\ldots,\boldsymbol{\theta}_{\pi^{k},j}^{p}\right] \\ &= \operatorname{Var}\left[\frac{\sum_{t=1}^{k}p_{j}(\pi^{t})V(\boldsymbol{\theta}_{\pi^{t},j})}{k}\middle|\boldsymbol{\bar{\theta}}_{\pi^{1},j},\ldots,\boldsymbol{\bar{\theta}}_{\pi^{k},j}\right] \\ &\geq \frac{p^{2}}{k^{2}}\sum_{t=1}^{k}\operatorname{Var}[V(\boldsymbol{\theta}_{\pi^{t},j})|\boldsymbol{\bar{\theta}}_{\pi^{t},j}] \\ &= \frac{\sum_{t=1}^{k}\Omega(k^{2})}{k^{2}} \\ &= \Omega(k) \;. \end{aligned}$$

C.4 USEFUL LEMMAS AND COROLLARIES

Lemma C.5. Let $z_t^* := \sum_{l=1}^t \mathbf{X}_{t,l} z_l$ where $z_l \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, k(C\sigma)^2 \mathbf{I}) \in \mathbb{R}^d$ are i.i.d. drawn. Among all matrices satisfying P1 and P2, \mathbf{X}^* as defined in Eq. (4) produces the lowest values of $N := \sum_{t=1}^k \mathbb{E}[||z_t^*||_2^2]$, $P := \sum_{t=1}^k \mathbb{E}[||z_t^*||_2^2]$, and $Q := \sum_{t=1}^k \sqrt{\mathbb{E}[||z_t^*||_2^2]}$. In particular, we have

$$\begin{split} N &= \mathcal{O}(k \log k) \ , \\ P &= \mathcal{O}(k^2) \ , \\ Q &= \mathcal{O}(k \log k) \ . \end{split}$$

Proof. Since $z_t^* \in \mathbb{R}^d$ and is diagonal, we may write $z_t^* = (\xi_{t,1}, \xi_{t,2}, \dots, \xi_{t,d})$, where $\forall t \in [k], r \in \mathcal{R}^d$ $[d], \xi_{t,r} \sim \mathcal{N}(0, k(C\sigma)^2 \sum_{l=1}^t X_{t,l}^2)$ since all z_l 's are independent. Then, for each z_t^* , there is

$$\mathbb{E}[\|z_t^*\|_2^2] = \mathbb{E}[\sum_{r=1}^d \xi_{t,r}^2] \\ = \mathbb{E}[\sum_{r=1}^d (k(C\sigma)^2 \sum_{l=1}^t X_{t,l}^2)\chi_1^2]$$

$$= \mathbb{E}[(k(C\sigma)^{2} \sum_{l=1}^{t} X_{t,l}^{2})\chi_{d}^{2}]$$

$$= dk(C\sigma)^{2} \sum_{l=1}^{t} X_{t,l}^{2}$$

where χ_d^2 refers to a chi-squared random variable with a degree of freedom d which satisfies $\mathbb{E}[\chi_d^2] = d$. Since N is the sum of $\mathbb{E}[||z_t^*||_2^2]$ over all $t \in [k]$, we can minimize N if $\exists \mathbf{X}^*$, s.t. $\forall t \in [k], \mathbf{X}^* = \operatorname{argmin}_{\mathbf{X}} \mathbb{E}[||z_t^*||_2^2]$. Fortunately, the optimization problem is a con-vex quadratic program. Specifically, let $w_t = (X_{t,1}, X_{t,2}, ..., X_{t,t})^{\top}$, we have the following optimisation problem for each $t \in [k]$:

$$\begin{array}{ll} \min_{\boldsymbol{w}_t} & \boldsymbol{w}_t^{\top} \boldsymbol{w}_t \\ \text{s.t.} & \mathbf{1}^{\top} \boldsymbol{w}_t = 1 \\ & \boldsymbol{w}_t \ge \mathbf{0} \ . \end{array} \tag{5}$$

To solve this optimization problem, we adopt the standard Lagrange multiplier approach. The Lagrangian is

 $L(\boldsymbol{w}_t, \lambda) = \boldsymbol{w}_t^{\top} \boldsymbol{w}_t - \lambda (1 - \mathbf{1}^{\top} \boldsymbol{w}_t).$

Then, we wish to have

 $\nabla_{w_i} L(\boldsymbol{w}_t, \lambda) = 0$ $\forall l \in [t], 2\boldsymbol{X}_{t,l} + \lambda = 0.$

This suggests that the optimal is obtained when $X_{t,1} = X_{t,2} = ... = X_{t,t} = -\lambda/2 = 1/t$ for each $t \in [k]$, which matches the form of X^* in Eq. 4. With X^* , $\xi_{t,r} \sim \mathcal{N}(0, k(C\sigma)^2/t)$, thus we can

1404 derive the lowest value of N as 1405

1406 1407

1410 1411

1418

where the second last step is by the inequality $\sum_{t=1}^{k} 1/t < 1 + \log k$. Next, we show that X^* is also a minimizer of P. Consider each summand of P can be expressed as 1420

 $= \mathcal{O}(k \log k)$

 $\min N = \sum_{k=1}^{k} \mathbb{E}[\|z_t^*\|_2^2]$

 $=\sum_{i=1}^{k}\mathbb{E}[\sum_{i=1}^{d}\xi_{t,r}^{2}]$

 $=\sum_{k=1}^{k}\frac{dk(C\sigma)^2}{t}$

 $=\sum_{k=1}^{k}\frac{k(C\sigma)^{2}}{t}\mathbb{E}[\chi_{d}^{2}]$

 $\leq dk (C\sigma)^2 (1 + \log k)$

1426 1427

1425
1426
1427
1428
1429
1430
1431

$$= \mathbb{E}[\sum_{r=1}^{d} (\xi_{t,r}^{2})^{2}] + \mathbb{E}[\sum_{r \neq r'} \xi_{t,r}^{2} \xi_{t,r'}^{2}]$$

$$= \sum_{r=1}^{d} k^{2} (C\sigma)^{4} (\sum_{l=1}^{t} X_{t,l}^{2})^{2} \mathbb{E}[(\chi_{1}^{2})^{2}] + d(d-1) \mathbb{E}[\xi_{t,1}^{2}]^{2}$$

 $\mathbb{E}[\|z_t^*\|_2^4] = \mathbb{E}[(\sum_{r=1}^u \xi_{t,r}^2)^2]$

1432
1433
1434
$$= \sum_{r=1}^{3} k^2 (C\sigma)^4 (\sum_{l=1}^{3} X_{t,l}^2)^2 \mathbb{E}[(\chi_1^2)^2] + d(d-1)(k(C\sigma)^2 (\sum_{l=1}^{3} X_{t,l}^2))^2 \mathbb{E}[\chi_1^2]^2$$

14

$$= (3d + d(d-1))k^2 (C\sigma)^4 (\sum_{l=1} X_{t,l}^2)^2$$

1439

$$= d(d+2)k^2(C\sigma)^4(\sum_{l=1}^{l} X_{t,l}^2)^2$$

1440 where in the 3rd step we use the fact $\mathbb{E}[\xi_{t,r}^2 \xi_{t,r'}^2] = \mathbb{E}[\xi_{t,r}^2] \mathbb{E}[\xi_{t,r'}^2] = \mathbb{E}[\xi_{t,1}^2]^2$ since $\xi_{t,r}$'s are i.i.d. 1441 for all $r \in [d]$. In the 4th step we use the fact that $\mathbb{E}[(\xi_1^2)^2] = 3$. Notice that since $\sum_{l=1}^t X_{t,l}^2 \ge 0$, 1442 minimizing $\mathbb{E}[||z_t^*||_2^4]$ is equivalent to minimizing $\sum_{l=1}^t X_{t,l}^2$. The minimizer is discussed in the 1443 above convex quadratic program where the optimal solution is X^* . As such, we can derive the 1444 minimum of P as 1445 k

1456 where the second last step is by the fact $\sum_{t=1}^{\infty} 1/t^2 = \pi^2/6$. Lastly, we show, in the same vein, that 1457 the minimizer of Q is also X^* . This is obvious since the minimizer of $\mathbb{E}[||z_t^*||_2^4]$ is obviously the 1458 minimizer of $\sqrt{\mathbb{E}[||z_t^*||_2^4]}$. The minimum of Q can hence be derived as

1460
1461
1462
1463
1464
1465
1466

$$\min Q = \sum_{t=1}^{k} \sqrt{\mathbb{E}[||z_t^*||_2^4]}$$

$$= \sum_{t=1}^{k} \sqrt{d(d+2)}k(C\sigma)^2/t$$

$$\leq \sqrt{d(d+2)}k(C\sigma)^2(1+\log k)$$

$$= \mathcal{O}(k \log k)$$
 .

1470 Lemma C.6. Let $z_t^* \coloneqq -\zeta_t + \sum_{l=1}^t X_{t,l}(z_l + \zeta_l)$ where $z_l \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, k(C\sigma)^2 \mathbf{I}) \in \mathbb{R}^d$ and $\zeta_l \stackrel{\text{i.i.d.}}{\sim}$ 1471 $\mathcal{N}(\mathbf{0}, \Sigma) \in \mathbb{R}^d$ are i.i.d. drawn and $\Sigma \in \mathbb{R}^{(d \times d)}$ is a diagonal matrix. Denote $\sigma_g^2 \coloneqq 1/d \sum_{r=1}^d \Sigma_{r,r}$. 1472 The matrix defined as $\forall t \in [k], \forall l \in [t-1], \mathbf{X}_{t,l} = t^{-1}(k(C\sigma)^2 + \sigma_g^2)^{-1}k(C\sigma)^2$ and $\forall t \in [k], \mathbf{X}_{t,t} = t^{-1}(k(C\sigma)^2 + \sigma_g^2)^{-1}(k(C\sigma)^2 + \sigma_g^2)^{-1}(k(C\sigma)^2)^{-1}(k(C\sigma)^2)^{-1}(k(C\sigma)^2)^{-1}(k(C\sigma)^2)^{-1}(k(C\sigma)^2)^{-1}(k(C\sigma)^2)^{-1}(k(C\sigma)^2)^$

1476
$$N = \mathcal{O}(k \log k + k\sigma_g^2),$$

1477
$$P = \mathcal{O}(k^2 + k\sigma_g^4 + \sigma_g^2 k \log k) ,$$
1478

$$Q = \mathcal{O}(k\log k + k\sigma_g^2)$$

Proof. We may write $z_t^* = (\xi_{t,1}, \xi_{t,2}, \dots, \xi_{t,d})$, where $\forall t \in [k], r \in [d], \xi_{t,r} \sim \mathcal{N}(0, k(C\sigma)^2 \sum_{l=1}^t \mathbf{X}_{t,l}^2 + \Sigma_{r,r} \sum_{l=1}^{t-1} \mathbf{X}_{t,l}^2)$, where $\forall t \in [k], r \in [d], \xi_{t,r} \sim \mathcal{N}(0, k(C\sigma)^2 \sum_{l=1}^t \mathbf{X}_{t,l}^2 + (1 - \mathbf{X}_{t,l})^2 \Sigma_{r,r})$ since all the z_l 's and ζ_l 's are independent. We may denote $\sigma_{t,r}^2 \coloneqq k(C\sigma)^2 \sum_{l=1}^t \mathbf{X}_{t,l}^2 + \Sigma_{r,r} \sum_{l=1}^{t-1} \mathbf{X}_{t,l}^2 + (1 - \mathbf{X}_{t,l})^2 \Sigma_{r,r}$. Then, for each z_t^* , there is

$$\mathbb{E}[\|z_t^*\|_2^2] = \mathbb{E}[\sum_{r=1}^d \xi_{t,r}^2]$$
$$= \sum_{r=1}^d \mathbb{E}[\sigma_{t,r}^2\chi_1^2]$$

where in the second step, we use the property that if $\xi \sim \mathcal{N}(0, (C\sigma)^2)$, then $\xi^2/(C\sigma)^2 = \chi_1^2$ where χ_1^2 refers to a chi-squared random variable with degree of freedom 1. Similar to Lemma C.5, let $\boldsymbol{w}_t = (\boldsymbol{X}_{t,1}, \boldsymbol{X}_{t,2}, ..., \boldsymbol{X}_{t,t})^\top$, we have the following optimization problem for each $t \in [k]$:

$$\min_{\substack{t \in \mathcal{I} \\ \text{s.t.}}} \sum_{t=1}^{d} \sigma_{t,r}^{2} \\ \mathbf{w}_{t} = 1 \\ \mathbf{w}_{t} \ge \mathbf{0} .$$
 (6)

To solve this optimization problem, we adopt the standard Lagrange multiplier approach. TheLagrangian is

$$L(\boldsymbol{w}_t, \lambda) = \sum_{r=1}^{d} \sigma_{t,r}^2 - \lambda (1 - \mathbf{1}^{\top} \boldsymbol{w}_t)$$

Then, we wish to have

$$\forall l \in [t], \nabla_{\boldsymbol{X}_{t,l}} L(\boldsymbol{w}_t, \lambda) = 0$$
,

1507 which leads to

1508
1509
1510
$$\forall l \in [t-1], 2(dk(C\sigma)^2 + \sum_{r=1}^{d} \Sigma_{r,r}) \boldsymbol{X}_{t,l} + \lambda = 0$$
1510
$$d = d$$

$$2(dk(C\sigma)^2 + \sum_{r=1}^{a} \Sigma_{r,r}) \boldsymbol{X}_{t,t} - 2\sum_{r=1}^{a} \Sigma_{r,r} + \lambda = 0.$$

1512 Solving this, we get the optimal solution as

1514
1515
$$\forall l \in [t-1], \mathbf{X}_{t,l} = \frac{dk(C\sigma)^2}{t(dk(C\sigma)^2 + \sum_{r=1}^d \Sigma_{r,r})} = \frac{k(C\sigma)^2}{t(k(C\sigma)^2 + \sigma_g^2)}$$

1515
1516
and
$$\mathbf{X}_{t,t} = 1 - \sum_{r=1}^{t-1} \mathbf{X}_{t,r} = \frac{dk(C\sigma)^2 + t \sum_{r=1}^d \Sigma_{r,r}}{dk(C\sigma)^2 + t \sum_{r=1}^d \Sigma_{r,r}} = \frac{k(C\sigma)^2 + t\sigma_g^2}{dk(C\sigma)^2 + t\sigma_g^2}$$

and $\mathbf{X}_{t,t} = 1 - \sum_{l=1}^{l} \mathbf{X}_{t,l} = \frac{1}{t(dk(C\sigma)^2 + \sum_{r=1}^{l} \Sigma_{r,r})} = \frac{1}{t(k(C\sigma)^2 + \sigma_g^2)}$. 1518 As compared to Lemma C.5, it is difficult to find a minimizer for h

1518 As compared to Lemma C.5, it is difficult to find a minimizer for P and Q with the Lagrange 1519 Multiplier approach as it involves a complicated quadratic equation about each $X_{t,l}$. Nevertheless, 1520 we show that with the current minimizer we have, we can derive a good enough upper bound for 1521 N, P, Q. First substitute the values of $X_{t,l}$ and we can express $\sigma_{t,r}^2$ in terms of $(C\sigma)^2$ and σ_g^2 as

1527 Now consider the values of N, P, Q with X:

$$N = \sum_{t=1}^{k} \mathbb{E}[\|z_t^*\|_2^2]$$

$$N = \sum_{t=1}^{k} \mathbb{E}[\|z_t^*\|_2^2]$$

$$= \sum_{t=1}^{k} \mathbb{E}[\sum_{r=1}^{d} \xi_{t,r}^2]$$

$$= \sum_{t=1}^{k} \sum_{r=1}^{d} \sigma_{t,r}^2 \mathbb{E}[\chi_1^2]$$

$$\leq \sum_{t=1}^{k} \left(\frac{3dk(C\sigma)^2}{t} + \frac{kd(C\sigma)^2\sigma_g^4}{(k(C\sigma)^2 + \sigma_g^2)^2} + d^2\sigma_g^2\right)$$

$$\leq 3dk(C\sigma)^2(1 + \log k) + \frac{k^2d(C\sigma)^2\sigma_g^4}{(k(C\sigma)^2 + \sigma_g^2)^2} + kd^2\sigma_g^2$$

$$= \mathcal{O}(k\log k + k\sigma_g^2)$$

where in the second last step we use the fact $\sum_{t=1}^{k} 1/t \le 1 + \log k$.

$$\begin{array}{ll} 1546 & P = \sum_{t=1}^{\kappa} \mathbb{E}[\|z_t^*\|_2^4] \\ 1547 & = \sum_{t=1}^{k} d(d+2)\sigma_{t,r}^4 \\ 1559 & = \sum_{t=1}^{k} d(d+2)\sigma_{t,r}^4 \\ 1550 & \leq d(d+2)\sum_{t=1}^{k} \left(\frac{3dk(C\sigma)^2}{t} + \frac{kd(C\sigma)^2\sigma_g^4}{(k(C\sigma)^2 + \sigma_g^2)^2} + d^2\sigma_g^2\right)^2 \\ 1553 & \leq d(d+2)\sum_{t=1}^{k} \left(\frac{9d^2k^2(C\sigma)^4}{t^2} + \left(\frac{kd(C\sigma)^2\sigma_g^4}{(k(C\sigma)^2 + \sigma_g^2)^2}\right)^2 + d^4\sigma_g^4 + \frac{6dk(C\sigma)^2}{t}\frac{kd(C\sigma)^2\sigma_g^4}{(k(C\sigma)^2 + \sigma_g^2)^2} \\ 1556 & + \frac{6d^3k(C\sigma)^2\sigma_g^2}{t} + \frac{2kd(C\sigma)^2\sigma_g^4}{(k(C\sigma)^2 + \sigma_g^2)^2}d^2\sigma_g^2\right) \\ 1561 & \leq d(d+2)\left(\frac{3d^2k^2(C\sigma)^4\pi^2}{2} + \frac{k^3d^2(C\sigma)^4\sigma_g^8}{(k(C\sigma)^2 + \sigma_g^2)^4} + kd^4\sigma_g^4 + 6d^2(C\sigma)^2(1 + \log k)\frac{k^2d(C\sigma)^2\sigma_g^4}{(k(C\sigma)^2 + \sigma_g^2)^2} \\ 1562 & + 6d^3k(C\sigma)^2\sigma_g^2(1 + \log k) + \frac{2k^2d^3(C\sigma)^2\sigma_g^6}{(k(C\sigma)^2 + \sigma_g^2)^2}\right) \\ 1565 & = \mathcal{O}(k^2 + k\sigma_q^4 + \sigma_q^2k\log k) \end{array}$$

where in the second last step we use the fact $\sum_{t=1}^{k} 1/t^2 \le \sum_{t=1}^{\infty} 1/t^2 = \pi^2/6$. $Q = \sum_{t=1}^{\kappa} \sqrt{\mathbb{E}[\|z_t^*\|_2^4]}$ $=\sqrt{d(d+2)}\sum_{t=1}^{k}\sigma_{t,r}^{2}$ $= \mathcal{O}(N)$ $= \mathcal{O}(k\log k + k\sigma_a^2) \; .$

1581 Corollary C.7. Let $z_t^* \coloneqq -\zeta_t + \sum_{l=1}^t X_{t,l}(z_l + \zeta_l)$ where $z_l \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, k(C\sigma)^2 I) \in \mathbb{R}^d$ and $\zeta_l \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \Sigma) \in \mathbb{R}^d$ are i.i.d. drawn and $\Sigma \in \mathbb{R}^{(d \times d)}$ is a diagonal matrix. Denote $\sigma_g^2 \coloneqq$ $1/d \sum_{r=1}^d \Sigma_{r,r}$. The matrix defined as $\forall t \in [k], \forall l \in [t-1], X_{t,l} = t^{-1}(k(C\sigma)^2 + \sigma_g^2)^{-1}k(C\sigma)^2$ 1585 and $\forall t \in [k], X_{t,t} = t^{-1}(k(C\sigma)^2 + \sigma_g^2)^{-1}(k(C\sigma)^2 + t\sigma_g^2)$ satisfies P1 and P2. Denote $N \coloneqq$ $\sum_{t=kq+1}^k \mathbb{E}[\|z_t^*\|_2^2], P \coloneqq \sum_{t=kq+1}^k \mathbb{E}[\|z_t^*\|_2^4]$, and $Q \coloneqq \sum_{t=kq+1}^k \sqrt{\mathbb{E}[\|z_t^*\|_2^4]}$. With X, we have N, P, Q are upper bounded by

$$N = \mathcal{O}(k \log \frac{1}{q} + k(1-q)\sigma_g^2) ,$$

1591
1592
$$P = \mathcal{O}(k^2 + (1-q)k\sigma_g^4 + \sigma_g^2 k \log \frac{1}{q})$$

$$Q = \mathcal{O}(k \log \frac{1}{q} + k(1-q)\sigma_g^2) .$$

Proof. First, note that the optimal matrix X is not related to where the summation starts. Therefore, 1599 the minimizer for N found in Lemma C.6 is also a minimizer for N in this lemma. Next, we find the 1600 upper bounds for N, P, Q using X:

$$=\sum_{t=kq+1}^{k}\sum_{r=1}^{d}\sigma_{t,r}^{2}\mathbb{E}[\chi_{1}^{2}]$$

$$^{k}_{2}dk(C\sigma)^{2}$$

 $N = \sum_{t=k,a+1}^{k} \mathbb{E}[\|z_t^*\|_2^2]$

 $=\sum_{k=1}^{k} \mathbb{E}[\sum_{i=1}^{d} \xi_{t,r}^{2}]$

 $\leq \sum_{t=ka+1}^{k} \frac{3dk(C\sigma)^{2}}{t} + \frac{kd(C\sigma)^{2}\sigma_{g}^{4}}{(k(C\sigma)^{2} + \sigma_{q}^{2})^{2}} + d^{2}\sigma_{g}^{2}$

$$\leq 3dk(C\sigma)^{2}(1+\log k - \log kq) + k(1-q)\frac{kd(C\sigma)^{2}\sigma_{g}^{4}}{(k(C\sigma)^{2} + \sigma_{g}^{2})^{2}} + k(1-q)d^{2}\sigma_{g}^{2}$$

1616
1617
$$= \mathcal{O}(k \log \frac{k}{kq} + k(1-q)\sigma_g^2)$$

1618
1619
$$= \mathcal{O}(k\log\frac{1}{q} + k(1-q)\sigma_g^2)$$

where in the 3rd last step we make use of the fact $\sum_{t=kq+1}^{k} 1/t = \sum_{t=1}^{k} 1/t - \sum_{t=1}^{kq} 1 + \log K - \log K$ $\log kq$.

$$=\sum_{k=1}^{k}$$

$$= \sum_{t=kq+1} d(d+2)\sigma_{t,r}^4$$

 $P = \sum_{t=kq+1}^{k} \mathbb{E}[\|z_t^*\|_2^4]$

$$\begin{split} &\leq d(d+2) \sum_{t=kq+1}^{k} \left(\frac{3dk(C\sigma)^{2}}{t} + \frac{kd(C\sigma)^{2}\sigma_{g}^{4}}{(k(C\sigma)^{2} + \sigma_{g}^{2})^{2}} + d^{2}\sigma_{g}^{2} \right)^{2} \\ &= d(d+2) \sum_{t=kq+1}^{k} \left(\frac{9d^{2}k^{2}(C\sigma)^{4}}{t^{2}} + \left(\frac{kd(C\sigma)^{2}\sigma_{g}^{4}}{(k(C\sigma)^{2} + \sigma_{g}^{2})^{2}} \right)^{2} + d^{4}\sigma_{g}^{4} \\ &+ \frac{6dk(C\sigma)^{2}}{t} \frac{kd(C\sigma)^{2}\sigma_{g}^{4}}{(k(C\sigma)^{2} + \sigma_{g}^{2})^{2}} + \frac{6d^{3}k(C\sigma)^{2}\sigma_{g}^{2}}{t} + \frac{2kd(C\sigma)^{2}\sigma_{g}^{4}}{(k(C\sigma)^{2} + \sigma_{g}^{2})^{2}} d^{2}\sigma_{g}^{2} \right) \\ &\leq d(d+2) \left(\frac{3d^{2}k^{2}(C\sigma)^{4}\pi^{2}}{2} + \frac{k^{3}d^{2}(C\sigma)^{4}\sigma_{g}^{8}}{(k(C\sigma)^{2} + \sigma_{g}^{2})^{4}} + (1-q)kd^{4}\sigma_{g}^{4} \\ &+ 6d^{2}(C\sigma)^{2}(1 + \log k - \log kq) \frac{k^{2}d(C\sigma)^{2}\sigma_{g}^{4}}{(k(C\sigma)^{2} + \sigma_{g}^{2})^{2}} + 6d^{3}k(C\sigma)^{2}\sigma_{g}^{2}(1 + \log k - \log kq) + \frac{2k^{2}d^{3}(C\sigma)^{2}\sigma_{g}^{6}}{(k(C\sigma)^{2} + \sigma_{g}^{2})^{2}} \right) \end{split}$$

where in the second last step we use the fact $\sum_{t=kq+1}^{k} 1/t^2 \le \sum_{t=1}^{\infty} 1/t^2 = \pi^2/6$.

1651

$$Q = \sum_{t=tq+1}^{k} \sqrt{\mathbb{E}[\|z_t^*\|_2^4]}$$

 1653
 $= \sqrt{d(d+2)} \sum_{t=tq+1}^{k} \sigma_{t,r}^2$

 1656
 $= \mathcal{O}(N)$

 1658
 $= \mathcal{O}(k \log \frac{1}{q} + k(1-q)\sigma_g^2)$.

 1660
 1661

 1662
 1663

Lemma C.8. Let $z_t^* := \sum_{l=1}^t \mathbf{X}_{t,l} z_l$ where $z_l \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, k(C\sigma)^2 \mathbf{I}) \in \mathbb{R}^d$ are i.i.d. drawn. Further let N, P, Q be defined equivalently as Lemma C.5. Then, we have $\mathbb{E}[(\sum_{t=1}^k ||z_t^*||_2^2)^2] \leq P + N^2 + Q^2$.

Proof. For this proof we make use of two facts, namely, for two random variables X, Y, there is $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y] - \operatorname{Cov}[X,Y]$ and $\operatorname{Cov}[X,Y] \leq \sqrt{\operatorname{Var}[X]\operatorname{Var}[Y]}$. Also notice that for any random variable X, $Var[X] \leq \mathbb{E}[X^2]$. With these, we have

$$\begin{split} & \mathbb{E}[(\sum_{t=1}^{k} \|z_{t}^{*}\|_{2}^{2})^{2}] = \sum_{t=1}^{k} \mathbb{E}[\|z_{t}^{*}\|_{2}^{4}] + \sum_{t=1}^{k} \sum_{t'=1}^{k} \mathbb{E}[\|z_{t}^{*}\|_{2}^{2}\|z_{t'}^{*}\|_{2}^{2}] \\ & = \sum_{t=1}^{k} \mathbb{E}[\|z_{t}^{*}\|_{2}^{4}] + \sum_{t=1}^{k} \sum_{t'=1}^{k} \left(\mathbb{E}[\|z_{t}^{*}\|_{2}^{2}]\mathbb{E}[\|z_{t'}^{*}\|_{2}^{2}] + \operatorname{Cov}[\|z_{t}^{*}\|_{2}^{2}, \|z_{t'}^{*}\|_{2}^{2}] \right) \\ & = \sum_{t=1}^{k} \mathbb{E}[\|z_{t}^{*}\|_{2}^{4}] + \sum_{t=1}^{k} \sum_{t'=1}^{k} \mathbb{E}[\|z_{t}^{*}\|_{2}^{2}] \mathbb{E}[\|z_{t'}^{*}\|_{2}^{2}] + \sum_{t=1}^{k} \sum_{t'=1}^{k} \operatorname{Cov}[\|z_{t}^{*}\|_{2}^{2}, \|z_{t'}^{*}\|_{2}^{2}] \\ & = \sum_{t=1}^{k} \mathbb{E}[\|z_{t}^{*}\|_{2}^{4}] + \sum_{t=1}^{k} \sum_{t'=1}^{k} \mathbb{E}[\|z_{t}^{*}\|_{2}^{2}] \mathbb{E}[\|z_{t'}^{*}\|_{2}^{2}] + \sum_{t=1}^{k} \sum_{t'=1}^{k} \operatorname{Cov}[\|z_{t}^{*}\|_{2}^{2}, \|z_{t'}^{*}\|_{2}^{2}] \\ & \leq \sum_{t=1}^{k} \mathbb{E}[\|z_{t}^{*}\|_{2}^{4}] + \sum_{t=1}^{k} \sum_{t'=1}^{k} \mathbb{E}[\|z_{t}^{*}\|_{2}^{2}] \mathbb{E}[\|z_{t'}^{*}\|_{2}^{2}] + \sum_{t=1}^{k} \sqrt{\operatorname{Var}[\|z_{t}^{*}\|_{2}^{2}]} \sum_{t'=1}^{k} \sqrt{\operatorname{Var}[\|z_{t'}^{*}\|_{2}^{2}]} \\ & \leq \sum_{t=1}^{k} \mathbb{E}[\|z_{t}^{*}\|_{2}^{4}] + \sum_{t=1}^{k} \mathbb{E}[\|z_{t}^{*}\|_{2}^{2}] \sum_{t'=1}^{k} \mathbb{E}[\|z_{t'}^{*}\|_{2}^{2}] + \sum_{t=1}^{k} \sqrt{\mathbb{E}}[\|z_{t}^{*}\|_{2}^{4}] \sum_{t'=1}^{k} \sqrt{\mathbb{E}}[\|z_{t''}^{*}\|_{2}^{4}] \\ & \leq \sum_{t=1}^{k} \mathbb{E}[\|z_{t}^{*}\|_{2}^{4}] + \sum_{t=1}^{k} \mathbb{E}[\|z_{t}^{*}\|_{2}^{2}] \sum_{t'=1}^{k} \mathbb{E}[\|z_{t''}^{*}\|_{2}^{2}] + \sum_{t=1}^{k} \sqrt{\mathbb{E}}[\|z_{t}^{*}\|_{2}^{4}] \sum_{t'=1}^{k} \sqrt{\mathbb{E}}[\|z_{t''}^{*}\|_{2}^{4}] \\ & = P + N^{2} + Q^{2} . \\ \\ & = P + N^{2} + Q^{2} . \\ \end{array}$$

1694 C.5 A HYPOTHETICAL CASE

1696 **Proposition C.9.** $\forall t \in [k]$ and, denote $\boldsymbol{\theta}_{\pi^t}^* \coloneqq \boldsymbol{\theta}_{\pi^t}^p - \alpha \sum_{l=1}^t \boldsymbol{X}_{t,l} \tilde{g}_{\pi^l}$. Assume that $\hat{g}_{\pi^1} = \hat{g}_{\pi^2} = \dots = \hat{g}_{\pi^k}$. Denote $m^*(\pi^t) \coloneqq V(\boldsymbol{\theta}_{\pi^t}^*) - V(\boldsymbol{\theta}_{\pi^t}^p)$ and $\psi^* \coloneqq k^{-1} \sum_{t=1}^k p(\pi^t) m^*(\pi^t)$. Then \boldsymbol{X}^* defined in Eq. (4) achieves an estimation uncertainty $\operatorname{Var}[\psi^*|\boldsymbol{\theta}_{\pi^1}^p, \boldsymbol{\theta}_{\pi^2}^p, \dots, \boldsymbol{\theta}_{\pi^k}^p] = (\log^2 k)$.

1700 1701 *Proof.* Denote $\bar{\theta}_{\pi^t,j} \coloneqq \theta_{\pi^t,j}^p - \alpha \hat{g}_{\pi^t,j}$. As explained in the proof for Prop. 5.1, notice that $\bar{\theta}_{\pi^t,j}$ is 1702 deterministic conditional on $\theta_{\pi^t,j}^p$. Then, with X^* ,

 $=(\boldsymbol{\theta}_{\pi^{t},j}^{p}-\alpha\sum_{l=1}^{t}\boldsymbol{X}_{t,l}^{*}\hat{g}_{\pi^{l},j})-\alpha\sum_{l=1}^{t}\boldsymbol{X}_{t,l}^{*}\alpha z_{l}$

 $\boldsymbol{\theta}_{\pi^{t},j}^{*} = \boldsymbol{\theta}_{\pi^{t},j}^{p} - \alpha \sum_{l=1}^{t} \boldsymbol{X}_{t,l}^{*} \tilde{g}_{\pi^{l},j}$

1703 1704

1693

1705

1706 1707

1708

1709 1710

1713

1714 1715

1717

1727

1718 Since $\forall t \in [k], z_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, k(C\sigma)^2 \mathbf{I})$, we can denote $z_t^* \coloneqq \frac{\sum_{l=1}^t z_l}{t} \sim \mathcal{N}(\mathbf{0}, \frac{k}{t}(C\sigma)^2 \mathbf{I})$. Then, 1719 $\boldsymbol{\theta}_{\pi^t,j}^* = \boldsymbol{\theta}_{\pi^t,j} - \alpha z_t^*$.

 $= (\boldsymbol{\theta}_{\pi^{t},j}^{p} - \alpha \hat{g}_{\pi^{t},j}) - \alpha \sum_{l=1}^{t} \boldsymbol{X}_{t,l}^{*} \alpha z_{l}$

 $= \bar{\boldsymbol{\theta}}_{\pi^t, j} - \alpha \sum_{l=1}^t \boldsymbol{X}_{t, l}^* \alpha z_l$

 $= \bar{\boldsymbol{\theta}}_{\pi^t, j} - \alpha \frac{\sum_{l=1}^t \alpha z_l}{t} \; .$

1721 Denote $\bar{p} \coloneqq \max_{\pi \in \Pi} p(\pi)$. Consider that

1722
1723
1724
1725
Var
$$\left[\frac{\sum_{t=1}^{k} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*})}{k} \middle| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] = \frac{1}{k^{2}} \operatorname{Var} \left[\sum_{t=1}^{k} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right].$$
1725

By definition of smoothness, if f is L-smooth, then $\forall x, y$,

$$f(y) \le f(x) + \nabla f(x)^{\top} (y - x) + \frac{L}{2} ||y - x||_2^2.$$

Since V is smooth, -V is (L-)smooth as well, which gives

 $\operatorname{Var}\left|\sum_{t=1}^{\kappa} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right|$

= Var $\left[\sum_{i=1}^{k} p(\pi^{t}) V(\bar{\theta}_{\pi^{t},j} - \alpha z_{t}^{*}) \middle| \bar{\theta}_{\pi^{1},j}, \dots, \bar{\theta}_{\pi^{k},j} \right]$

 $= \mathbb{E} \left[\left(\sum_{t=1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{*}) \right)^{2} \left| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right| \right]$

 $-\mathbb{E}\left|\sum_{t=1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right|^{2}$

 $\leq \mathbb{E}\left[\left(\sum_{t=1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{*})\right)^{2} \left| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right]\right]$

$$-V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{*}) \leq -V(\bar{\boldsymbol{\theta}}_{\pi^{t},j}) + \nabla_{\boldsymbol{\theta}}V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{\top}\alpha z_{t}^{*} + \frac{L}{2}\|\alpha z_{t}^{*}\|_{2}^{2}$$

Thus, we can bound the variance as

 Denote $M \coloneqq \max_{\pi \in \Pi} -V(\boldsymbol{\theta}_{\pi,j})$ and $D \coloneqq \max_{\pi \in \Pi} \|\nabla_{\boldsymbol{\theta}} V(\boldsymbol{\theta}_{\pi,j})\|_2^2$. With these results, the first part can be expanded as

 $\leq \mathbb{E}\left[\left(p(\pi^t)\sum_{t=1}^k -V(\bar{\boldsymbol{\theta}}_{\pi^t,j}) + \nabla_{\boldsymbol{\theta}}V(\bar{\boldsymbol{\theta}}_{\pi^t,j})^\top \alpha z_t^* + \frac{L}{2} \|\alpha z_t^*\|_2^2\right)^2 \left|\bar{\boldsymbol{\theta}}_{\pi^1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^k,j}\right] \right].$

$$=\bar{p}^{2}(\sum_{t=1}^{k}V(\bar{\theta}_{\pi^{t},j}))^{2}+\bar{p}^{2}\sum_{t=1}^{k}\mathbb{E}[\|\nabla_{\theta}V(\bar{\theta}_{\pi^{t},j})\|_{2}^{2}\|\alpha z_{t}^{*}\|_{2}^{2}]+\frac{\bar{p}^{2}L^{2}}{4}\mathbb{E}[(\sum_{t=1}^{k}\|\alpha z_{t}^{*}\|_{2}^{2})^{2}]-\bar{p}^{2}kL\sum_{t=1}^{k}V(\bar{\theta}_{\pi^{t},j})\mathbb{E}[\|\alpha z_{t}^{*}\|_{2}^{2}].$$

Here, applying Lemma C.5 and let $N = \sum_{t=1}^{k}\mathbb{E}[\|z_{t}^{*}\|_{2}^{2}], P = \sum_{t=1}^{k}\mathbb{E}[\|z_{t}^{*}\|_{2}^{4}], Q =$

For every apprying Lemma C.5 and let $N = \sum_{t=1}^{k} \mathbb{E}[\|z_t\|_2], P = \sum_{t=1}^{k} \mathbb{E}[\|z_t\|_2], Q = \sum_{t=1}^{k} \sqrt{\mathbb{E}[\|z_t^*\|_2^4]}$. Additionally, by Lemma C.8, we have $\mathbb{E}[(\sum_{t=1}^{k} \|\alpha z_t^*\|_2^2)^2] \le P + N^2 + Q^2$. Therefore, the above inequality can be bounded by

1774
1775
$$E_1 \le \bar{p}^2 k^2 M^2 + \bar{p}^2 DN + \frac{\bar{p}^2 L^2}{4} (P + N^2 + Q^2) + \bar{p}^2 k LMN$$
1776
$$-\mathcal{O}(k^2 + N + P + N^2 + Q^2 + kN)$$

1776
1777
$$= \mathcal{O}(k^2 + N + P + N^2 + Q^2 + kN)$$

where the minimum is attained at X^* , with $N = O(k \log k)$, $P = O(k^2)$, $Q = O(k \log k)$ and thus $\mathcal{O}(k^2 + N + P + N^2 + Q^2 + kN) = \mathcal{O}(k^2 \log^2 k)$. So, the variance of the average utility is

1780
1781
$$\operatorname{Var}\left[\frac{\sum_{t=1}^{k} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*})}{k} \middle| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] = \frac{E_{1}}{k^{2}} = \frac{1}{k^{2}} \mathcal{O}(k^{2} \log^{2} k) = \mathcal{O}(\log^{2} k)$$

by the multiplicative property of big- \mathcal{O} notations. As such, for the semivalue estimator, we have

$$\begin{aligned}
& \text{Var}[\psi_{j}^{*}|\theta_{\pi^{1},j}^{p},\theta_{\pi^{2},j}^{p},\ldots,\theta_{\pi^{k},j}^{p}] = \text{Var}\left[\frac{\sum_{t=1}^{k} p(\pi^{t})[V(\theta_{\pi^{t},j}^{*}) - V(\theta_{\pi^{t},j}^{p})]}{k} \middle| \theta_{\pi^{1},j}^{p},\theta_{\pi^{2},j}^{p},\ldots,\theta_{\pi^{k},j}^{p} \right] \\
& = \text{Var}\left[\frac{\sum_{t=1}^{k} p(\pi^{t})V(\theta_{\pi^{t},j}^{*})}{k} \middle| \bar{\theta}_{\pi^{1},j},\ldots,\bar{\theta}_{\pi^{k},j} \right] \\
& = \mathcal{O}(\log^{2} k) .
\end{aligned}$$

1793 C.6 PROOF OF PROP. 5.3.

1809 1810

1811 1812

1815 1816 1817

1835

Proposition C.10. (Reproduced from Prop. 5.3, formal) Let $\tilde{g}_{\pi^{l}}, l \in [t]$ be perturbed gradients using the Gaussian mechanism that satisfies (ϵ, δ) -DP. $\forall t \in [k]$, denote $\theta_{\pi^{t}}^{*} \coloneqq \theta_{\pi^{t}}^{p} - \alpha \sum_{l=1}^{t} X_{t,l} \tilde{g}_{\pi^{l}},$ $m^{*}(\pi^{t}) \coloneqq V(\theta_{\pi^{t}}^{*}) - V(\theta_{\pi^{t}}^{p})$ and $\psi^{*} \coloneqq k^{-1} \sum_{t=1}^{k} p(\pi^{t}) m^{*}(\pi^{t})$. Assume that $\forall t \in [k], \hat{g}_{\pi^{t}} - \mathbb{E}[\hat{g}_{\pi^{t}}]$ i.i.d. follow an diagonal multivariate sub-Gaussian distribution with covariance $\Sigma \in \mathbb{R}^{(d \times d)}$ and let $\sigma_{g}^{2} \coloneqq 1/d \sum_{r=1}^{d} \Sigma_{r,r}$. Then the matrix satisfying $\forall t \in [k], \forall l \in [t-1], X_{t,l} = t^{-1}(k(C\sigma)^{2} + \sigma_{g}^{2})^{-1}k(C\sigma)^{2}$ and $\forall t \in [k], X_{t,t} = t^{-1}(k(C\sigma)^{2} + \sigma_{g}^{2})^{-1}(k(C\sigma)^{2} + t\sigma_{g}^{2})$ produces an estimation uncertainty $\operatorname{Var}[\psi^{*}|\theta_{\pi^{1}}^{p}, \theta_{\pi^{2}}^{p}, \dots, \theta_{\pi^{k}}^{p}] = \mathcal{O}(\log^{2} k + \sigma_{g}^{4})$ and $\mathbb{E}[\psi^{*} - \psi] = \mathcal{O}(\log k + \sigma_{g}^{2})$ while satisfying (ϵ, δ) -DP. Moreover, as $k \to \infty, X \to X^{*}$.

1804 1805 Proof. Denote $\bar{\theta}_{\pi^t,j} := \theta_{\pi^t,j}^p - \alpha \hat{g}_{\pi^t,j}$. Note that as compared to the case of Prop. C.9, $\hat{g}_{\pi^t,j}$ is now a random variable (coming from an diagonal sub-Gaussian distribution). We may denote the mean 1806 of the distribution as μ_g . Then $\hat{g}_{\pi^t,j} = \mu_g + \zeta'_t$ where $\zeta'_t \stackrel{\text{i.i.d.}}{\sim}$ sub-Gaussian $(\mathbf{0}, \Sigma)$ with Σ a diagonal 1808 covariance matrix with $\sigma_g^2 = \max_{r \in [d]}(\Sigma)_{r,r}$. As such, $\forall t \in [k]$, denote

$$z_t'^* \coloneqq (\bar{\theta}_{\pi^t,j} - \theta_{\pi^t,j}^*) / \alpha = -\hat{g}_{\pi^t,j} + \sum_{l=1}^t X_{t,l} \tilde{g}_{t,l} = -\zeta_t' + \sum_{l=1}^t X_{t,l} (z_l + \zeta_l') .$$

We use the Gaussian distribution to bound the moments of $z_t^{\prime*}$. To do this, first define a Gaussian counterpart of ζ_t^{\prime} , $\zeta_t \sim \mathcal{N}(\mathbf{0}, \Sigma)$. Then denote

$$z_t^* \coloneqq -\zeta_t + \sum_{l=1}^t \boldsymbol{X}_{t,l}(z_l + \zeta_l)$$

By properties of sub-Gaussian distribution, we have that $z_t^{\prime*}$ follows a sub-Gaussian with the same mean and variance as z_t^* . Hence, $\mathbb{E}[\|z_t^{\prime*}\|_2^2] = \mathbb{E}[\|z_t^{\ast}\|_2^2]$ and $\mathbb{E}[\|z_t^{\prime*}\|_2^4] \le \frac{16}{3}\mathbb{E}[\|z_t^{\ast}\|_2^4]$. The second result is a direct consequence of (Rigollet, 2015, Lemma 1.4) and (Balakrishnan, 2016, Section 7.2). Specifically, let the variance of $z_t^{\prime*}$ and z_t^* be σ_z^2 . (Rigollet, 2015, Lemma 1.4) states that $\mathbb{E}[\|z_t^{\prime*}\|_2^4] \le 16\sigma_z^4$ and (Balakrishnan, 2016, Section 7.2) states that $\mathbb{E}[\|z_t^{\ast}\|_2^4] = 3\sigma_z^4$.

1823 Further denote $\bar{p} := \max_{\pi \in \Pi} p(\pi)$. Consider that 1824

$$\begin{aligned} & \operatorname{Var}\left[\frac{\sum_{t=1}^{k} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*})}{k} \middle| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] = \frac{1}{k^{2}} \operatorname{Var}\left[\sum_{t=1}^{k} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\ & = \frac{1}{k^{2}} \operatorname{Var}\left[\sum_{t=1}^{k} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\ & = \frac{1}{k^{2}} \operatorname{Var}\left[\sum_{t=1}^{k} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*} - \alpha z_{t}^{*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\ & = \frac{1}{k^{2}} \operatorname{Var}\left[\sum_{t=1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \end{aligned}$$

1834 By definition of smoothness, if f is L-smooth, then $\forall x, y$,

$$f(y) \le f(x) + \nabla f(x)^{\top} (y - x) + \frac{L}{2} \|y - x\|_2^2$$

Since V is smooth, -V is (L-)smooth as well, which gives $0 \leq -V(\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t^{\prime*}) \leq -V(\bar{\boldsymbol{\theta}}_{\pi^t,j}) + \nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^t,j})^\top \alpha z_t^{\prime*} + \frac{L}{2} \|\alpha z_t^{\prime*}\|_2^2.$ As such, we have, for the bias, we have $\mathbb{E}[\psi^* - \psi] = \mathbb{E}\left[\mathbb{E}\left[\frac{\sum_{t=1}^k p(\pi^t)[V(\boldsymbol{\theta}^*_{\pi^t,j}) - V(\bar{\boldsymbol{\theta}}_{\pi^t,j})]}{k} \middle| \bar{\boldsymbol{\theta}}_{\pi^1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^k,j}\right]\right]$ $\leq \mathbb{E}\left[\mathbb{E}\left[\mathbb{E}\left[\frac{\sum_{t=1}^{k} \nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{\top} \alpha z_{t}^{\prime*} + \frac{L}{2} \|\alpha z_{t}^{\prime*}\|_{2}^{2}}{k} \left| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right]\right]\right]$ $= \mathbb{E} \left[\sum_{i=1}^{k} \frac{\frac{L}{2} \|\alpha z_t^{\prime *}\|_2^2}{k} \right] .$ By Lemma C.6, we have $\mathbb{E}\left[\sum_{t=1}^{k} \|\alpha z_t^{\prime *}\|_2^2\right] = N = \mathcal{O}(k \log k + k\sigma_g^2) \,.$ Therefore, we have the bound on the bias as $\mathbb{E}[\psi^* - \psi] = \mathbb{E}[\sum_{t=1}^k \|\alpha z_t^{**}\|_2^2]/k = \mathcal{O}(\log k + \sigma_q^2)$. Consider that the variance is upper-bounded as $\operatorname{Var}\left|\sum_{i=1}^{\kappa} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime *})\right| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right|$ $= \mathbb{E}\left[\left(\sum_{t=1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*})\right)^{2} \left| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] - \mathbb{E}\left[\sum_{t=1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \left| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right|^{2}\right]$ $\leq \mathbb{E}\left[\left(\sum_{t=1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*})\right)^{2} \left| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right].$

Since V is strictly non-positive, $-V(\bar{\theta}_{\pi^t,j} - \alpha z_t^{\prime*}) \geq 0$. Denote $M \coloneqq \max_{\pi \in \Pi} -V(\theta_{\pi,j})$ and $D \coloneqq \max_{\pi \in \Pi} \| \nabla_{\boldsymbol{\theta}} V(\boldsymbol{\theta}_{\pi,j}) \|_2^2$. We have $E_{1} \leq \mathbb{E} \left| \left(\sum_{t=1}^{k} -p(\pi^{t})V(\bar{\theta}_{\pi^{t},j}) + p(\pi^{t})\nabla_{\theta}V(\bar{\theta}_{\pi^{t},j})^{\top}\alpha z_{t}^{\prime*} + \frac{p(\pi^{t})L}{2} \|\alpha z_{t}^{\prime*}\|_{2}^{2} \right)^{2} \left| \bar{\theta}_{\pi^{1},j}, \dots, \bar{\theta}_{\pi^{k},j} \right|$ $= \left(\sum_{t=1}^{k} -p(\pi^{t})V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})\right)^{2} + \mathbb{E} \left[\left(\sum_{t=1}^{k} p(\pi^{t})\nabla_{\boldsymbol{\theta}}V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{\top}\alpha z_{t}^{\prime *}\right)^{2} \left| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right| \right]$ $+ \mathbb{E}\left[\left(\frac{p(\pi^{t})L}{2}\sum_{t=1}^{k} \|\alpha z_{t}^{\prime*}\|_{2}^{2}\right)^{2}\right] - kL\sum_{t=1}^{k} p(\pi^{t})^{2}V(\bar{\theta}_{\pi^{t},j})\mathbb{E}[\|\alpha z_{t}^{\prime*}\|_{2}^{2}]$ $\leq \bar{p}^{2} (\sum_{i=1}^{k} -V(\bar{\theta}_{\pi^{t},j}))^{2} + \bar{p}^{2} \sum_{i=1}^{k} \mathbb{E}[\|\nabla_{\theta} V(\bar{\theta}_{\pi^{t},j})\|_{2}^{2} \|\alpha z_{t}^{\prime *}\|_{2}^{2} |\bar{\theta}_{\pi^{1},j}, \dots, \bar{\theta}_{\pi^{k},j}]$ $+ \frac{\bar{p}^{2}L^{2}}{4} \mathbb{E} \left[\left(\sum_{i=1}^{k} \|\alpha z_{t}^{\prime*}\|_{2}^{2} \right)^{2} \right] - \bar{p}^{2} kL \sum_{i=1}^{k} V(\bar{\theta}_{\pi^{t},j}) \mathbb{E}[\|\alpha z_{t}^{\prime*}\|_{2}^{2}]$ $\leq \bar{p}^{2} (\sum_{i=1}^{k} - V(\bar{\theta}_{\pi^{t},j}))^{2} + \bar{p}^{2} D \sum_{i=1}^{k} \mathbb{E}[\|\alpha z_{t}^{\prime *}\|_{2}^{2}] + \frac{\bar{p}^{2} L^{2}}{4} \mathbb{E}\left[\left(\sum_{i=1}^{k} \|\alpha z_{t}^{\prime *}\|_{2}^{2}\right)^{2}\right] - \bar{p}^{2} k L \sum_{i=1}^{k} V(\bar{\theta}_{\pi^{t},j}) \mathbb{E}[\|\alpha z_{t}^{\prime *}\|_{2}^{2}] + \frac{\bar{p}^{2} L^{2}}{4} \mathbb{E}\left[\left(\sum_{i=1}^{k} \|\alpha z_{t}^{\prime *}\|_{2}^{2}\right)^{2}\right] - \bar{p}^{2} k L \sum_{i=1}^{k} V(\bar{\theta}_{\pi^{t},j}) \mathbb{E}[\|\alpha z_{t}^{\prime *}\|_{2}^{2}] + \frac{\bar{p}^{2} L^{2}}{4} \mathbb{E}\left[\left(\sum_{i=1}^{k} \|\alpha z_{t}^{\prime *}\|_{2}^{2}\right)^{2}\right] - \bar{p}^{2} k L \sum_{i=1}^{k} V(\bar{\theta}_{\pi^{t},j}) \mathbb{E}[\|\alpha z_{t}^{\prime *}\|_{2}^{2}] + \frac{\bar{p}^{2} L^{2}}{4} \mathbb{E}\left[\left(\sum_{i=1}^{k} \|\alpha z_{t}^{\prime *}\|_{2}^{2}\right)^{2}\right] + \frac{\bar{p}^{2} L^{2}}{4} \mathbb{E}\left[\left(\sum_{i=1}^{k} \|\alpha z_{t}^{\prime *}\|_{2}^{2}\right)^{2}\right] - \bar{p}^{2} k L \sum_{i=1}^{k} V(\bar{\theta}_{\pi^{t},j}) \mathbb{E}[\|\alpha z_{t}^{\prime *}\|_{2}^{2}] + \frac{\bar{p}^{2} L^{2}}{4} \mathbb{E}\left[\left(\sum_{i=1}^{k} \|\alpha z_{t}^{\prime *}\|_{2}^{2}\right)^{2}\right] - \bar{p}^{2} k L \sum_{i=1}^{k} V(\bar{\theta}_{\pi^{t},j}) \mathbb{E}\left[\left(\sum_{i=1}^{k} \|\alpha z_{t}^{\prime *}\|_{2}^{2}\right)^{2}\right] + \frac{\bar{p}^{2} L^{2}}{4} \mathbb{E}\left[\left(\sum_{i=1}^{k} \|\alpha z_{t}^{\prime}\|_{2}^{2}\right)^{2}\right] + \frac{\bar{p}^{2} L^{2}}{4} \mathbb{E}\left[\left(\sum_{i=1}^{k} \|\alpha z_{t}^{$ $\leq \bar{p}^2 k^2 M^2 + \bar{p}^2 (D + kLM) \sum_{k=1}^{k} \mathbb{E}[\|\alpha z_t^{\prime *}\|_2^2] + \frac{\bar{p}^2 L^2}{4} \mathbb{E}[(\sum_{k=1}^{k} \|\alpha z_t^{\prime *}\|_2^2)^2].$

Let $N = \sum_{t=1}^{k} \mathbb{E}[\|z_t^*\|_2^2] \ge \sum_{t=1}^{k} \mathbb{E}[\|z_t'^*\|_2^2], P = \sum_{t=1}^{k} \mathbb{E}[\|z_t^*\|_2^4] \ge \sum_{t=1}^{k} \mathbb{E}[\|z_t'^*\|_2^4], Q = \sum_{t=1}^{k} \sqrt{\mathbb{E}}[\|z_t^*\|_2^4] \ge \sqrt{3/16} \sum_{t=1}^{k} \sqrt{\mathbb{E}}[\|z_t'^*\|_2^4].$ Additionally, by Lemma C.8, we have $\mathbb{E}[(\sum_{t=1}^{k} \|\alpha z_t'^*\|_2^2)^2] \le \frac{16}{3} \mathbb{E}[(\sum_{t=1}^{k} \|\alpha z_t^*\|_2^2)^2] \le \frac{16}{3} (P + N^2 + Q^2).$ By Lemma C.6, we have the matrix satisfying $\forall t \in [k], \forall l \in [t-1], \mathbf{X}_{t,l} = t^{-1}(k(C\sigma)^2 + \sigma_g^2)^{-1}k(C\sigma)^2$ and $\forall t \in [k], \mathbf{X}_{t,t} = t^{-1}(k(C\sigma)^2 + \sigma_g^2)^{-1}(k(C\sigma)^2 + t\sigma_g^2)$ produces the upper bound of N, P, Qwith

$$N = \mathcal{O}(k\sigma_g^2 + k\log k) ,$$
$$P = \mathcal{O}(k^2 + k\sigma_g^4 + \sigma_g^2 k\log k) ,$$

$$Q = \mathcal{O}(k\log k + k\sigma_g^2) \,.$$

As such, we can further simply E_1 as

$$E_1 = \mathcal{O}(k^2 + kN + P + N^2 + Q^2) = \mathcal{O}(k^2 \sigma_g^4 + k^2 \log^2 k) .$$

Therefore, we can bound the variance as

$$\operatorname{Var}\left[\frac{\sum_{t=1}^{k} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*})}{k} \middle| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \leq \frac{1}{k^{2}} E_{1}$$
$$= \frac{1}{k^{2}} \mathcal{O}(k^{2} \sigma_{g}^{4} + k^{2} \log^{2} k) \quad = \mathcal{O}(\log^{2} k + \sigma_{g}^{4})$$

by the multiplicative property of big- \mathcal{O} notation. As such, for the semivalue estimator, we have

1938
1939
$$\operatorname{Var}[\psi_{j}^{*}|\theta_{\pi^{1},j}^{p},\theta_{\pi^{2},j}^{p},\ldots,\theta_{\pi^{k},j}^{p}] = \operatorname{Var}\left[\frac{\sum_{t=1}^{k} p(\pi^{t})[V(\theta_{\pi^{t},j}^{*}) - V(\theta_{\pi^{t},j}^{p})]}{k} \middle| \theta_{\pi^{1},j}^{p},\theta_{\pi^{2},j}^{p},\ldots,\theta_{\pi^{k},j}^{p}\right]$$
1940
1941
1942
$$= \operatorname{Var}\left[\frac{\sum_{t=1}^{k} p(\pi^{t})V(\theta_{\pi^{t},j}^{*})}{k} \middle| \bar{\theta}_{\pi^{1},j},\bar{\theta}_{\pi^{2},j},\ldots,\bar{\theta}_{\pi^{k},j}\right]$$
1943

 $= \mathcal{O}(\log^2 k + \sigma_a^4)$

where the change in conditioning in the 2nd step is because $\bar{\theta}_{\pi^t,j}$ is deterministic conditional on $\theta^p_{\pi^t,j}$, explained in more detail in the proof for Prop. 5.1. As the operations are all applied to \tilde{g} , by the post-processing immunity of DP, the same (ϵ, δ) -DP guarantee holds. With the given X, it is easy to see that as $k \to \infty$, $X \to X^*$.

1949 1950

1951

1952 **Proposition C.11.** (Reproduced from Prop. 5.4, *formal*) Let $\tilde{g}_{\pi^{l}}, l \in [t]$ be perturbed gra-1953 dients using the Gaussian mechanism that satisfies (ϵ, δ) -DP. $\forall t \in \{kq + 1, \dots, k\}$, denote 1954 $\boldsymbol{\theta}_{\pi^t}^* \coloneqq \boldsymbol{\theta}_{\pi^t}^p - \alpha \sum_{l=1}^t \boldsymbol{Y}_{t-kq,l} \tilde{g}_{\pi^l}.$ Denote $m^*(\pi^t) \coloneqq V(\boldsymbol{\theta}_{\pi^t}^*) - V(\boldsymbol{\theta}_{\pi^t}^p)$ and $\psi^* \coloneqq (k-kq)^{-1} \sum_{t=kq+1}^k p(\pi^t) m^*(\pi^t)$ for $q \in (0,1).$ Assume that $\forall t \in [k], \hat{g}_{\pi^t} - \mathbb{E}[\hat{g}_{\pi^t}]$ i.i.d. follow an diagonal multivariate sub-Gaussian distribution with covariance $\Sigma \in \mathbb{R}^{(d \times d)}$ and let $\sigma_g^2 \coloneqq 1/d \sum_{r=1}^d \Sigma_{r,r}.$ 1955 1956 1957 Then the matrix satisfying $\forall t \in [k], \forall l \in [t-1], \mathbf{X}_{t,l} = t^{-1}(k(C\sigma)^2 + \sigma_g^2)^{-1}k(C\sigma)^2$ and $\forall t \in [k], \mathbf{X}_{t,t} = t^{-1}(k(C\sigma)^2 + \sigma_g^2)^{-1}(k(C\sigma)^2 + t\sigma_g^2)$ produces an estimation uncertainty 1958 1959 1960 $\operatorname{Var}[\psi^* | \theta_{\pi^1}^p, \theta_{\pi^2}^p, \dots, \theta_{\pi^k}^p] = \mathcal{O}\left((1 - q)^{-2} \log^2 (1/q) + \sigma_g^{\overline{4}} \right) \text{ and } \mathbb{E}[\psi^* - \psi] = \mathcal{O}((1 - q) \log 1/q + \sigma_g^{\overline{4}})$ 1961 σ_q^2) while satisfying (ϵ, δ) -DP. 1962

1963 1964 1965 1966 1966 1966 1967 1968 Proof. The proof largely follows the proof for Prop. 5.3. Denote $\bar{\theta}_{\pi^t,j} := \theta_{\pi^t,j}^p - \alpha \hat{g}_{\pi^t,j}$. Note that as compared to the case of Prop. C.9, $\hat{g}_{\pi^t,j}$ is now a random variable (coming from an diagonal sub-Gaussian distribution). We may denote the mean of the distribution as μ_g . Then $\hat{g}_{\pi^t,j} = \mu_g + \zeta'_t$ where $\zeta'_t \stackrel{\text{i.i.d.}}{\sim}$ sub-Gaussian $(\mathbf{0}, \Sigma)$ with Σ a diagonal covariance matrix with $\sigma_g^2 = \max_{r \in [d]}(\Sigma)_{r,r}$. As such, $\forall t \in [k]$, denote

 $z_t^{\prime*} := (\bar{\boldsymbol{\theta}}_{\pi^t,j} - \boldsymbol{\theta}_{\pi^t,j}^*) / \alpha = -\hat{g}_{\pi^t,j} + \sum_{l=1}^t \boldsymbol{X}_{t,l} \tilde{g}_{t,l} = -\zeta_t^{\prime} + \sum_{l=1}^t \boldsymbol{X}_{t,l} (z_l + \zeta_l^{\prime}) .$ We use the Gaussian distribution to bound the moments of $z_t^{\prime*}$. To do this, first define a Gaussian

 $z_t^* \coloneqq -\zeta_t + \sum_{l=1}^t \boldsymbol{X}_{t,l}(z_l + \zeta_l) .$

1973 counterpart of ζ'_t , $\zeta_t \sim \mathcal{N}(\mathbf{0}, \Sigma)$. Then denote 1974

1975

1978

1979

1981

1982

1995

By properties of sub-Gaussian distribution, we have that z_t^{**} follows a sub-Gaussian with the same mean and variance as z_t^* . Hence, $\mathbb{E}[||z_t^{**}||_2^2] = \mathbb{E}[||z_t^{**}||_2^2]$ and $\mathbb{E}[||z_t^{**}||_2^4] \le \frac{16}{3}\mathbb{E}[||z_t^{**}||_2^4]$. The second result is a direct consequence of (Rigollet, 2015, Lemma 1.4) and (Balakrishnan, 2016, Section 7.2). Specifically, let the variance of z_t^{**} and z_t^* be σ_z^2 . (Rigollet, 2015, Lemma 1.4) states that $\mathbb{E}[||z_t^{**}||_2^4] \le 16\sigma_z^4$ and (Balakrishnan, 2016, Section 7.2) states that $\mathbb{E}[||z_t^{**}||_2^4] = 3\sigma_z^4$.

1983 Further denote $\bar{p} \coloneqq \max_{\pi \in \Pi} p(\pi)$. Consider that

$$\begin{aligned}
& \text{1984} \\
& \text{1985} \\
& \text{1986} \\
& \text{Var} \left[\frac{\sum_{t=kq+1}^{k} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*})}{k-kq} \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\
& = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\
& = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*} - \alpha z_{t}^{\prime*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\
& = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\
& = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\
& = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\
& = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\
& = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\
& = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\
& = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\
& = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\
& = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\
& = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \middle| \bar{\boldsymbol{\theta}}_{\pi^{k}q+1,j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \\
& = \frac{1}{(k-kq)^{2}} \text{Var} \left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},$$

Since -V is L-smooth, following the proof for Prop. C.9, we have

1996
1997
$$0 \le -V(\bar{\theta}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \le -V(\bar{\theta}_{\pi^{t},j}) + \nabla_{\theta} V(\bar{\theta}_{\pi^{t},j})^{\top} \alpha z_{t}^{\prime*} + \frac{L}{2} \|\alpha z_{t}^{\prime*}\|_{2}^{2}$$
1997



 $\mathbb{E}[\psi^* - \psi] = \mathbb{E}\left[\mathbb{E}\left[\frac{\sum_{t=kq+1}^k p(\pi^t)[V(\boldsymbol{\theta}_{\boldsymbol{\pi}^t, j}^*) - V(\bar{\boldsymbol{\theta}}_{\boldsymbol{\pi}^t, j})]}{k - kq} \middle| \bar{\boldsymbol{\theta}}_{\pi^1, j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^k, j}\right]\right]$ $\leq \mathbb{E}\left[\mathbb{E}\left[\frac{\sum_{t=kq+1}^{k} \nabla_{\boldsymbol{\theta}} V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{\top} \alpha z_{t}^{\prime*} + \frac{L}{2} \|\alpha z_{t}^{\prime*}\|_{2}^{2}}{k-kq} \left| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right| \right]\right]$ $= \mathbb{E}\left[\sum_{k=1}^{k} \frac{\frac{L}{2} \|\alpha z_t^{\prime *}\|_2^2}{k - kq}\right] .$ By Corollary C.7, we have $\mathbb{E}\left[\sum_{t=1}^{\kappa} \|\alpha z_t^{\prime *}\|_2^2\right] = N = \mathcal{O}(k \log 1/q + k(1-q)\sigma_g^2) .$ Therefore, we have the bound on the bias as $\mathbb{E}[\psi^* - \psi] = \mathbb{E}[\sum_{t=ka+1}^k \|\alpha z_t^{**}\|_2^2]/(k-kq) =$ $\mathcal{O}((1-q)\log 1/q + \sigma_a^2).$ Consider that the variance is upper-bounded as $\operatorname{Var}\left|\sum_{t=t-1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*})\right| \bar{\boldsymbol{\theta}}_{\pi^{kq+1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right|$ $= \mathbb{E}\left| \left(\sum_{t=kq+1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \right)^{2} \left| \bar{\boldsymbol{\theta}}_{\pi^{kq+1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right| - \mathbb{E}\left[\sum_{t=kq+1}^{k} p(\pi^{t}) V(\bar{\boldsymbol{\theta}}_{\pi^{t},j} - \alpha z_{t}^{\prime*}) \left| \bar{\boldsymbol{\theta}}_{\pi^{kq+1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right|^{2} \right] \right|$ $\leq \mathbb{E} \left| \left(\sum_{t=kq+1}^{k} p(\pi^t) V(\bar{\boldsymbol{\theta}}_{\pi^t,j} - \alpha z_t'^*) \right)^2 \left| \bar{\boldsymbol{\theta}}_{\pi^{kq+1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^k,j} \right| \right|.$

Since V is strictly non-positive, $-V(\bar{\theta}_{\pi^t,j} - \alpha z_t'^*) \ge 0$. Denote $M \coloneqq \max_{\pi \in \Pi} -V(\theta_{\pi,j})$ and $D \coloneqq \max_{\pi \in \Pi} \| \nabla_{\boldsymbol{\theta}} V(\boldsymbol{\theta}_{\pi,j}) \|_2^2$. We have $E_{1} \leq \mathbb{E} \left[\left(\sum_{t=t,s+1}^{k} -p(\pi^{t})V(\bar{\theta}_{\pi^{t},j}) + p(\pi^{t})\nabla_{\theta}V(\bar{\theta}_{\pi^{t},j})^{\top}\alpha z_{t}^{\prime*} + \frac{p(\pi^{t})L}{2} \|\alpha z_{t}^{\prime*}\|_{2}^{2} \right)^{2} \left| \bar{\theta}_{\pi^{k}q+1,j}, \dots, \bar{\theta}_{\pi^{k},j} \right| \right]$ $= \left(\sum_{t=k_{d+1}}^{k} - p(\pi^{t})V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})\right)^{2} + \mathbb{E} \left[\left(\sum_{t=k_{d+1}}^{k} p(\pi^{t})\nabla_{\boldsymbol{\theta}}V(\bar{\boldsymbol{\theta}}_{\pi^{t},j})^{\top}\alpha z_{t}^{\prime*} \right)^{2} \left| \bar{\boldsymbol{\theta}}_{\pi^{k_{d+1}},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right| \right]$ $+ \mathbb{E} \left| \left(\frac{p(\pi^t)L}{2} \sum_{t=ka+1}^k \|\alpha z_t'^*\|_2^2 \right)^2 \right| - kL \sum_{t=kq+1}^k p(\pi^t)^2 V(\bar{\theta}_{\pi^t,j}) \mathbb{E}[\|\alpha z_t'^*\|_2^2]$ $\leq \bar{p}^{2} (\sum_{k=1}^{k} -V(\bar{\theta}_{\pi^{t},j}))^{2} + \bar{p}^{2} \sum_{k=1}^{k} \mathbb{E}[\|\nabla_{\theta}V(\bar{\theta}_{\pi^{t},j})\|_{2}^{2} \|\alpha z_{t}^{\prime*}\|_{2}^{2} |\bar{\theta}_{\pi^{k}q+1,j}, \dots, \bar{\theta}_{\pi^{k},j}]$ $+\frac{\bar{p}^{2}L^{2}}{4}\mathbb{E}\left[\left(\sum_{t=ka+1}^{k}\|\alpha z_{t}^{\prime*}\|_{2}^{2}\right)^{2}\right]-\bar{p}^{2}kL\sum_{t=ka+1}^{k}V(\bar{\theta}_{\pi^{t},j})\mathbb{E}[\|\alpha z_{t}^{\prime*}\|_{2}^{2}]$ $\leq \bar{p}^{2} (\sum_{t=t_{\alpha+1}}^{k} -V(\bar{\theta}_{\pi^{t},j}))^{2} + \bar{p}^{2} D \sum_{t=k\alpha+1}^{k} \mathbb{E}[\|\alpha z_{t}'^{*}\|_{2}^{2}] + \frac{\bar{p}^{2} L^{2}}{4} \mathbb{E} \left| \left(\sum_{t=kq+1}^{k} \|\alpha z_{t}'^{*}\|_{2}^{2} \right) \right|$ $-\bar{p}^2kL\sum_{t=1}^{n}V(\bar{\boldsymbol{\theta}}_{\pi^t,j})\mathbb{E}[\|\alpha z_t^{\prime*}\|_2^2]$ $\leq \bar{p}^{2}(k-kq)^{2}M^{2} + \bar{p}^{2}(D+(k-kq)LM) \sum_{k=1}^{k} \mathbb{E}[\|\alpha z_{t}^{\prime*}\|_{2}^{2}] + \frac{\bar{p}^{2}L^{2}}{4}\mathbb{E}[(\sum_{k=1}^{k} \|\alpha z_{t}^{\prime*}\|_{2}^{2})^{2}].$ Let $N = \sum_{t=kq+1}^{k} \mathbb{E}[\|z_t^*\|_2^2] \ge \sum_{t=kq+1}^{k} \mathbb{E}[\|z_t^*\|_2^2], P = \sum_{t=kq+1}^{k} \mathbb{E}[\|z_t^*\|_2^4] \ge \sum_{t=kq+1}^{k} \mathbb{E}[\|z_t^*\|_2^4], Q = \sum_{t=kq+1}^{k} \sqrt{\mathbb{E}[\|z_t^*\|_2^4]} \ge \sqrt{3/16} \sum_{t=kq+1}^{k} \sqrt{\mathbb{E}[\|z_t^*\|_2^4]}.$ Additionally, by Lemma C.8, we have $\mathbb{E}[(\sum_{t=kq+1}^{k} \|\alpha z_t^*\|_2^2)^2] \le \frac{16}{3} \mathbb{E}[(\sum_{t=kq+1}^{k} \|\alpha z_t^*\|_2^2)^2] \le \frac{16}{3} (P + N^2 + Q^2).$ By Corollary C.7, we have the matrix satisfying $\forall t \in [k], \forall l \in [t-1], \mathbf{X}_{t,l} = t^{-1}(k(C\sigma)^2 + \sigma_q^2)^{-1}k(C\sigma)^2$ and $\forall t \in [k], \mathbf{X}_{t,t} = t^{-1}(k(C\sigma)^2 + \sigma_g^2)^{-1}(k(C\sigma)^2 + t\sigma_g^2)$ produces an upper bound of N, P, Q with $N = \mathcal{O}(k \log \frac{1}{q} + k(1-q)\sigma_g^2)$ $P = \mathcal{O}(k^2 + (1-q)k\sigma_g^4 + \sigma_g^2 k \log \frac{1}{q})$

 With these, we can derive the big- \mathcal{O} bound for E_1 as

$$E_1 = \mathcal{O}(k^2(1-q)^2 + k(1-q)N + P + N^2 + Q^2) = \mathcal{O}(k^2 \log^2 \frac{1}{q} + k^2(1-q)^2 \sigma_g^4) .$$

 $Q = \mathcal{O}(k \log \frac{1}{q} + k(1-q)\sigma_g^2) .$

Therefore, we can bound the variance as

$$\operatorname{Var}\left[\frac{\sum_{t=kq+1}^{k} p(\pi^{t}) V(\boldsymbol{\theta}_{\pi^{t},j}^{*})}{k-kq} \middle| \bar{\boldsymbol{\theta}}_{\pi^{1},j}, \dots, \bar{\boldsymbol{\theta}}_{\pi^{k},j} \right] \quad \leq \frac{1}{(k-kq)^{2}} E_{1} = \mathcal{O}\left(\frac{\log^{2} \frac{1}{q}}{(1-q)^{2}} + \sigma_{g}^{4}\right) \;.$$

As such, for the semivalue estimator, we have

where the change in conditioning random variable in the 2nd step is because $\bar{\theta}_{\pi^t,j}$ is deterministic conditional on $\theta^p_{\pi^t,j}$, explained in more detail in the proofs for the above propositions. As all operations are applied to \tilde{g} , by the post-processing immunity of DP, the same (ϵ, δ) -DP guarantee holds.

D EXPERIMENTS

Hardware and software details. We perform all our experiments on Nvidia L40 GPUs using the
PyTorch (Paszke et al., 2019) deep learning framework with the Opacus (Opacus) implementation of
private random variables (PRV) (Gopi et al., 2021) for privacy accounting. Source codes are included
in supplementary materials.

2128 Hyperparameter settings. For experiments using Algorithm 1, we follow (Ghorbani & Zou, 2019) 2129 and apply hyper-parameter search to find a suitable α which produces the best model performance 2130 with one pass of training examples. For a given number of evaluations k and ϵ , Opacus automatically adjusts the noise multiplier σ to achieve the $(\epsilon, 5 \times 10^{-5})$ -DP guarantee. Hence, we focus on the 2131 analysis of the interaction between DP requirements and data valuation by studying ϵ and k. In all 2132 experiments, we choose the value of ϵ such that we can observe a degradation of the performance 2133 of the estimate with i. i. d. noise as compared to the no-DP estimate given the limited budget. This 2134 approach of setting ϵ allows us to highlight the advantage of using correlated noise. 2135

Model specification. We specify the parameterization of the models used in our experiments. We adopt the following notations: ReLU denotes a rectified linear unit activation function; Linear(x, y)denotes a linear layer (i.e. a matrix of dimension x by y); Sigmoid denotes a Sigmoid activation function; Softmax denotes a Softmax activate function; Conv(x, y, z) denotes a convolutional layer with input size x, output size y, and kernel size z (with stride 1 and padding 0); Pool(z, w) denotes a pooling layer with kernel size z and stride w (with padding 0). The neural network models used in our experiments are parameterized as follows:

2143 2144

2145

2121

2122

Logistic Regression(x) := Sigmoid \circ Linear(no. features, no. classes)(x);

 $CNN(x) \coloneqq Softmax \circ Linear(\cdot, no. classes) \circ Pool(2, 2) \circ ReLU \circ Conv(1, 16, 3)(x);$

ResNet18 and ResNet34 are standard. We follow the implementation of these two networks in
 PyTorch's torchvision library.

2148

2149 D.1 MIA ATTACK

2150 We follow the MIA attack proposed in (Wang et al., 2023), which constructs a likelihood ratio test 2151 based on the estimated data values. As a setup, we select 25 data points as the members and 25 data 2152 points as non-members. We further select 200 data points to subsample "shadow dataset". We choose 2153 k = 200 for breast cancer and diabetes datasets and k = 500 for Covertype dataset. As shown in 2154 Table 5, our method offers privacy protection against MIA by reducing the AUROC of the attack to 2155 around 0.5 (same level as random guess).

2156

2158

2157 D.2 GENERALIZING TO OTHER USE CASES

Notes on dataset valuation. For dataset valuation on ResNet18 and ResNet34, we apply some engineering tricks to encourage faster convergence so that we obtain meaningful marginal contributions

Table 5: Mean (std. errors) of AUC on Covertype, breast cancer, and diabetes datasets trained with LR. $\epsilon = 1.0$.

		Covertype	breast cancer	diabetes
r	10 DP	0.554 (1.61e-02)	0.583 (2.57e-02)	0.567 (1.59e-02)
	Ours	0.490 (3.76e-02)	0.477 (2.62e-02)	0.513 (3.68e-02)

in each iteration using the G-Shapley framework. These include 1) freezing part of the network and only fine-tuning the last residual block and the fully connected layer; 2) reinitializing the model with the pre-trained weight after each iteration.

Federated learning setup. (Wang et al., 2020) proposed using an average of Shapley value in multiple rounds of FL as the data value metric

$$\psi_j \coloneqq k^{-1} \sum_{t=1}^k \nu_j^t \, ,$$

where ν_i^t refers to the data Shapley of party j at round t using its gradients released at round t (which we estimate with 100 samples of permutations in each participation round). Note that this formulation is analogous to the general semivalue definition of Eq. (1). As such, we may treat each ν_i^t as a marginal contribution of party j which is averaged over all participation rounds. Note that in a total of k rounds, each party still needs to release the gradients k times, thus incurring the problem of linearly scaling variance of the Gaussian noise which can be alleviated with our method. However, the setting of FL is not identical to the conventional data valuation setting as the global model keeps updating which causes the gradients to change drastically, especially in the first few rounds. To address this challenge, we use a weighted sum that puts more emphasis on the more current gradients in the first few rounds: $X_{t,t} = 0.75 - 0.7 \times t/k$. Note that this choice is consistent with Prop. 5.3 which suggests setting a larger $X_{t,t}$ when σ_g^2 is large. This can be implemented algorithmically by updating \tilde{g}_{π^t} with $\tilde{g}^*_{\pi^t} = (0.25 + 0.7 \times t/k)\tilde{g}^*_{\pi^{t-1}} + (0.75 - 0.7 \times t/k)\tilde{g}_{\pi^t}$. Moreover, as the global model converges with more collaboration rounds, more information about the data values is revealed in the first few rounds. Therefore, setting q too large causes ψ_i to miss important information about the data value, leading to inaccurate data value estimates, even though the estimation uncertainty is lowered. As such, we set a moderately small burn-in ratio q = 0.2.

D.3 ADDITIONAL EXPERIMENTAL RESULTS FOR SEC. 6.1

Additional results for s^2 and μ . We additionally plot a counterpart version of Fig. 2 with the ψ 's evaluated with no DP added in Fig. 5. It can be observed that the variance of ψ 's computed with no DP noise is almost the same as that computed with correlated noise. The means μ 's are also similar.



Figure 5: (Left) $n^{-1} \sum_{j \in [n]} s_j^2 / |\mu_j|$ and (right) μ_j vs. k using i.i.d. noise, correlated noise, and no DP noise.

2241

2242

2243

2249

2251

2252

2253

2214 Additional results for data selections. We supplement the data selection experiment with results 2215 for adding/removing data with high/low ψ 's shown in Fig. 6. For i.i.d. noise, as k increases, the 2216 test accuracy approaches that of random selection, suggesting that the data value estimates are less 2217 reflective of the true worth of data. On the other hand, ψ 's computed with correlated noise exhibit 2218 a test accuracy curve close to ψ 's computed without DP, implying that the data value estimates computed with correlated noise are reflective of the true worth of data. 2219



Figure 6: Data selection task where data with high/low ψ 's are added/removed from the training dataset with ψ 's computed using (top) i.i.d. noise with different k and using (bottom) i.i.d. noise, correlated noise, and no DP noise with k = 1000. y-axis represents test accuracy and x-axis represents percentage of data added/removed.

Results with regression tasks. We notice that the noise due to DP has a much less significant 2248 impact on the data selection task performance for regression models. This may be attributable to the lack of a critical decision boundary which reduces the importance of individual data points. We conduct an experiment with the wine quality dataset (Cortez et al., 2009) where we randomly select 400 training examples and 1000 test examples trained with a linear regression using the average negated mean squared error as V. The results are shown in Fig. 7. It can be seen that the mean absolute errors are almost the same for different k and with different methods.



Figure 7: Data selection task where data with high/low ψ 's are added/removed from the training 2264 dataset with ψ 's computed using (a)(b) i.i.d. noise with different k and (c)(d) i.i.d. noise, correlated 2265 noise, and no DP noise with k = 240. y-axis represents mean absolute error and x-axis represents 2266 percentage of data added/removed. 2267



Figure 8: AUC v.s. $q \in [0, 1)$ using negated loss as V. Lines (shades) represent mean (std).

2281 D.4 ADDITIONAL EXPERIMENTAL RESULTS FOR SEC. 6.2

We include in Fig. 8 a counterpart of AUC plots for the noisy label detection task to that shown in Fig. 3. A similar trend can be observed as in the main text.

We also include the AUC plots for the noisy label detection task with other semivalues in Fig. 9. Similar trends as discussed in the main text can be observed in the figures. Moreover, we provide a result for Leave-one-out (LOO) (Cook, 1977) which is not a regular semivalue but also enjoys the improvement offered by our approach, as demonstrated in Table 6. However, the performance is worse than regular semivalues demonstrated in the main content in Table 2, which is expected since LOO does not take into account the marginal contribution over different subsets, as explained in (Ghorbani & Zou, 2019).

Table 6: Average (standard errors) of AUC on Covertype trained with logistic regression (top) and MNIST trained with CNN (bottom). The best score (except for "no DP" which is the baseline to be approximated) is highlighted. Higher is better.



2317

2319

2268

2270 2271

2272 2273

2274

2277 2278

2279

2282

²³¹⁸ D.5 Additional Experimental Results for Different Values of ϵ .

2320 We follow the discussion in App. B and show experimental results for the noisy label detection task 2321 with different ϵ values: 0.1 and 10 ($\epsilon = 1$ is provided in the main text). The results are shown in App. D.5. It can be observed that for $\epsilon = 10$, using i.i.d. noise can perform on par with using no 2322 DP, e.g., on Covertype dataset with logistic regression using data Shapley and Beta(4, 1). On the 2323 other hand, when $\epsilon = 0.1$, both i.i.d. noise and correlated noise have performance close to random 2324 selection on all ML tasks and semivalues. 2325 k=200, CNN, MNIST k=200, CNN, MNIST k=200, CNN, MNIST k=200, CNN, MNIST 2326 1.0 T 1.0 1.0 1.0 2327 0.8 0.8 0.8 0.8 2328 2329 0.6 0.6 0.6 0.6 2330 0.40.4 $\overrightarrow{0.0}$ $0.4_{0.0}^{11}$ $0.4_{0.0}^{11}$ 03 2331 0 6 0 0 0.3 0.6 0.9 0.3 0.6 0.9 0.3 0.6 0.9 2332 (a) MNIST + CNN Shap-(b) MNIST + (c) MNIST + CNN (d) MNIST + CNN CNN ley $\epsilon = 10$ 2333 Banzhaf $\epsilon = 10$ $Beta(4, 1) \epsilon = 10$ Beta(16, 1) $\epsilon = 10$ 1.0 k=200, CNN, MNIST k=200, CNN, MNIST k=200, CNN, MNIST k=200, CNN, MNIST 2334 1.0-1.0 -1.0 =2335 0.8 0.8 **0.**8 0.8 2336 2337 0.6 0.6 0.6 0.6 2338 $0.4\frac{11}{0.0}$ $0.4\frac{1}{0.0}$ $0.4\frac{1}{0.0}$ $0.4\frac{11}{0.0}$ 2339 0.3 0.6 0.9 0.3 0.6 0.9 0.3 0.6 0.9 0.3 0.6 0.9 2340 (e) MNIST + CNN Shap-(f) MNIST + CNN (g) MNIST + CNN (h) MNIST + CNN 2341 ley $\epsilon = 0.1$ Banzhaf $\epsilon = 0.1$ $Beta(4, 1) \epsilon = 0.1$ Beta(16, 1) $\epsilon = 0.1$ k=200, LR, Covertype k=200, LR, Covertype k=200, LR, Covertype 2342 k=200, LR, Covertype **1.0** 1.01 **1.0** 1.0 2343 0.8 0.8 0.8 0.8 2344 2345 0.6 0.6 0.6 0.6 2346 $0.4^{\perp \downarrow}_{0.0}$ 0.4 $\underbrace{1}_{0.0}$ $0.4^{\perp \perp}_{0.0}$ 0.4 $\stackrel{\square}{0.0}$ 2347 03 0 6 0 9 03 0 6 0 9 03 0.6 0 9 03 0 6 0 9 2348 (j) Covertype + (l) Covertype + (i) Covertype + LR Shap-LR (k) Covertype + LR LR 2349 Banzhaf $\epsilon = 10$ $Beta(4, 1) \epsilon = 10$ Beta(16, 1) $\epsilon = 10$ ley $\epsilon = 10$ 2350 k=200, LR, Covertype k=200, LR, Covertype k=200, LR, Covertype k=200, LR, Covertype 1.0 1.0 1.0 1.0 2351 0.8 0.8 0.8 0.8 2352 2353 0.6 0.6 0.6 0.6 2354 0.4 $\stackrel{\square}{0.0}$ 0.4 $\underbrace{1}_{0.0}$ 2355 0.4 $\stackrel{\square}{0.0}$ 0.40.3 0.6 0.3 0.3 0.3 0.6 0.9 0.9 0.6 0.9 0.6 0.9 2356 (m) Covertype LR (n) Covertype + LR (o) Covertype + LR (p) Covertype + LR + 2357 Shapley $\epsilon = 0.1$ Banzhaf $\epsilon = 0.1$ $Beta(4, 1) \epsilon = 0.1$ $Beta(16, 1) \epsilon = 0.1$





2358

D.6 ADDITIONAL EXPERIMENTAL RESULTS ON RUNTIME AND MEMORY

We assess with experimental results the memory overhead of saving and computing correlated gradients $\tilde{g}_{\pi,i}^*$ for party *i*. Theoretically, the memory overhead is $\mathcal{O}(n)$ for a valuation with *n* parties as each party needs to store its rolling gradient. The runtime overhead is $\mathcal{O}(1)$ as each party only needs to update the rolling gradient each time. The memory overhead results tabulated in Table 7 show a minor (GPU) memory overhead results tabulated in Table 8 demonstrate *no obvious difference* in the computational time with or without correlated noise.

2370 2371

2371

2373

2374

2375

Table 7: GPU peak memory usage with and without (shown in bracket) correlated noise in megabytes. Results are obtained on the diabetes and MNIST datasets with logistic regression (LR) and CNN w.r.t. various number of parties n.

$n \mid$	diabetes+LR	MNIST+CNN
100	18.231 (18.129)	120.788 (109.729)
200	18.329 (18.124)	131.847 (109.729)
300	18.427 (18.120)	142.906 (109.729)

2407Table 8: Program runtime with and without (shown in bracket) correlated noise in seconds. Results2408are obtained on the diabetes (top) and MNIST (bottom) datasets with logistic regression (top) and2409CNN (bottom) w.r.t. various number of parties n and evaluation budget k.

$\begin{array}{c} k \\ n \end{array}$	120	240	360
100	78.466 (78.896)	154.954 (159.989)	262.166 (257.588)
200	113.612 (125.255)	256.431 (252.465)	407.126 (408.107)
300	158.565 (154.425)	344.286 (365.348)	530.579 (519.320)
$\begin{array}{ c c c } k \\ n \end{array}$	120	240	360
$\frac{k}{100}$	120 99.992 (100.140)	240	360 328.204 (323.646)
	120 99.992 (100.140) 158.511 (151.589)	240 199.471 (195.137) 347.525 (329.267)	360 328.204 (323.646) 527.296 (531.594)