
Correlated Noise Provably Beats Independent Noise for Differentially Private Learning

Christopher A. Choquette-Choo* Krishnamurthy (Dj) Dvijotham*
Krishna Pillutla* Arun Ganesh
Thomas Steinke Abhradeep Guha Thakurta

Google

*Equal contribution; alphabetical ordering.

Abstract

Differentially private learning algorithms inject noise into the learning where the most common private learning algorithm, DP-SGD, adds independent Gaussian noise in each iteration. Motivated by the practical considerations in federated learning, recent work on matrix factorization mechanisms has shown empirically that introducing correlations in the noise can greatly improve their utility. We characterize the asymptotic objective suboptimality for any choice of the correlation function, giving precise analytical bounds for linear regression. We show, using these bounds, how correlated noise provably improves upon vanilla DP-SGD as a function of problem parameters such as the effective dimension and condition number. Moreover, our analytical expression for the near-optimal correlation function circumvents the cubic complexity of the semi-definite program used to optimize the noise correlation in prior work. We validate these theoretical results with experiments on private deep learning in both centralized and federated settings. Our work matches or outperforms prior work while being efficient both in terms of computation and memory.

1 Introduction

The broad adoption of deep learning using sensitive data has led to the increasing popularity of rigorous frameworks for privacy preservation, such as differential privacy [20]. The workhorse of private learning, a differentially private variant of stochastic gradient descent called DP-SGD [2, 9, 48], clips per-example gradients to some ℓ_2 norm and adds *independent* Gaussian noise. DP-SGD has been used in a range of applications from learning with medical images [3] to finetuning large language models with $O(100B)$ parameters [22].

A recent line of work instead proposes a family of algorithms called DP-FTRL that add *correlated* Gaussian noise to each clipped gradient [17, 19, 33, 47]. This work is motivated by the fact that DP-SGD and its federated variant DP-FedAvg [41] require amplification via sampling/shuffling to achieve competitive privacy-utility tradeoffs. This in turn requires uniformly random data samples, an assumption that is violated in federated learning due to complicated client availability [32].

DP-FTRL, on the other hand, has been showing to compete with amplified DP-SGD, and has been used for private federated learning at industrial scale [40, 50]. In fact, by solving an expensive semi-definite program (SDP) to find the noise correlations, Choquette-Choo et al. [16] demonstrated *empirically* that DP-FTRL is never worse and often *much better* than DP-SGD in its privacy-utility tradeoff across multiple modalities.

Table 1: Comparison of Noisy-FTRL (and Noisy-SGD), the variants of DP-FTRL (and DP-SGD) without gradient clipping. We give their asymptotic suboptimality for linear regression with Gaussian inputs with covariance \mathbf{H} and noise multiplier $\sigma_{\text{dp}}^2 = \gamma_\infty(\boldsymbol{\beta})^2/(2\rho)$ based on the limiting sensitivity (defined in §2) in terms of the learning rate η , dimension d , the effective dimension $d_{\text{eff}} = \text{Tr}[\mathbf{H}]/\lambda_{\max}\mathbf{H}$, the strong convexity $\mu = \lambda_{\min}(\mathbf{H})$ and the noise variance ρ^{-1} representing the privacy level. We take the gradient norm $G = 1$ (used to scale the noise) and $\lambda_{\max}(\mathbf{H}) = 1$ w.l.o.g. and only show the term depending on ρ . Since $1 \leq d_{\text{eff}} \leq d$, Noisy-FTRL is better than Noisy-SGD at smaller learning rates η or when the effective dimension d_{eff} is small (e.g., when \mathbf{H} is close to low rank).

Algorithm	Asymptotic Suboptimality F_∞	Ratio w/ Lower Bound	Remark
Lower Bound	$\Omega(\eta^2 \rho^{-1} d_{\text{eff}})$	1	for all $\boldsymbol{\beta}$ with finite $\ \boldsymbol{\beta}\ _1$
Noisy-SGD	$\Theta(\eta \rho^{-1} d)$	$\frac{d}{\eta d_{\text{eff}}}$	$\Theta(\cdot)$ denotes matching upper & lower bounds
Noisy-FTRL	$O\left(\eta^2 \rho^{-1} d_{\text{eff}} \log^2 \frac{1}{\eta \mu}\right)$	$\log^2 \frac{1}{\eta \mu}$	The bound is attained for $\boldsymbol{\beta}$ in (7)

However, several questions remain open. Does DP-FTRL **provably improve** over DP-SGD in its expected utility? Further, can we design a more **computationally efficient** procedure to find the noise correlations for DP-FTRL without significantly worsening the privacy-utility tradeoff?

We answer both questions affirmatively by (1) providing a sharp theoretical characterization of the noisy training dynamics of DP-FTRL (summarized in Table 1), and (2) leveraging these analytical tools to circumvent the SDP required in past work. Our experiments on private image classification and language modeling demonstrate that the proposed approach is competitive with the state-of-the-art correlated noise approaches while being significantly more efficient.

2 Problem Setup and Background

Let $\mathcal{D} = \{z_0, \dots, z_{T-1}\}$ be a dataset of T datapoints, where each datapoint is sampled i.i.d. from an underlying distribution \mathbb{P}_{data} . Our learning objective is to minimize:

$$F(\boldsymbol{\theta}) = \mathbb{E}_{z \sim \mathbb{P}_{\text{data}}} [f(\boldsymbol{\theta}; z)] + r(\boldsymbol{\theta}), \quad (1)$$

where $f(\boldsymbol{\theta}; z)$ is the loss incurred by model parameters $\boldsymbol{\theta} \in \mathbb{R}^d$ on a datapoint z , and $r(\cdot)$ is data-independent regularization. We aim to minimize F while satisfying differential privacy with respect to the dataset \mathcal{D} . We assume that F has a unique minimizer denoted $\boldsymbol{\theta}_*$.

We focus on variants of stochastic gradient descent with a batch size of 1 for data arriving in a stream.¹ DP-FTRL with a noise correlation matrix $\mathbf{B} \in \mathbb{R}^{T \times T}$ (which is lower triangular) iterates

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \left(\text{clip}(\nabla f(\boldsymbol{\theta}_t; z_t), G) + \nabla r(\boldsymbol{\theta}_t) + \sum_{\tau=0}^t \mathbf{B}_{t,\tau} \mathbf{w}_\tau \right) \quad (2)$$

for Gaussian noise $\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \sigma_{\text{dp}}^2 G^2 \mathbf{I}_d)$, where $\text{clip}(\cdot, G)$ denotes projection onto an ℓ_2 ball of radius G . We define Noisy-FTRL to be DP-FTRL without clipping. Taking $\mathbf{B} = \mathbf{I}$ as the identity matrix recovers DP-SGD (with clipping) and Noisy-SGD (without clipping), and other choices give rise to alternate algorithms. We restate a result from prior work showing that DP-FTRL is differentially private for any choice of \mathbf{B} , provided the noise multiplier is scaled appropriately.

Theorem 2.1 ([10, 19]). *DP-FTRL (2) satisfies ρ -zero concentrated differential privacy (zCDP) if the $\sigma_{\text{dp}}^2 = \gamma_T^2(\mathbf{B})/(2\rho)$ where $\gamma_T(\mathbf{B}) = \max_{t < T} \|(\mathbf{B}^{-1})_{:,t}\|_2$ is the sensitivity of \mathbf{B}^{-1} .*

Although Noisy-FTRL is not differentially private, it lets us analyze the noise dynamics of DP-FTRL without technicalities associated with clipping. We sharply characterize the asymptotic utility of Noisy-FTRL for linear regression and use this to give bounds for DP-FTRL.

Provable separation between DP-SGD and DP-FTRL: The best-known separation between DP-SGD and DP-FTRL in the literature is due to Kairouz et al. [33]. For G -Lipschitz convex losses, DP-FTRL at a privacy level of ρ -zCDP achieves a suboptimality of $O(Gd^{1/4}/\sqrt{\rho T})$ compared to DP-SGD's $O(Gd^{1/4}/\sqrt{\rho^2 T})$. The only improvement here is in terms of the privacy parameter ρ .

¹We focus on the centralized setting with a batch size of 1 for simplicity. The analysis extends straightforwardly to larger batch sizes.

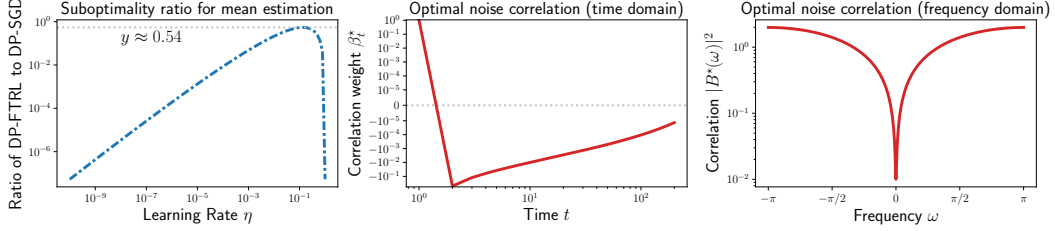


Figure 1: **Left:** The ratio of the asymptotic suboptimality of DP-FTRL to DP-SGD for mean estimation vs. the learning rate η . DP-FTRL is never worse but is orders of magnitude better at $\eta \rightarrow 0$ or 1. **Middle/Right:** Time- and frequency-domain descriptions of the optimal correlations for mean estimation (cf. Theorem 3.1).

This theory fails to reflect the large margin by which DP-FTRL empirically outperforms DP-SGD across the board [16], and a precise characterization is missing.

Computationally efficient DP-FTRL: Prior work on DP-FTRL utilizes the noise correlation matrix B that minimizes the squared error in the *gradient prefix sums* [19, 33]:

$$\varphi(B) = \sum_{t=0}^{T-1} \mathbb{E} \left\| \sum_{\tau=0}^t \tilde{\mathbf{g}}_\tau - \sum_{\tau=0}^t \mathbf{g}_\tau \right\|_2^2 \quad (3)$$

where \mathbf{g}_t is the clipped gradient applied in iteration t and $\tilde{\mathbf{g}}_t$ is its noisy counterpart (cf. Algorithm 1). This was, in turn, obtained as an upper bound on the regret in an adversarial online learning setting [33, Thm. C.1]. The most potent algorithm from the previous work gave B as the solution of a semidefinite program with matrix variables of size $O(T^2)$, requiring $O(T^3)$ time [19, Eq. 4]. This cost is prohibitive for large learning problems. Moreover, there is a mismatch between the objective (3) used to find the noise correlations and the final learning objective $F(\theta_T)$. In particular, Koloskova et al. [35] give two noise correlation matrices B_1, B_2 with equal squared error $\varphi(B_1) = \varphi(B_2)$ such that DP-FTRL with B_1 diverges while DP-FTRL with B_2 converges.

Our approach: We study the suboptimality in the final objective $\mathbb{E}[F(\theta_T) - F(\theta_*)]$. We work in the asymptotic $T \rightarrow \infty$ regime to allow the use of analytic tools, but also to derive results that apply regardless of the dataset size. Second, we restrict the search over B to *Toeplitz* matrices $B_{t,\tau} = \beta_{t-\tau}$ generated by a sequence $\beta = (\beta_0, \beta_1, \dots)$ of reals, but a stronger motivation is that they are **anytime**, i.e., they do not be recomputed for each value of T and easily apply as $T \rightarrow \infty$. Toeplitz B were previously considered for their computational efficiency in learning [17] and for optimal ℓ_2 error (including constants) in linear counting queries [24]. Thus, our goal is to characterize the asymptotic suboptimality

$$F_\infty(\beta) := \lim_{T \rightarrow \infty} \mathbb{E}[F(\theta_T) - F(\theta_*)] \quad (4)$$

for θ_T produced by Noisy-FTRL or DP-FTRL under noise correlation weights β . We analyze this in the frequency domain using the **discrete-time Fourier transform** $B(\omega) = \sum_{t=0}^{\infty} \beta_t \exp(i\omega t)$, with i the imaginary unit. Further, we define the limiting sensitivity associated with B as the limiting value of γ_T , which, using standard Fourier analysis tools, can be expressed as

$$\gamma_\infty(B)^2 := \lim_{T \rightarrow \infty} \gamma_T(B)^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} |B(\omega)|^{-2} d\omega. \quad (5)$$

3 Conceptual Overview: Private Mean Estimation in One Dimension

We begin with the simplest objective function, the squared error for a mean estimation problem on the real line. This setting captures the core intuition and ideas used to derive further results. We emphasize that we do not aim to derive optimal rates for private mean estimation (or later, for private linear regression). Instead, our goal is to demonstrate a separation between DP-FTRL and DP-SGD.

Consider a distribution \mathbb{P}_{data} with $|z - \mathbb{E}[z]| \leq \sigma_{\text{sgd}}$ and $|z| \leq 1$ a.s. for $z \sim \mathbb{P}_{\text{data}}$. Our objective now is to estimate the mean by minimizing the objective

$$F(\theta) = \frac{1}{2} \mathbb{E}_{z \sim \mathbb{P}_{\text{data}}} (\theta - z)^2 \quad \text{with} \quad f(\theta; z) = \frac{z^2}{2} - z\theta, \quad \text{and} \quad r(\theta) = \frac{\theta^2}{2}. \quad (6)$$

We show a strict separation between DP-FTRL and DP-SGD for this simple minimization problem.

Theorem 3.1. Consider the setting above with learning rate $\eta \leq 1$ and clip norm $G = 1$. Then, the asymptotic suboptimality of a ρ -zCDP sequence $(\theta_t)_{t=0}^\infty$ obtained via DP-SGD is $F_\infty(\beta_{\text{dpsgd}}) = \Theta(\eta\rho^{-1} + \eta\sigma_{\text{sgd}}^2)$. Further, the asymptotic suboptimality of any ρ -zCDP sequence $(\theta_t)_{t=0}^\infty$ from DP-FTRL is

$$\inf_{\beta} F_\infty(\beta) = F_\infty(\beta^*) = \Theta\left(\eta^2\rho^{-1}\log^2(1/\eta) + \eta\sigma_{\text{sgd}}^2\right).$$

The infimum above is attained by $\beta_t^* = (-1)^t \binom{1/2}{t} (1-\eta)^t$, where $\binom{1/2}{t} = \prod_{k=0}^{t-1} \frac{1/2-k}{t-k}$.

Proof Sketch. Using tools from frequency domain analysis of linear time invariant systems [43], we can show that the asymptotic variance $\mathbb{E}(\theta_t - \mathbb{E}[z])^2$ is an integral of $|B(\omega)|^2$. The sensitivity is an integral of $|B(\omega)|^{-2}$ (cf. (5)) so that F_∞ is a product of these integrals. Its minimizer B^* can be analytically computed in the Fourier domain (Figure 1, right). An inverse Fourier transform yields the claimed expression for β^* (Figure 1, center). \square

The optimal ρ^{-1} coefficient is $\eta^2\log^2(1/\eta)$ improves over DP-SGD's η (Figure 1, left).

ν -DP-FTRL/ ν -Noisy-FTRL: Theorem 3.1 gives an analytical expression for the optimal noise correlation weights for DP-FTRL for this simplified setting. Using a parameter $0 < \nu < 1$, we define

$$\hat{\beta}_t^\nu := (-1)^t \binom{1/2}{t} (1-\nu)^t. \quad (7)$$

We analyze this choice theoretically for the setting of Noisy-FTRL and demonstrate near optimality. Later, for our experiments with DP-FTRL, we tune ν as a hyperparameter to tune. We call this approach (with clipping) ν -DP-FTRL and (without clipping) ν -Noisy-FTRL.

4 Analysis for Linear Regression

We now give a precise analysis of F_∞ for linear regression with ν -Noisy-FTRL. We consider (unregularized) linear regression with loss function $f(\theta; (\mathbf{x}, y)) = \frac{1}{2}(y - \langle \theta, \mathbf{x} \rangle)^2$

$$F(\theta) = \frac{1}{2} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{\text{data}}} (y - \langle \theta, \mathbf{x} \rangle)^2. \quad (8)$$

We assume d -dimensional Gaussian covariates $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{H})$ and independent Gaussian residuals $y - \langle \theta_*, \mathbf{x} \rangle \sim \mathcal{N}(0, \sigma_{\text{sgd}}^2)$ where $\theta_* = \arg \min F$. We make these assumptions for ease of presentation; we state and prove our results under weaker assumptions in the supplement. Further, we assume that F is L -smooth and μ -strongly convex (equivalently, $\mu\mathbf{I} \preceq \mathbf{H} \preceq L\mathbf{I}$ since the input covariance \mathbf{H} is also the Hessian of F). We express the bounds on F_∞ in terms of the correlation weights β and the problem parameters ρ, G which, for DP-FTRL, denote the target privacy level and the gradient clip norm respectively. See §C for proofs.

Theorem 4.1. Let c, C_1, C_2 denote universal constants. For $\eta \leq c/\text{Tr}[\mathbf{H}]$, we have

$$\begin{aligned} \text{(Noisy-SGD)} \quad & F_\infty(\beta^{\text{sgd}}) = \Theta(\eta d G^2 \rho^{-1} + \eta \sigma_{\text{sgd}}^2 \text{Tr}[\mathbf{H}]) \quad \text{with } \beta^{\text{sgd}} = (1, 0, \dots), \\ \text{(\nu-Noisy-FTRL)} \quad & F_\infty(\hat{\beta}^\nu) \leq C_1 (\eta^2 G^2 \rho^{-1} \log^2 \frac{1}{\nu} + \eta \sigma_{\text{sgd}}^2) \text{Tr}[\mathbf{H}] \quad \text{with } \nu \leq \eta \mu, \text{ and} \\ \text{(Lower bound)} \quad & F_\infty(\beta) \geq C_2 (\eta^2 G^2 \rho^{-1} + \eta \sigma_{\text{sgd}}^2) \text{Tr}[\mathbf{H}] \quad \text{for all } \beta \text{ with } \|\beta\|_1 < \infty. \end{aligned}$$

Observe that our bounds separate the contributions arising from correlated noise (ρ^{-1} term) and those from the inherent noise in the linear model (σ_{sgd}^2 term). We focus on the effect of correlation because the effect of latter noise is the same across all choices of β . We plot this in Figure 2.

Exponential separation between Noisy-SGD and Noisy-FTRL: Noisy-SGD's stationary error depends on the ambient dimension d , while the lower bound depends on the *effective dimension* $d_{\text{eff}} = \text{Tr}[\mathbf{H}] / \|\mathbf{H}\|_2$ of the covariance \mathbf{H} . We have, $d_{\text{eff}} \leq d$ with equality when all the eigenvalues of \mathbf{H} are equal but $d_{\text{eff}} \ll d$ when the eigenvalues of \mathbf{H} decay rapidly or it is nearly low rank. This is true particularly for overparameterized models where the features may be highly correlated resulting in an approximately low-rank covariance. For instance, if the eigenvalues of \mathbf{H} are $(1, 1/d, \dots, 1/d)$, then $d_{\text{eff}} \leq 2$. Then, Noisy-FTRL's error of $O(\eta^2 \rho^{-1} \log^2(d/\eta))$ is exponentially better than Noisy-SGD's $\Theta(\eta \rho^{-1} d)$. A similar advantage also holds when eigenvalues of \mathbf{H} decay

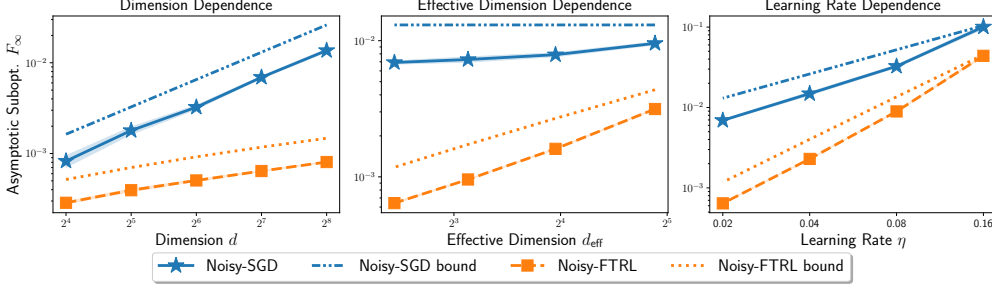


Figure 2: **Linear regression simulations:** We plot the empirically observed asymptotic suboptimality of ν -Noisy-FTRL/Noisy-SGD and their theoretical bounds with $d = 128$ (varied in the left plot) where the Hessian \mathbf{H} has eigenvalues $\lambda_k = 1/k$ (varied as $k^{-\alpha}$ for $\alpha \in [0.4, 1]$ in the middle plot), and learning rate $\eta = 0.02$ (varied in the right plot). The **slope** of the corresponding empirical and theoretical lines are nearly equal, showing the **tightness of the theory**. In particular, we observe that Noisy-SGD has a linear dependence on the dimension (slope 1.00) and is nearly constant w.r.t. the effective dimension (slope 0.18) while Noisy-FTRL has a near-linear dependence on the effective dimension (slope 0.94). Noisy-FTRL (slope 2.03) also has a better dependence on the learning rate than Noisy-SGD (slope 1.27).

at various rates (see Table 2 in §C). The learning rate dependence of Noisy-SGD is also suboptimal, similar to §3. This result is also confirmed empirically in Figure 2 (right).

Finite-time privacy-utility analysis: Noisy-FTRL, which we analyzed so far, is not differentially private. Differential privacy requires gradient clipping that significantly complicates the analysis. However, for a finite time horizon T , we can argue using concentration that $\nabla f(\theta; z)$ is bounded with high probability and clipping can be avoided. Formal statements and proofs for the finite-time analysis are given in Appendix D.

Consider DP-FTRL with noise correlation $\hat{\beta}^\nu$ from (7) with $\nu = \eta\mu$ and gradients clipped to any ℓ_2 -norm G . As mentioned in §2, the outputs $(\theta_1, \dots, \theta_T)$ of DP-FTRL are ρ -zCDP. For an appropriate η , we give utility bounds in terms of the effective dimension d_{eff} and the condition number $\kappa = L/\mu$:

(a) For η small enough, we have with probability at least $1 - p$ that

$$\max_{t < T} \|\mathbf{g}_t\|_2 \leq c \max \left\{ \text{Tr}[\mathbf{H}] \|\theta_0 - \theta_\star\|_2, \sigma_{\text{sgd}} \sqrt{\text{Tr}[\mathbf{H}]} \right\} \text{polylog}(T/p) =: G. \quad (9)$$

Let \mathcal{E} denote this event. If it holds, no gradients are clipped and DP-FTRL coincides with Noisy-FTRL.

(b) For $T \geq \tilde{\Omega}(\kappa^2 d_{\text{eff}}^2 d / \rho)$, we have (omitting log factors and $o(1/T^2)$ terms and assuming $\|\mathbf{H}\|_2 = 1$):

$$\mathbb{E}[(F(\theta_t) - F(\theta_\star)) \cdot \mathbb{1}(\mathcal{E})] \lesssim \begin{cases} \kappa d_{\text{eff}} \left(\frac{d d_{\text{eff}} \|\theta_0 - \theta_\star\|_2^2}{\rho T} + \frac{d \sigma_{\text{sgd}}^2}{\rho T} + \frac{\sigma_{\text{sgd}}^2}{T} \right) & \text{for DP-SGD,} \\ \kappa d_{\text{eff}} \left(\frac{\kappa d_{\text{eff}}^2 \|\theta_0 - \theta_\star\|_2^2}{\rho T^2} + \frac{\kappa d_{\text{eff}} \sigma_{\text{sgd}}^2}{\rho T^2} + \frac{\sigma_{\text{sgd}}^2}{T} \right) & \text{for } \nu\text{-DP-FTRL.} \end{cases}$$

Thus, the dimension d in DP-SGD's bound effectively becomes $\kappa d_{\text{eff}}/T$ for DP-FTRL, leading to a better dimension dependence. While faster $1/(\rho T^2)$ rates are known for DP-SGD-style algorithms for linear regression [39, 49], they require sophisticated adaptive clipping strategies. Our algorithms use a fixed clipping norm G and a fixed noise multiplier σ_{dp} independent of T ; the bounds presented above are, to the best of our knowledge, the best known in the literature for DP-SGD in this setting. We leave the exploration of combining adaptive clipping with correlated noise for future work.

Extensions to this analysis: We extend this analysis in several directions:

- We give bounds on the asymptotic suboptimality of other algorithms such as anti-PGD [44] that fall into the Noisy-FTRL framework; see Table 3 in §C.
- We give bounds on the asymptotic suboptimality for general strongly convex functions as the solution to a second-order cone program in §E. These bounds show that ν -DP-FTRL has a better asymptotic error as a function of the condition number than DP-SGD.

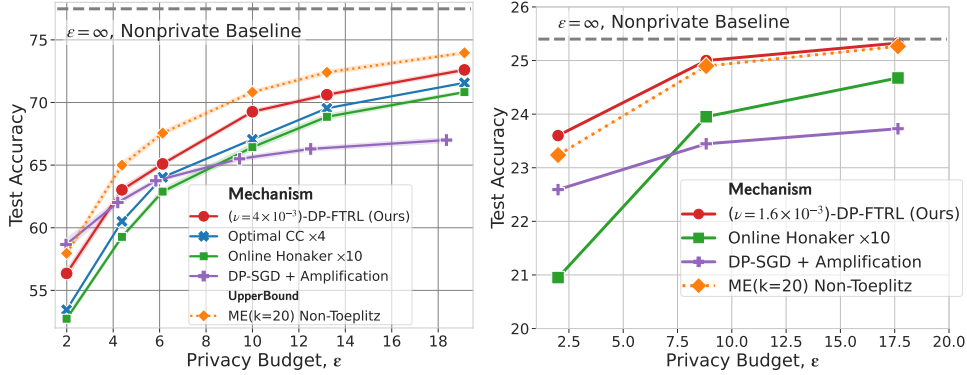


Figure 3: **Left:** Example-level DP on CIFAR-10 (non-federated image classification). **Right:** User-level DP on StackOverflow (federated language modeling). **Our ν -DP-FTRL matrices outperform all other efficient and anytime mechanisms.** They also achieve or slightly outperform state-of-the-art mechanisms that require significantly more compute (cf. Table 4 in Appendix G).

5 Experiments

We demonstrate the practical benefits and versatility of ν -DP-FTRL for two deep learning tasks across modalities and settings: (1) example-level DP for image classification on CIFAR-10 (centralized learning), and (2) federated user-level DP for language modeling with the StackOverflow dataset. The detailed setup is given in §G.

Baselines: In addition to DP-SGD (with amplification), we compare ν -DP-FTRL (per-step cost = T) with the following anytime DP-FTRL baselines: (a) the TreeAgg/Online Honaker method [33] (per-step cost $\log T$), and (b) OptimalCC [24] (per-step cost T), which corresponds to (7) with $\nu = 0$. In addition, we compare to “ME” the state-of-the-art full DP-FTRL approach [17] where computing the B matrix takes $O(T^3)$ cost and the per-step training cost is $O(T^2)$. This is infeasible² as T becomes large in modern models [5, 34], but we use it as the ceiling of DP-FTRL without amplification. Finally, we stamp/restart all baselines as suggested in [17]. This gives the baselines an advantage of an additional tuning parameter (tuned to minimize the squared error (3)), but does not affect their per-step training cost. We denote this by the suffix “ $\times S$ ” for $S > 1$ in the plot.

Main Results: Across both datasets, ν -DP-FTRL outperforms all existing anytime mechanisms by a significant margin (Figure 3, left). We find an average 3pp improvement that grows as ϵ becomes small. Indeed, the proposed ν -DP-FTRL makes up 30-80% of the gap between previous efficient approaches and the state-of-the-art and computationally intense ME approach. For instance, at $\epsilon = 10$, we have ν -DP-FTRL at 69.26% nearly matches ME at 70.83%. For StackOverflow, we find that ν -DP-FTRL matches the state-of-the-art ME at $\epsilon = 8$ and slightly exceeds it at $\epsilon = 2$ (Figure 3 right; 23.6% for our approach vs. 23.2% at $\epsilon = 2$).

As ϵ becomes small, DP-SGD can outperform DP-FTRL due to privacy amplification. We find that ν -DP-FTRL outperforms DP-SGD for $\epsilon \geq 4$ on CIFAR-10 (63.02% vs. 62.02%) and around $\epsilon \approx 2$ for StackOverflow (23.6% versus 22.6%), showing its broad applicability. Finally, we observe that that our mechanism achieves near non-private baselines on StackOverflow. A model trained via ν -DP-FTRL gets 24.8% validation accuracy at $\epsilon = 8$, a mere 0.6% off from the nonprivate baseline.

6 Conclusion

This work shows a clear separation between the noisy training dynamics with uncorrelated (DP-SGD) and correlated noise (DP-FTRL) for linear regression. The matching upper/lower bounds reveal that DP-FTRL has a better dependence than DP-SGD on problem parameters such as the effective dimension and condition number. Inspired by the theory, we propose ν -DP-FTRL and validated its empirical performance on two DP tasks spanning image and language modalities. We

²Generating B for $T = 10^4$ takes around 24 hours [17].

found it can compete the state-of-the-art while circumventing the need for any expensive computations like the semi-definite programs used in prior work. This work opens up several exciting directions including leveraging correlated-noise mechanisms for instance-optimal bounds and further improving the computational efficiency to enable private training of foundation models.

Acknowledgements

The authors thank H. Brendan McMahan, Fabian Pedregosa, Ian R. Manchester, Keith Rush, and Rahul Kidambi for fruitful discussions and helpful comments.

References

- [1] *NIST Digital Library of Mathematical Functions*. <https://dlmf.nist.gov/>, Release 1.1.10 of 2023-06-15. URL <https://dlmf.nist.gov/>. F. W. J. Olver, A. B. Olde Daalhuis, D. W. Lozier, B. I. Schneider, R. F. Boisvert, C. W. Clark, B. R. Miller, B. V. Saunders, H. S. Cohl, and M. A. McClain, eds.
- [2] Martín Abadi, Andy Chu, Ian J. Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proc. of the 2016 ACM SIGSAC Conf. on Computer and Communications Security (CCS'16)*, pp. 308–318, 2016.
- [3] Mohammed Adnan, Shivam Kalra, Jesse C Cresswell, Graham W Taylor, and Hamid R Tizhoosh. Federated learning and differential privacy for medical image analysis. *Scientific reports*, 12(1):1953, 2022.
- [4] Rafik Aguech, Eric Moulines, and Pierre Priouret. On a Perturbation Approach for the Analysis of Stochastic Tracking Algorithms. *SIAM J. Control. Optim.*, 39(3):872–899, 2000.
- [5] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.
- [6] Julyan Arbel, Olivier Marchal, and Hien D Nguyen. On strict sub-Gaussianity, optimal proxy variance and symmetry for bounded random variables. *ESAIM: Probability and Statistics*, 24: 39–55, 2020.
- [7] Francis R. Bach and Eric Moulines. Non-Strongly-Convex Smooth Stochastic Approximation with Convergence Rate $O(1/n)$. In *NeurIPS*, pp. 773–781, 2013.
- [8] Borja Balle, Gilles Barthe, Marco Gaboardi, Justin Hsu, and Tetsuya Sato. Hypothesis Testing Interpretations and Rényi Differential Privacy. In *AISTATS*, pp. 2496–2506, 2020.
- [9] Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *Proc. of the 2014 IEEE 55th Annual Symp. on Foundations of Computer Science (FOCS)*, pp. 464–473, 2014.
- [10] Mark Bun and Thomas Steinke. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Theory of Cryptography Conference*, pp. 635–658. Springer, 2016.
- [11] Paul F Byrd and Morris D Friedman. *Handbook of Elliptic Integrals for Engineers and Scientists*, volume 67. Springer, 2013.
- [12] Andrea Caponnetto and Ernesto De Vito. Optimal Rates for the Regularized Least-Squares Algorithm. *Foundations of Computational Mathematics*, 7:331–368, 2007.
- [13] Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer: Evaluating and testing unintended memorization in neural networks. In *28th USENIX Security Symposium (USENIX Security 19)*, pp. 267–284, 2019.
- [14] Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pp. 2633–2650, 2021.

- [15] Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. Quantifying memorization across neural language models. *arXiv preprint arXiv:2202.07646*, 2022.
- [16] Christopher A Choquette-Choo, Arun Ganesh, Ryan McKenna, H Brendan McMahan, Keith Rush, Abhradeep Guha Thakurta, and Zheng Xu. (amplified) banded matrix factorization: A unified approach to private training. *arXiv preprint arXiv:2306.08153*, 2023. URL <https://arxiv.org/abs/2306.08153>.
- [17] Christopher A. Choquette-Choo, Hugh Brendan McMahan, J Keith Rush, and Abhradeep Guha Thakurta. Multi-Epoch Matrix Factorization Mechanisms for Private Machine Learning. In *ICML*, volume 202, pp. 5924–5963, 23–29 Jul 2023.
- [18] Alexandre Défossez and Francis R. Bach. Averaged Least-Mean-Squares: Bias-Variance Trade-offs and Optimal Sampling Distributions. In *AISTATS*, volume 38, 2015.
- [19] Sergey Denisov, H Brendan McMahan, John Rush, Adam Smith, and Abhradeep Guha Thakurta. Improved Differential Privacy for SGD via Optimal Private Linear Operators on Adaptive Streams. *NeurIPS*, 35:5910–5924, 2022.
- [20] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating Noise to Sensitivity in Private Data Analysis. In *Proc. of the Third Conf. on Theory of Cryptography (TCC)*, pp. 265–284, 2006. URL http://dx.doi.org/10.1007/11681878_14.
- [21] Hendrik Fichtenberger, Monika Henzinger, and Jalaj Upadhyay. Constant Matters: Fine-grained Error Bound on Differentially Private Continual Observation. In *ICML*, 2023.
- [22] Jiyan He, Xuechen Li, Da Yu, Huishuai Zhang, Janardhan Kulkarni, Yin Tat Lee, Arturs Backurs, Nenghai Yu, and Jiang Bian. Exploring the Limits of Differentially Private Deep Learning with Group-wise Clipping. In *ICLR*, 2023.
- [23] William Paul Heath and Adrian G Wills. Zames-Falb multipliers for quadratic programming. In *Proceedings of the 44th IEEE Conference on Decision and Control*, pp. 963–968. IEEE, 2005.
- [24] Monika Henzinger, Jalaj Upadhyay, and Sarvagya Upadhyay. A unifying framework for differentially private sums under continual observation. *arXiv preprint arXiv:2307.08970*, 2023.
- [25] James Honaker. Efficient use of differentially private binary trees. *Theory and Practice of Differential Privacy (TPDP 2015)*, London, UK, 2015.
- [26] Daniel J. Hsu, Sham M. Kakade, and Tong Zhang. Random Design Analysis of Ridge Regression. *Found. Comput. Math.*, 14(3):569–600, 2014.
- [27] Daphne Ippolito, Florian Tramèr, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee, Christopher A Choquette-Choo, and Nicholas Carlini. Preventing verbatim memorization in language models gives a false sense of privacy. *arXiv preprint arXiv:2210.17546*, 2022.
- [28] Palak Jain, Sofya Raskhodnikova, Satchit Sivakumar, and Adam Smith. The Price of Differential Privacy under Continual Observation. In *ICML*, pp. 14654–14678. PMLR, 2023.
- [29] Prateek Jain, Sham M. Kakade, Rahul Kidambi, Praneeth Netrapalli, Venkata Krishna Pillutla, and Aaron Sidford. A Markov Chain Theory Approach to Characterizing the Minimax Optimality of Stochastic Gradient Descent (for Least Squares). In *FSTTCS*, volume 93, pp. 2:1–2:10, 2017.
- [30] Prateek Jain, Sham M. Kakade, Rahul Kidambi, Praneeth Netrapalli, and Aaron Sidford. Parallelizing Stochastic Gradient Descent for Least Squares Regression: Mini-batching, Averaging, and Model Misspecification. *J. Mach. Learn. Res.*, 18:223:1–223:42, 2017.
- [31] Prateek Jain, Sham M. Kakade, Rahul Kidambi, Praneeth Netrapalli, and Aaron Sidford. Accelerating Stochastic Gradient Descent for Least Squares Regression. In *COLT*, volume 75, pp. 545–604, 2018.

- [32] Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Keith Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, Rafael G. L. D’Oliveira, Salim El Rouayheb, David Evans, Josh Gardner, Zachary Garrett, Adrià Gascón, Badih Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaid Harchaoui, Chaoyang He, Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi, Gauri Joshi, Mikhail Khodak, Jakub Konečný, Aleksandra Korolova, Farinaz Koushanfar, Sanmi Koyejo, Tancrède Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock, Ayfer Özgür, Rasmus Pagh, Mariana Raykova, Hang Qi, Daniel Ramage, Ramesh Raskar, Dawn Song, Weikang Song, Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr, Praneeth Vepakomma, Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu, and Sen Zhao. Advances and open problems in federated learning. *CoRR*, abs/1912.04977, 2019. URL <http://arxiv.org/abs/1912.04977>.
- [33] Peter Kairouz, Brendan McMahan, Shuang Song, Om Thakkar, Abhradeep Thakurta, and Zheng Xu. Practical and private (deep) learning without sampling or shuffling. In *ICML*, 2021.
- [34] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [35] Anastasia Koloskova, Ryan McKenna, Zachary Charles, Keith Rush, and Brendan McMahan. Convergence of Gradient Descent with Linearly Correlated Noise and Applications to Differentially Private Learning. *arXiv Preprint*, 2023.
- [36] Dan Kucеровsky, Kaveh Mousavand, and Aydin Sarraf. On some properties of Toeplitz matrices. *Cogent Mathematics*, 3(1):1154705, 2016.
- [37] Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier Garcia, Christopher A Choquette-Choo, Katherine Lee, Derrick Xin, Aditya Kusupati, Romi Stella, Ankur Bapna, et al. Madlad-400: A multilingual and document-level large audited dataset. *arXiv preprint arXiv:2309.04662*, 2023.
- [38] Chao Li, Gerome Miklau, Michael Hay, Andrew McGregor, and Vibhor Rastogi. The matrix mechanism: optimizing linear counting queries under differential privacy. *The VLDB journal*, 24:757–781, 2015.
- [39] Xiyang Liu, Prateek Jain, Weihao Kong, Sewoong Oh, and Arun Sai Suggala. Near Optimal Private and Robust Linear Regression. *arXiv preprint arXiv:2301.13273*, 2023.
- [40] Brendan McMahan and Abhradeep Thakurta. Federated learning with formal differential privacy guarantees. *Google AI Blog*, 2022.
- [41] H Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning Differentially Private Recurrent Language Models. In *ICLR*, 2018.
- [42] Kamyar Moshksar. On the Absolute Constant in Hanson-Wright Inequality. *arXiv preprint*, 2021.
- [43] Alan V Oppenheim, Alan S Willsky, and Nawab. *Signals and Systems*, volume 2. 1997.
- [44] Antonio Orvieto, Hans Kersting, Frank Proske, Francis Bach, and Aurelien Lucchi. Anticorrelated Noise Injection for Improved Generalization. In *ICML*, pp. 17094–17116, 2022.
- [45] Krishna Pillutla, Yassine Laguel, Jérôme Malick, and Zaid Harchaoui. Federated learning with superquantile aggregation for heterogeneous data. *Machine Learning*, pp. 1–68, 2023.
- [46] Mark Rudelson and Roman Vershynin. Hanson-Wright Inequality and Sub-Gaussian Concentration, 2013.
- [47] Adam Smith and Abhradeep Thakurta. (nearly) optimal algorithms for private online learning in full-information and bandit settings. In *Advances in Neural Information Processing Systems*, pp. 2733–2741, 2013.

- [48] Shuang Song, Kamalika Chaudhuri, and Anand D Sarwate. Stochastic gradient descent with differentially private updates. In *2013 IEEE Global Conference on Signal and Information Processing*, pp. 245–248. IEEE, 2013.
- [49] Prateek Varshney, Abhradeep Thakurta, and Prateek Jain. (nearly) optimal private linear regression via adaptive clipping. *arXiv preprint arXiv:2207.04686*, 2022.
- [50] Zheng Xu, Yanxiang Zhang, Galen Andrew, Christopher A Choquette-Choo, Peter Kairouz, H Brendan McMahan, Jesse Rosenstock, and Yuanbo Zhang. Federated Learning of Gboard Language Models with Differential Privacy. *arXiv preprint arXiv:2305.18465*, 2023.

Appendix

Table of Contents

A Further Background on DP-FTRL	12
A.1 DP-FTRL: The Matrix Mechanism for Private Learning	12
A.2 Differential Privacy in Adaptive Streams	13
B Asymptotics of DP-FTRL for Mean Estimation	13
C Asymptotics of DP-FTRL for Linear Regression	14
C.1 Setup, Assumptions, and Notation	15
C.2 Proof of the Upper Bound on the Asymptotic Suboptimality	18
C.3 Proofs of Lower Bounds on the Asymptotic Suboptimality	23
C.4 Asymptotics of ν -Noisy-FTRL	26
C.5 Asymptotics of Anti-PGD	27
C.6 Proofs of Technical Lemmas	27
D Finite-Time Privacy-Utility Tradeoffs for Linear Regression	28
D.1 Setup, Assumptions and Notation	28
D.2 High-Probability Bounds on Noisy-FTRL	30
D.3 Expected Bounds on Noisy-FTRL	36
D.4 Privacy-Utility Guarantees of DP-FTRL	39
E Proofs for General Strongly Convex Functions	41
E.1 Proofs	42
F Technical Definitions and Lemmas	44
F.1 Linear Time-Invariant (LTI) Systems	44
F.2 Stationary Covariance of Stochastic Gradient Descent for Linear Regression	45
F.3 Concentration of Measure	46
F.4 Review of Elliptic Integrals	46
F.5 Useful Integrals	47
F.6 Other Helper Results	49
G Empirical Details	50
G.1 Image classification	51
G.2 Language modeling	51

Algorithm 1 The DP-FTRL and Noisy-FTRL algorithms with a noise correlation matrix $\mathbf{B} \in \mathbb{R}^{T \times T}$

Input: $\mathbf{B} \in \mathbb{R}^{T \times T}$, initial iterate $\boldsymbol{\theta}_0 \in \mathbb{R}^d$, ℓ_2 clip norm G , noise multiplier σ_{dp} , learning rate η , dataset \mathcal{D}

1: **for** $t = 0, \dots, T - 1$ **do**

2: Compute the gradient $\mathbf{g}_t = \begin{cases} \nabla f(\boldsymbol{\theta}_t; \mathbf{z}_t) + \nabla r(\boldsymbol{\theta}) & \text{for Noisy-FTRL,} \\ \text{clip}(\nabla f(\boldsymbol{\theta}_t; \mathbf{z}_t), G) + \nabla r(\boldsymbol{\theta}) & \text{for DP-FTRL} \end{cases}$

3: Sample $\mathbf{w}_t \sim \mathcal{N}(0, \sigma_{\text{dp}}^2 G^2 \mathbf{I}_d)$ and calculate the correlated noise $\tilde{\mathbf{w}}_t = \sum_{\tau=0}^t \mathbf{B}_{t,\tau} \mathbf{w}_\tau$

4: Update $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \tilde{\mathbf{g}}_t$ for the noisy gradient $\tilde{\mathbf{g}}_t = \mathbf{g}_t + \tilde{\mathbf{w}}_t$

Return $\boldsymbol{\theta}_T$

A Further Background on DP-FTRL

In this appendix, we give a more detailed background of DP-FTRL, and its exact notion of DP.

A.1 DP-FTRL: The Matrix Mechanism for Private Learning

The DP-FTRL algorithm [19, 33] is obtained by adapting the matrix mechanism, originally designed for linear counting queries [38], to optimization with a sequence $(\mathbf{g}_0, \dots, \mathbf{g}_{T-1})$ of gradient vectors.

Algorithm 1 gives a detailed description of DP-FTRL. We given an alternate description of DP-FTRL with an invertible lower-triangular noise correlation matrix $\mathbf{B} \in \mathbb{R}^{T \times T}$. Denoting $\mathbf{C} = \mathbf{B}^{-1}$, the iterates of DP-FTRL are generated by the update

$$\begin{pmatrix} \boldsymbol{\theta}_1 \\ \vdots \\ \boldsymbol{\theta}_T \end{pmatrix} = \begin{pmatrix} \boldsymbol{\theta}_0 \\ \vdots \\ \boldsymbol{\theta}_{T-1} \end{pmatrix} - \eta \mathbf{B} \left(\mathbf{C} \begin{pmatrix} \mathbf{g}_0 \\ \vdots \\ \mathbf{g}_{T-1} \end{pmatrix} + \begin{pmatrix} \mathbf{w}_0 \\ \vdots \\ \mathbf{w}_{T-1} \end{pmatrix} \right) \quad (10)$$

where η is a learning rate and $\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, G^2 \sigma_{\text{dp}}^2 \mathbf{I}_d)$ is i.i.d. Gaussian noise with a noise multiplier σ_{dp} and G is the ℓ_2 clip norm. It is common to refer to \mathbf{C} as the *encoder*, while \mathbf{B} is referred to as the *decoder*.

This privacy of (10) can be seen a postprocessing of a single application of the Gaussian mechanism. Let $\mathbf{G}, \mathbf{W} \in \mathbb{R}^{T \times d}$ denote the matrix where each row is the gradient \mathbf{g}_t (and respectively the noise \mathbf{w}_t). Then, (10) is effectively the postprocessing of one run one run of the Gaussian mechanism $\mathbf{C}\mathbf{G} + \mathbf{W}$. Under a neighborhood model that can change one row of \mathbf{G} , it can be seen that the maximum sensitivity of this operation is $\max_t \|\mathbf{C}_{:,t}\|_2^2$ [19]. This sensitivity logic also holds for adaptively chosen gradients; we postpone a formal description to Appendix A.2.

Connection to the exposition in prior work: Prior work introduced DP-FTRL differently. Letting $\mathbf{A} \in \mathbb{R}^{T \times T}$ denote the lower triangular matrix of all ones, update (10) can also be written as

$$\begin{pmatrix} \boldsymbol{\theta}_1 - \boldsymbol{\theta}_0 \\ \vdots \\ \boldsymbol{\theta}_T - \boldsymbol{\theta}_0 \end{pmatrix} = -\eta \tilde{\mathbf{B}} \left(\mathbf{C} \begin{pmatrix} \mathbf{g}_0 \\ \vdots \\ \mathbf{g}_{T-1} \end{pmatrix} + \begin{pmatrix} \mathbf{w}_0 \\ \vdots \\ \mathbf{w}_{T-1} \end{pmatrix} \right), \quad (11)$$

where $\tilde{\mathbf{B}} = \mathbf{A}\mathbf{B}$. The equivalence between (10) and (11) can be seen by multiplying (10) by \mathbf{A} , which is also equivalent to taking the cumulative sum of the rows of a matrix. In this notation, the objective from (3) used in previous work to find the matrix \mathbf{B} can equivalently be written as

$$\varphi(\mathbf{B}) = \|\tilde{\mathbf{B}}\|_F^2 = \|\mathbf{A}\mathbf{B}\|_F^2.$$

DP-FTRL with Toeplitz matrices: We focus on the class of lower-triangular and Toeplitz matrices \mathbf{B} . That is, $[\mathbf{B}]_{t,t'} = \beta_{t-t'}$ for all $t \geq t'$ where $\boldsymbol{\beta} = (\beta_0, \dots, \beta_{T-1})$ is the first column of \mathbf{B} .³ In this case, (10) reduces to this simple update:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \left(\mathbf{g}_t + \sum_{\tau=0}^t \beta_\tau \mathbf{w}_{t-\tau} \right). \quad (12)$$

³This implies that $\mathbf{C} = \mathbf{B}^{-1}$ is also lower-triangular and Toeplitz [36, Prop. 2.2 & Rem. 2.3].

This lets us study DP-FTRL as a time-invariant stochastic process and characterize its stationary behavior.

A.2 Differential Privacy in Adaptive Streams

Neighboring streams: We consider learning algorithms as operating over streams of gradients $\mathbf{g}_0, \mathbf{g}_1, \dots \in \mathbb{R}^d$. We consider differential privacy (DP) under the “zero-out” notion of neighborhood [33]. Two streams $\mathbf{G} = (\mathbf{g}_0, \dots, \mathbf{g}_{T-1})$ and $\mathbf{G}' = (\mathbf{g}'_0, \dots, \mathbf{g}'_{T-1})$ of length T are said to be neighbors if $\mathbf{g}_\tau = \mathbf{g}'_\tau$ for all positions $\tau \leq T - 1$ except possibly one position t where one of \mathbf{g} or \mathbf{g}' is the zero vector.

DP with adaptive continual release: It is customary to formalize DP with adaptive streams as a privacy game between a mechanism \mathcal{M} and a privacy adversary \mathcal{A} . This is known as the *adaptive continual release setting* [28]. The game makes a binary choice $b \in \{0, 1\}$ ahead of time — this remains fixed throughout and is not revealed to either \mathcal{M} or \mathcal{A} . Each round t consists of four steps:

- \mathcal{M} sends the current model parameters $\boldsymbol{\theta}_t$ to the adversary \mathcal{A} ;
- \mathcal{A} generates two gradient vectors $\mathbf{g}_t, \mathbf{g}'_t$ (e.g. as $\nabla f(\boldsymbol{\theta}_t; \mathbf{z}_t)$ for $\mathbf{z}_t \sim \mathbb{P}_{\text{data}}$ or simply the zero vector);
- the game accepts these inputs if the partial streams $(\mathbf{g}_0, \dots, \mathbf{g}_t)$ and $(\mathbf{g}'_0, \dots, \mathbf{g}'_t)$ are neighbors;
- \mathcal{M} receives \mathbf{g}_t if $b = 0$ else \mathbf{g}'_t .

DP in this setting requires that the adversary cannot infer the value of b , i.e., the distribution of $\boldsymbol{\theta}_{0:T}|b = 0$ to be “close” to that of $\boldsymbol{\theta}_{0:T}|b = 1$ (where the definition of “closeness” depends on the DP variant). For instance, (ε, δ) -DP [20] requires for each $b \in \{0, 1\}$ and any outcome set S that

$$\mathbb{P}(\boldsymbol{\theta}_{0:T} \in S | b) \leq \exp(\varepsilon) \mathbb{P}(\boldsymbol{\theta}_{0:T} \in S | 1 - b) + \delta.$$

Similarly, ρ -zCDP [10] in this setting requires that the Rényi α -divergence between the distribution P_0 of $\boldsymbol{\theta}_{0:T}|b = 0$ and the distribution P_1 of $\boldsymbol{\theta}_{0:T}|b = 1$ are close:

$$D_\alpha(P_0 || P_1) \leq \rho\alpha$$

for all $\alpha \in (0, \infty)$. Following standard arguments [e.g. 8], ρ -zCDP in this setting implies $(\varepsilon_\delta, \delta)$ -DP with

$$\varepsilon_\delta \leq \inf_{\alpha > 1} \left\{ \rho\alpha + \frac{1}{\alpha - 1} \log \left(\frac{1}{\alpha\delta} \right) + \log(1 - \alpha^{-1}) \right\}$$

DP-FTRL satisfies a zCDP guarantee as described in Theorem 2.1 in §1. This guarantee is equivalent to the one obtained by interpreting (10) as the postprocessing of one run one run of the Gaussian mechanism $\mathbf{C}\mathbf{G} + \mathbf{W}$.

B Asymptotics of DP-FTRL for Mean Estimation

We now prove Theorem 3.1 on mean estimation.

Proof of Theorem 3.1. Since $|\nabla f(\boldsymbol{\theta}; \mathbf{z})| = |\mathbf{z}| \leq 1$ and $G \geq 1$, there is no gradient clipping. Setting $\delta_t = \boldsymbol{\theta}_t - \boldsymbol{\mu}$ and $w_{\text{sgd}} = (\mathbf{z} - \mathbb{E}[\mathbf{z}])/\sigma_{\text{sgd}}$, the updates (2) can be written:

$$\delta_{t+1} = (1 - \eta)\delta_t + \eta\sigma_{\text{sgd}}w_{\text{sgd}} - \eta\sigma_{\text{dp}}G \sum_{\tau=0}^t \beta_\tau w_{t-\tau}.$$

We can then use the results from §F.1 (Theorem F.2) to obtain an expression for the steady state error:

$$F_\infty(B) = \frac{1}{2\pi} \eta^2 \int_{-\pi}^{\pi} \frac{|B(\omega)|^2 G^2 \sigma_{\text{dp}}^2 \gamma_\infty(B)^2 + \sigma_{\text{sgd}}^2}{|1 - \eta - \exp(i\omega)|^2} d\omega.$$

where we also used the fact that the SGD noise $\mathbf{z} - \mathbb{E}[\mathbf{z}]$ is iid in each step and independent of the DP noise \mathbf{w} , so that the power spectral density of the sum of these two noise sources is simply the sum of the power spectral densities of the individual sources, with the power spectral density of

Table 2: Asymptotic suboptimality of Noisy-SGD and Noisy-FTRL for linear regression with Gaussian inputs based on the eigenvalues λ_k of the Hessian \mathbf{H} . We give the bounds in terms of the learning rate η , dimension d , the effective dimension $d_{\text{eff}} = \text{Tr}[\mathbf{H}] / \|\mathbf{H}\|_2$, and the noise variance ρ^{-1} representing the privacy level. We take $G = 1$ and $\|\mathbf{H}\|_2 = 1$ w.l.o.g. Noisy-FTRL is always better at large dimension d or small learning rate η .

Eigenvalues of the Hessian \mathbf{H}	Effective dimension d_{eff}	Noisy-SGD	Noisy-FTRL	Ratio of $\frac{\text{Noisy-FTRL}}{\text{Noisy-SGD}}$
$\lambda_k = 1$	d	$\eta d \rho^{-1}$	$\eta^2 d \rho^{-1} \log^2(\frac{1}{\eta})$	$\eta \log^2(\frac{1}{\eta})$
$\lambda_k = 1/\sqrt{k}$	\sqrt{d}	$\eta d \rho^{-1}$	$\eta^2 \sqrt{d} \rho^{-1} \log^2(\frac{d}{\eta})$	$\frac{\eta}{\sqrt{d}} \log^2(\frac{d}{\eta})$
$\lambda_k = 1/k$	$\log d$	$\eta d \rho^{-1}$	$\eta^2 \rho^{-1} \log^3(\frac{d}{\eta})$	$\frac{\eta}{d} \log^3(\frac{d}{\eta})$
$\lambda_k = 1/k^2$	constant	$\eta d \rho^{-1}$	$\eta^2 \rho^{-1} \log^2(\frac{d}{\eta})$	$\frac{\eta}{d} \log^3(\frac{d}{\eta})$

the SGD noise being constant and equal to σ_{sgd} . Furthermore, we also have the expression for the sensitivity.

$$\gamma_{\infty}(B)^2 = \frac{1}{2\pi} \left(\int_{-\pi}^{\pi} |B(\omega)|^2 \right).$$

Thus F_{∞} is a product of a linear function of $|B|^2$ and $\frac{1}{|B|^2}$ (through $\gamma(B)$). By the Cauchy-Schwarz inequality, the produce is minimized when

$$\frac{|B(\omega)|^2}{|\exp(i\omega) - (1 - \eta)|^2} = \frac{1}{|B(\omega)|^2} \implies |B^*(\omega)| = |\sqrt{\exp(i\omega) - (1 - \eta)}| = |\sqrt{1 - (1 - \eta) \exp(-i\omega)}|$$

The rest of the proof follows by computing standard integrals (See §F.4 and Lemma F.15 for details). \square

C Asymptotics of DP-FTRL for Linear Regression

The goal of this section is to prove Theorem 4.1. The proof relies heavily on the following matching upper and lower bounds on the stationary error of Noisy-FTRL with any noise correlations β in the frequency domain using its discrete-time Fourier transform (DTFT) B as:

$$F_{\infty}(B) = \Theta \left(\eta \sigma_{\text{sgd}}^2 \text{Tr}[\mathbf{H}] + \eta^2 G^2 \rho^{-1} \gamma_{\infty}(B)^2 \int_{-\pi}^{\pi} |B(\omega)|^2 h(\omega) d\omega \right), \quad (13)$$

where the function $h : [-\pi, \pi] \rightarrow \mathbb{R}$ depends on the eigenvalues $\lambda_1, \dots, \lambda_d$ of the input covariance \mathbf{H} :

$$h(\omega) = \sum_{j=1}^d \frac{\lambda_j}{|1 - \exp(i\omega) - \eta \lambda_j|^2}. \quad (14)$$

The outline of the section is

- **Appendix C.1:** Setup, including notation, and assumptions.
- **Appendix C.2:** Proofs of the upper bound of (13), specifically Theorem C.15 (see also Theorem C.14 for the time-domain description).
- **Appendix C.3:** Proofs of the lower bound of (13), specifically Theorem C.17.
- **Appendix C.4:** Asymptotics of ν -Noisy-FTRL.
- **Appendix C.5:** Asymptotics of anti-PGD (see Table 3).
- **Appendix C.6:** Proofs of intermediate technical results.

The separation between Noisy-SGD and ν -Noisy-FTRL is further illustrated in Table 2. Following common practice [e.g. 12], we compare the rates for various regimes of eigenvalue decays for \mathbf{H} .

Our analysis can also give bounds for other related approaches such as anti-PGD [44].

Table 3: **Comparison to prior work:** We apply our theory to compute F_∞ for linear regression given choices of \mathbf{B} used in prior work. Though certain choices of the noise correlation β may be optimal for finite linear counting queries [21], our results show that they have $F_\infty = \infty$ because the sensitivity diverges as $T \rightarrow \infty$. ν -Noisy-FTRL effectively introduces an additional damping term $(1 - \nu)^t$ in the correlations of [21] to achieve near-optimality for linear regression. Damping similarly helps for anti-PGD [44], where the resulting error is the geometric mean of the lower bound and the bound of Noisy-SGD from Theorem 4.1.

Algorithm	Noise Correlation Weights β	Sensitivity in T steps $\gamma_T(\beta)^2$	Asymptotic Suboptimality $F_\infty(\beta)$
[21]	Eq. (7) with $\nu = 0$	$\log T$	∞
ν -Noisy-FTRL (Ours)	Eq. (7) with $0 < \nu \leq \eta\mu$	$\log(1/\nu)$	$\eta^2 G^2 \rho^{-1} \text{Tr}[\mathbf{H}] \log^2(1/\nu)$
Anti-PGD [44]	$(1, -1, 0, \dots)$	T	∞
Anti-PGD + Damping	$(1, -(1 - \nu), 0, \dots)$	$1/\nu$	$\eta^{3/2} G^2 \rho^{-1} \sqrt{d \text{Tr}[\mathbf{H}]}$

C.1 Setup, Assumptions, and Notation

C.1.1 Setup

Recall that we wish to minimize the objective

$$F(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{\text{data}}} [(y - \langle \boldsymbol{\theta}, \mathbf{x} \rangle)^2]. \quad (15)$$

Stochastic gradients: Given $(\mathbf{x}, y) \sim \mathbb{P}_{\text{data}}$, the vector

$$\mathbf{g} := (\mathbf{x} \otimes \mathbf{x}) \boldsymbol{\theta} - y \mathbf{x} = (\mathbf{x} \otimes \mathbf{x})(\boldsymbol{\theta} - \boldsymbol{\theta}_*) - \xi \mathbf{x}$$

is a stochastic gradient of F at $\boldsymbol{\theta}$, i.e., $\mathbb{E}[\mathbf{g}] = \nabla F(\boldsymbol{\theta})$.

Noisy-FTRL Iterations: We specialize the Noisy-FTRL algorithm with Toeplitz noise correlations. Let T denote the number of iterations and $\beta_{:T} = (\beta_0, \dots, \beta_{T-1})$ denote the first column of the Toeplitz matrix $\mathbf{B} = \text{Toeplitz}(\beta_{:T}) \in \mathbb{R}^{T \times T}$. Starting from a given $\boldsymbol{\theta}_0 \in \mathbb{R}^d$, Noisy-FTRL samples a fresh input-output pair $(\mathbf{x}_t, y_t) \sim \mathbb{P}_{\text{data}}$ and noise \mathbf{w}_t to set

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta ((\mathbf{x}_t \otimes \mathbf{x}_t) \boldsymbol{\theta}_t - y_t \mathbf{x}_t) - \eta \sum_{\tau=0}^t \beta_\tau \mathbf{w}_{t-\tau}. \quad (16)$$

Recall that the sensitivity $\gamma_T(\beta)$ equals to the maximum columns norm of $\mathbf{B}^{-1} = (\text{Toeplitz}(\beta))^{-1}$:

$$\gamma_T(\beta) = \max_{\tau=0, \dots, T-1} \|\mathbf{B}^{-1} \mathbf{e}_\tau\|_2, \quad (17)$$

where $\mathbf{e}_\tau = (\mathbb{I}(j = \tau))_{\tau=0}^{T-1} \in \mathbb{R}^T$ is a standard basis vector. Note that the submatrix $\mathbf{B}_{0:m, 0:m}$ of the first m rows and columns of \mathbf{B} equals $(\text{Toeplitz}(\beta_0, \dots, \beta_{m-1}))^{-1}$. Thus, the sensitivity $\gamma_t(\beta)$ is an increasing function of t always.

Infinite-time limit of Noisy-FTRL: We study the Noisy-FTRL error under the limit $T \rightarrow \infty$ with an infinite sequence $\beta = (\beta_0, \beta_1, \dots)$ of weights.

It is also convenient to re-index time to start from $t = -\infty$ and consider the sequence $(\boldsymbol{\theta})_{t=-\infty}^\infty$ produced by analogue of Equation (16), which reads

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta ((\mathbf{x}_t \otimes \mathbf{x}_t) \boldsymbol{\theta}_t - y_t \mathbf{x}_t) - \eta \sum_{\tau=0}^\infty \beta_\tau \mathbf{w}_{t-\tau}. \quad (18)$$

Note that this includes a summation over all previous DP noise $(\mathbf{w}_\tau)_{\tau=-\infty}^t$. For this sum to have finite variance, we require $\sum_{\tau=0}^\infty \beta_\tau^2 < \infty$ or that $\beta \in \ell^2$, the space of all square-summable infinite sequences. We will assume this holds throughout.

Sensitivity in the infinite limit: We define the sensitivity $\gamma_\infty(\beta)$ by consider the linear operator $\mathbf{B} = \text{Toeplitz}(\beta)$ as the convolution operator $[\mathbf{B}\mathbf{w}]_t = \sum_{\tau=0}^\infty \beta_\tau \mathbf{w}_{t-\tau}$ on input $\mathbf{w} = (\mathbf{w}_\tau)_{\tau=-\infty}^\infty$. Let \mathbf{B}^{-1} be the inverse operator to \mathbf{B} , assuming it exists. Note that the column norms $\|\mathbf{B}^{-1} \mathbf{e}_\tau\|_2$

from (17) become equal for all τ as $T \rightarrow \infty$. Thus, we get that the limiting sensitivity in the infinite time limit equals

$$\gamma_\infty(\boldsymbol{\beta}) = \|\mathbf{B}^{-1}\mathbf{e}_0\|_2 \quad (19)$$

for $\mathbf{B} = \text{Toeplitz}(\boldsymbol{\beta})$ and $\mathbf{e}_0 = (\mathbb{1}(\tau = 0))_{\tau=0}^\infty \in \ell^2$. If $\mathbf{e}_0 \notin \text{Range}(\mathbf{B})$, then we take $\gamma_\infty(\boldsymbol{\beta}) = \infty$.

Frequency-domain description: Our analysis relies on the frequency-domain representation $B : [-\pi, \pi] \rightarrow \mathbb{C}$ of $\boldsymbol{\beta}$ obtained via a discrete-time Fourier transform (DTFT) and defined as

$$B(\omega) = \sum_{t=0}^{\infty} \beta_t \exp(i\omega t). \quad (20)$$

The sequence $\boldsymbol{\beta}$ can be recovered from $B(\omega)$ using the inverse Fourier transform. Note that $\boldsymbol{\beta} \in \ell^2$ is equivalent to $B \in L^2$, the space of square integrable functions, by Parseval's theorem. The sensitivity (19) can be defined in the Fourier domain as follows.

Property C.1. Let $B(\omega)$ denote the DTFT of $\boldsymbol{\beta} \in \ell^2$. Then, we have

$$\gamma_\infty(\boldsymbol{\beta})^2 = \gamma_\infty(B)^2 := \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{d\omega}{|B(\omega)|^2}. \quad (21)$$

Proof. Let $\mathbf{z} = \mathbf{B}^{-1}\mathbf{e}_0$ be the solution of the linear system $\mathbf{B}\mathbf{z} = \mathbf{e}_0$. Let $Z(\omega)$ denote the DTFT of \mathbf{z} . Since the linear operator \mathbf{B} is a convolution with the weights of $\boldsymbol{\beta}$, this system can be expressed in the Fourier domain as

$$B(\omega)Z(\omega) = \sum_{\tau=0}^{\infty} [\mathbf{e}_0]_\tau \exp(-i\omega\tau) = 1.$$

Thus, $Z(\omega) = 1/B(\omega)$. We complete the proof with Parseval's theorem: $\|\mathbf{z}\|_2^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} |Z(\omega)|^2 d\omega$. \square

C.1.2 Assumptions

We prove the stationary error bounds under a relaxation of the assumptions in §4.

Assumption C.2. The data distribution \mathbb{P}_{data} satisfies the following:

- (B1) **Input Mean and Covariance:** The inputs have mean $\mathbb{E}[\mathbf{x}] = \mathbf{0}$ and covariance $\mathbb{E}[\mathbf{x} \otimes \mathbf{x}] =: \mathbf{H}$. Further, $L = \lambda_1 \geq \dots \geq \lambda_d =: \mu > 0$ are the eigenvalues of \mathbf{H} .
- (B2) **Noise Mean and Variance:** There exists a $\boldsymbol{\theta}_* \in \mathbb{R}^d$ such that $y = \langle \boldsymbol{\theta}_*, \mathbf{x} \rangle + \xi$ where ξ is independent of \mathbf{x} with $\mathbb{E}[\xi] = 0$ and $\mathbb{E}[\xi^2] \leq \sigma_{\text{sgd}}^2$.
- (B3) **Input Kurtosis:** There exists $R^2 < \infty$ such that $\mathbb{E}[\|\mathbf{x}\|_2^2 (\mathbf{x} \otimes \mathbf{x})] \preceq R^2 \mathbf{H}$. Moreover, for every PSD $\mathbf{P} \in \mathbb{S}_+^d$ that commutes with \mathbf{H} (i.e., $\mathbf{P}\mathbf{H} = \mathbf{H}\mathbf{P}$), we have $\mathbb{E}[(\mathbf{x} \otimes \mathbf{x})\mathbf{H}^{-1/2}\mathbf{P}\mathbf{H}^{-1/2}(\mathbf{x} \otimes \mathbf{x})] \preceq C_{\text{kurt}} \text{Tr}[\mathbf{P}] \mathbf{H}$ for some $C_{\text{kurt}} < \infty$.

These assumptions are fairly standard in the context of linear regression. Assumption (B1) implies that the Hessian matrix of objective $F(\boldsymbol{\theta})$ is $\mathbf{H} \succ 0$. Thus, F is L -smooth and μ -strongly convex. Assumption (B2) implies that $\boldsymbol{\theta}_*$ is the unique global minimizer of F and that the linear model is well-specified. The upper bounds we prove continue to hold in the case where the linear model is mis-specified (i.e. ξ is not independent of \mathbf{x}) but we still have $\mathbb{E}[\xi^2 (\mathbf{x} \otimes \mathbf{x})] \preceq \sigma_{\text{sgd}}^2 \mathbf{H}$.

Assumption (B3) is a kurtosis (i.e. 4th moment) assumption on the input distribution; we will momentarily show that it follows with absolute constants when $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{H})$. More generally, by taking a trace, we get from Jensen's inequality that $\text{Tr}[\mathbf{H}] \leq R^2$. The case of $\mathbf{P} = \mathbf{I}$ of the second part of Assumption (B3) has a special significance in the literature [e.g. 26, 31] as $C_{\text{kurt}} \text{Tr}[\mathbf{I}] = C_{\text{kurt}} d$ is the number of samples that allows the spectral concentration of the empirical covariance to the population covariance \mathbf{H} .

Property C.3. if $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{H})$, we have that Assumption (B3) holds with $R^2 \leq 3 \text{Tr}[\mathbf{H}]$ and $C_{\text{kurt}} \leq 3$.

Proof. Let $\mathbf{z} = \mathbf{H}^{-1/2}\mathbf{x}$ is element-wise independent and distributed as a standard Gaussian. For the first part, denote $\mathbf{M} = \mathbf{H}^{-1/2} \mathbb{E}[\|\mathbf{x}\|_2^2 \mathbf{x} \otimes \mathbf{x}] \mathbf{H}^{-1/2} = \mathbb{E}[\langle \mathbf{z}, \mathbf{H}\mathbf{z} \rangle \mathbf{z} \otimes \mathbf{z}]$. Elementary properties of the standard Gaussian distribution give

$$\mathbb{E}[z_k z_l z_j^2] = \begin{cases} 3, & \text{if } k = l = j \\ 1, & \text{if } k = l \neq j \\ 0, & \text{if } k \neq l, \end{cases} \quad \text{and} \quad \mathbb{E}[z_k z_l z_j z_{j'}] = \begin{cases} 1, & \text{if } k = j \text{ and } l = j' \\ 1, & \text{if } k = j' \text{ and } l = j \\ 0, & \text{else} \end{cases}$$

for $j \neq j'$. Thus, we have $\mathbf{M} = 2\mathbf{H} + \text{Tr}[\mathbf{H}]\mathbf{I}$. This gives

$$\mathbb{E}[\|\mathbf{x}\|_2^2 \mathbf{x} \otimes \mathbf{x}] = \mathbf{H}^{1/2} \mathbf{M} \mathbf{H}^{1/2} = 2\mathbf{H}^2 + \text{Tr}[\mathbf{H}]\mathbf{H} \preceq 3\text{Tr}[\mathbf{H}]\mathbf{H}.$$

For the second part, let $\mathbf{H} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$ and $\mathbf{P} = \mathbf{U}\mathbf{\Sigma}\mathbf{U}^\top$ be the eigenvalue decomposition of \mathbf{H}, \mathbf{P} respectively (since they commute, they are simultaneously diagonalized in the same basis given by the columns of \mathbf{U}). Since $\mathbf{U}^\top \mathbf{z}$ has the same distribution as \mathbf{z} by the spherical invariance of Gaussians, we have,

$$\mathbf{H}^{-1/2} \mathbb{E}[(\mathbf{x} \otimes \mathbf{x}) \mathbf{H}^{-1/2} \mathbf{P} \mathbf{H}^{-1/2} (\mathbf{x} \otimes \mathbf{x})] \mathbf{H}^{-1/2} = \mathbb{E}[(\mathbf{z} \otimes \mathbf{z}) \mathbf{P} (\mathbf{z} \otimes \mathbf{z})] = \mathbf{U} \mathbb{E}[(\mathbf{z} \otimes \mathbf{z}) \mathbf{\Sigma} (\mathbf{z} \otimes \mathbf{z})] \mathbf{U}^\top. \quad (22)$$

Each off-diagonal entry of $\mathbb{E}[(\mathbf{z} \otimes \mathbf{z}) \mathbf{\Sigma} (\mathbf{z} \otimes \mathbf{z})]$ is zero since it involves expected odd powers of Gaussians. Its j^{th} diagonal entry equals (denoting $\sigma_j := [\mathbf{\Sigma}]_{j,j}$)

$$\mathbb{E}\left[z_j^2 \sum_{k=1}^d \sigma_k z_k^2\right] = \sigma_j \mathbb{E}[z_j^4] + \sum_{k \neq j} \sigma_k \mathbb{E}[z_j^2 z_k^2] = 2\sigma_j + \text{Tr}[\mathbf{\Sigma}].$$

This gives $\mathbb{E}[(\mathbf{z} \otimes \mathbf{z}) \mathbf{\Sigma} (\mathbf{z} \otimes \mathbf{z})] = 2\mathbf{\Sigma} + \text{Tr}[\mathbf{\Sigma}]\mathbf{I} \preceq 3\text{Tr}[\mathbf{\Sigma}]\mathbf{I}$ since $\mathbf{\Sigma} \succeq \mathbf{0}$. Plugging this back into (22) and rearranging completes the proof. \square

C.1.3 Notation

We set up some notation, that we use throughout this section.

- It is convenient to rewrite the Noisy-FTRL recursion in terms of the difference $\boldsymbol{\theta}'_t := \boldsymbol{\theta}_t - \boldsymbol{\theta}_*$. We can rewrite the Noisy-FTRL recursion (18) as

$$\boldsymbol{\theta}'_{t+1} = (\mathbf{I} - \eta(\mathbf{x}_t \otimes \mathbf{x}_t))\boldsymbol{\theta}'_t + \eta \xi_t \mathbf{x}_t - \eta \sum_{\tau=0}^{\infty} \beta_\tau \mathbf{w}_{t-\tau}. \quad (23)$$

We will analyze this recursion.

- We describe the asymptotic suboptimality in terms of the self-adjoint linear operator $\mathbf{T} : \ell^2 \rightarrow \ell^2$ defined by

$$[\mathbf{T}\boldsymbol{\beta}]_t = \sum_{\tau=0}^{\infty} \beta_\tau \sum_{j=1}^d (1 - \eta\lambda_j)^{|t-\tau|}. \quad (24)$$

This operator is positive semi-definite, as we show in Lemma C.6 below. In the finite time setting, we could represent \mathbf{T} by the matrix

$$\mathbf{T} = \begin{bmatrix} d & \sum_{j=1}^d (1 - \eta\lambda_j) & \sum_{j=1}^d (1 - \eta\lambda_j)^2 & \cdots \\ \sum_{j=1}^d (1 - \eta\lambda_j) & d & \sum_{j=1}^d (1 - \eta\lambda_j) & \cdots \\ \sum_{j=1}^d (1 - \eta\lambda_j)^2 & \sum_{j=1}^d (1 - \eta\lambda_j) & d & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

We only consider step-size $0 < \eta < 1/R^2$, which implies that $1 - \eta\lambda_j \in (0, 1)$ for all j .

- For $j = 1, \dots, d$, define $\mathbf{T}_j : \ell^2 \rightarrow \ell^2$ as the linear operator

$$[\mathbf{T}_j \boldsymbol{\beta}]_t = \sum_{\tau=0}^{\infty} \beta_\tau (1 - \eta\lambda_j)^{|t-\tau|}. \quad (25)$$

Note that $[\mathbf{T}_j \boldsymbol{\beta}]_t < \infty$ always since

$$\sum_{\tau=0}^{\infty} \beta_{\tau} (1 - \eta \lambda_j)^{|\tau-t|} \leq \frac{2 \|\boldsymbol{\beta}\|_{\infty}}{\eta \lambda_j} < \infty,$$

since $0 < \eta \lambda < 1$. Thus, we have that $\mathbf{T} = \sum_{j=1}^d \mathbf{T}_j$ by the bounded convergence theorem. Further, we show in the upcoming Lemma C.6 that each \mathbf{T}_j are PSD.

- Define $\boldsymbol{\Sigma}_{\boldsymbol{\beta}}, \mathbf{P}_{\boldsymbol{\beta}} \in \mathbb{S}^d$ as

$$\boldsymbol{\Sigma}_{\boldsymbol{\beta}} := \text{diag} \left((\langle \boldsymbol{\beta}, \mathbf{T}_j \boldsymbol{\beta} \rangle)_{j=1}^d \right), \quad \text{and} \quad \mathbf{P}_{\boldsymbol{\beta}} = \mathbf{U} \boldsymbol{\Sigma}_{\boldsymbol{\beta}} \mathbf{U}^{\top}, \quad (26)$$

where \mathbf{U} is the eigen-basis of $\mathbf{H} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^{\top}$. By definition, $\mathbf{P}_{\boldsymbol{\beta}}$ commutes with \mathbf{H} since $\mathbf{P}_{\boldsymbol{\beta}} \mathbf{H} = \mathbf{H} \mathbf{P}_{\boldsymbol{\beta}} = \mathbf{U} (\boldsymbol{\Lambda} \boldsymbol{\Sigma}_{\boldsymbol{\beta}}) \mathbf{U}^{\top}$. Further, since each \mathbf{T}_j is PSD (Lemma C.6), we have that $\boldsymbol{\Sigma}_{\boldsymbol{\beta}}$ and $\mathbf{P}_{\boldsymbol{\beta}}$ are PSD as well. We also have

$$\text{Tr} [\mathbf{P}_{\boldsymbol{\beta}}] = \text{Tr} [\boldsymbol{\Sigma}_{\boldsymbol{\beta}}] = \langle \boldsymbol{\beta}, \mathbf{T} \boldsymbol{\beta} \rangle. \quad (27)$$

- Define the matrix $\mathbf{M}_{\omega} \in \mathbb{C}^{d \times d}$ as

$$\mathbf{M}_{\omega} = ((1 - \exp(i\omega)) \mathbf{I} - \eta \mathbf{H})^{-1}. \quad (28)$$

Throughout, we assume that Assumption C.2 holds.

Preliminary lemmas: This lemma helps us move back and forth between the time-domain and frequency-domain representations. See Appendix C.6 for a proof.

Lemma C.4. Consider $\boldsymbol{\beta} \in \ell^2$ and its DTFT $B(\omega)$. If $0 < \eta < 1/\lambda_j$, we have

$$\frac{1}{2} \langle \boldsymbol{\beta}, \mathbf{T}_j \boldsymbol{\beta} \rangle \leq \frac{\eta \lambda_j}{2\pi} \int_{-\pi}^{\pi} \frac{|B(\omega)|^2 d\omega}{|1 - \eta \lambda_j - \exp(i\omega)|^2} \leq \langle \boldsymbol{\beta}, \mathbf{T}_j \boldsymbol{\beta} \rangle.$$

Setting $B(\omega) = 1$ and $\boldsymbol{\beta} = (1, 0, \dots)$ gives the next corollary.

Corollary C.5. If $0 < \eta < 1/\lambda_j$, we have,

$$\frac{1}{2} \leq \frac{\eta \lambda_j}{2\pi} \int_{-\pi}^{\pi} \frac{d\omega}{|1 - \eta \lambda_j - \exp(i\omega)|^2} \leq 1.$$

Lemma C.6. The operators \mathbf{T}_j defined in (25) and \mathbf{T} defined in (24) are both positive semi-definite for $\eta < 1/\max_{j \in [d]} \lambda_j$.

Proof. Consider any $\boldsymbol{\beta} \in \ell^2$ and its DTFT $B(\omega)$. We have from Lemma C.4 that

$$0 \leq \int_{-\pi}^{\pi} \frac{|B(\omega)|^2 d\omega}{|1 - \eta \lambda_j - \exp(i\omega)|^2} \leq \frac{2\pi}{\eta \lambda_j} \langle \boldsymbol{\beta}, \mathbf{T}_j \boldsymbol{\beta} \rangle,$$

or that $\langle \boldsymbol{\beta}, \mathbf{T}_j \boldsymbol{\beta} \rangle \geq 0$. □

C.2 Proof of the Upper Bound on the Asymptotic Suboptimality

The key tool in the analysis is the use of linear time-invariant (LTI) input-output systems to relate the output covariance to the input covariance using its *transfer function* (see Appendix F.1 for a summary). The Noisy-FTRL recursion is not trivial to characterize in this manner because the update (18) is not LTI. Instead, we decompose it into an infinite sequence of LTI systems and carefully analyze the error propagation.

This consists of the following steps:

- Part 1: Decompose the Noisy-FTRL recursion into a sequence of LTI systems.
- Part 2: Compute the transfer function of each LTI system.
- Part 3: Compute the stationary covariance for each LTI system from the previous one.
- Part 4: Combine the stationary covariances to get the stationary error of the original iterate.

C.2.1 Part 1: Decomposition into a Sequence of LTI Systems

A challenge in analyzing the stationary error of Equation (23) in the frequency domain is that it is not an LTI system. Replacing $\mathbf{x}_t \otimes \mathbf{x}_t$ by \mathbf{H} in Equation (23) results in a LTI update; this system is quite similar to fixed design linear regression. However, this leads to an error in the general case, which satisfies a recursion of the same form as (23). We can repeat the same technique of replacing $\mathbf{x}_t \otimes \mathbf{x}_t$ by \mathbf{H} and repeat this process indefinitely. This proof technique has been used in [4] to analyze stochastic tracking algorithms and [7] to analyze iterate-averaged SGD for linear regression. We adopt this technique to analyze the stationary covariance of DP mechanisms with correlated noise.

We define sequences $(\boldsymbol{\theta}_t^{(r)})_{t=-\infty}^{\infty}$ and $(\boldsymbol{\delta}_t^{(r)})_{t=-\infty}^{\infty}$ for $r \geq 0$ as follows:

$$\begin{aligned}\boldsymbol{\theta}_{t+1}^{(0)} &= (\mathbf{I} - \eta\mathbf{H})\boldsymbol{\theta}_t^{(0)} + \eta\xi_t\mathbf{x}_t - \eta\sum_{\tau=0}^{\infty}\beta_{\tau}\mathbf{w}_{t-\tau}, \\ \boldsymbol{\theta}_{t+1}^{(r)} &= (\mathbf{I} - \eta\mathbf{H})\boldsymbol{\theta}_t^{(r)} + \eta(\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t)\boldsymbol{\theta}_t^{(r-1)} \text{ for } r > 0, \\ \boldsymbol{\delta}_{t+1}^{(r)} &= (\mathbf{I} - \eta\mathbf{x}_t \otimes \mathbf{x}_t)\boldsymbol{\delta}_t^{(r)} + \eta(\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t)\boldsymbol{\theta}_t^{(r)}.\end{aligned}\tag{29}$$

These recursions are assumed to start at $t = -\infty$ from $\boldsymbol{\theta}_t^{(0)} = \boldsymbol{\theta}'_t$, $\boldsymbol{\delta}_t^{(r)} = \mathbf{0}$ for $r \geq 0$ and $\boldsymbol{\theta}_t^{(r)} = \mathbf{0}$ for $r > 0$. These recursions are a decomposition of (23) as we define below.

Property C.7. For each iteration t and any integer $m \geq 0$, we have $\boldsymbol{\theta}'_t = \sum_{r=0}^m \boldsymbol{\theta}_t^{(r)} + \boldsymbol{\delta}_t^{(m)}$.

Proof. We prove this by induction. The base case at $t = -\infty$ holds by definition. Assume that this is true for some integer t . Then, we have

$$\begin{aligned}\sum_{r=0}^m \boldsymbol{\theta}_{t+1}^{(r)} + \boldsymbol{\delta}_{t+1}^{(m)} &= (\mathbf{I} - \eta\mathbf{x}_t \otimes \mathbf{x}_t) \left(\sum_{r=0}^m \boldsymbol{\theta}_t^{(r)} + \boldsymbol{\delta}_t^{(m)} \right) + \eta\xi_t\mathbf{x}_t - \eta\sum_{\tau=0}^{\infty}\beta_{\tau}\mathbf{w}_{t-\tau} \\ &= (\mathbf{I} - \eta\mathbf{x}_t \otimes \mathbf{x}_t)\boldsymbol{\theta}'_t + \eta\xi_t\mathbf{x}_t - \eta\sum_{\tau=0}^{\infty}\beta_{\tau}\mathbf{w}_{t-\tau} = \boldsymbol{\theta}'_{t+1}.\end{aligned}$$

□

The idea behind the proof is to show that $\mathbb{E}[\boldsymbol{\delta}_0^{(m)} \otimes \boldsymbol{\delta}_0^{(m)}] \rightarrow \mathbf{0}$ as $m \rightarrow \infty$. Then, we can use the triangle inequality to bound

$$\|\boldsymbol{\theta}'_t\| \leq \sum_{r=0}^{\infty} \|\boldsymbol{\theta}_t^{(r)}\|,$$

where the stationary error of the right side can be obtained from analyzing the LTI systems defined in (29).

C.2.2 Part 2: Characterize the Transfer Function of each LTI System

There are two LTI systems. First, $\boldsymbol{\theta}_t^{(r)}$ for $r > 0$ is an LTI system

$$\mathbf{z}_{t+1} = (\mathbf{I} - \eta\mathbf{H})\mathbf{z}_t + \eta\mathbf{u}_t\tag{30}$$

with input $\mathbf{u}_t \in \mathbb{R}^d$ and output $\mathbf{z}_t \in \mathbb{R}^d$. Second, $\boldsymbol{\theta}_t^{(0)}$ satisfies satisfies an LTI system

$$\mathbf{z}_{t+1} = (\mathbf{I} - \eta\mathbf{H})\mathbf{z}_t + \eta\mathbf{u}_t - \eta\sum_{\tau=0}^{\infty}\beta_{\tau}\mathbf{w}_{t-\tau}\tag{31}$$

with inputs $(\mathbf{u}_t, \mathbf{w}_t) \in \mathbb{R}^d \times \mathbb{R}^d$ and output $\mathbf{z}_t \in \mathbb{R}^d$ where the weights $\beta \in \ell^2$ is assumed given.

We now characterize the transfer functions of these LTI systems; see Appendix F.1 for a review.

Property C.8. The LTI system (30) is $\mathbf{G}(\omega) = -\eta\mathbf{M}_{\omega} \in \mathbb{C}^{d \times d}$, where \mathbf{M}_{ω} is defined in Equation (28). Moreover, this system is asymptotically stable as long as $\mathbf{0} \prec \eta\mathbf{H} \prec \mathbf{I}$.

Proof. Let $\mathbf{U}(\omega) \in \mathbb{C}^d$ and $\mathbf{Z}(\omega) \in \mathbb{C}^d$ be the Fourier transforms of \mathbf{u}_t and \mathbf{z}_t respectively. The transfer function must hold for any input-output sequences, so we can choose some sequences and solve for the transfer functions. It is convenient to consider the delta spike on a standard basis (upto scaling), i.e., $\mathbf{U} = 2\pi\delta_\omega \mathbf{e}_j$, where δ_ω is the Dirac delta at ω , and \mathbf{e}_j is the j^{th} standard basis vector in \mathbb{R}^d . This gives $\mathbf{Z} = 2\pi\mathbf{g}_j\delta_\omega$ where $\mathbf{g}_j(\cdot)$ is the j^{th} column of $\mathbf{G}(\cdot)$.

To move back to time domain, we take an inverse Fourier transform to get $\mathbf{u}_t = \exp(i\omega t)\mathbf{e}_j$ and $\mathbf{z}_t = \mathbf{g}_j(\omega)\exp(i\omega t)$. Plugging this into the update (30) gives and solving for $\mathbf{g}_j(\omega)$ gives $\mathbf{g}_j(\omega) = -\eta\mathbf{M}_\omega\mathbf{e}_j$. Stacking these into a matrix gives the expression.

If $\mathbf{u}_t \equiv \mathbf{0}$ for all t , then $\|\mathbf{z}_{t+1}\|_2 \leq \|\mathbf{I} - \eta\mathbf{H}\|_2\|\mathbf{z}_t\|_2 < \|\mathbf{z}_t\|_2$ since $\|\mathbf{I} - \eta\mathbf{H}\|_2 < 1$. Hence, $\|\mathbf{z}_t\|_2 \rightarrow 0$ giving the asymptotic stability of the system. \square

Property C.9. *The transfer function of the LTI system (31) is*

$$\tilde{\mathbf{G}}(\omega) = [\mathbf{G}(\omega) \quad \mathbf{G}'(\omega)] \in \mathbb{C}^{d \times 2d}$$

where $\mathbf{G}(\omega) = -\eta\mathbf{M}_\omega$ and $\mathbf{G}'(\omega) = \eta B(\omega)\mathbf{M}_\omega$ with $B(\omega)$ as the DTFT of β . Moreover, this system is asymptotically stable as long as $\mathbf{0} \prec \eta\mathbf{H} \prec \mathbf{I}$.

Proof. The expression for $\mathbf{G}(\omega)$ is the same as in Property C.8. To find \mathbf{G}' , we set the Fourier transforms $\mathbf{U} \equiv \mathbf{0}$, $\mathbf{W} = 2\pi\delta_\omega \mathbf{e}_j$ so that $\mathbf{Z} = 2\pi\delta_\omega \mathbf{g}'_j$, where $\mathbf{g}'_j(\cdot)$ is the j^{th} column of $\mathbf{G}'(\cdot)$.

An inverse Fourier transforms gives the time domain versions $\mathbf{w}_t = \exp(i\omega t)\mathbf{u}_t \equiv \mathbf{0}$, $\mathbf{z}_t = \exp(i\omega t)\mathbf{g}'_j(\omega)$. Plugging these into (31) and plugging in the definition of $B(\omega)$ gives the expression for the transfer function. Its asymptotic stability holds similar to Property C.8. \square

C.2.3 Part 3: Compute the Stationary Covariance of each LTI System

The stationary covariance of an LTI system driven by white noise. A sequence (\mathbf{u}_t) is said to be a white noise process if it is mean zero and $\mathbb{E}[\mathbf{u}_t\mathbf{u}_\tau] = \mathbf{0}$ for $t \neq \tau$. This is true for both $\theta_t^{(0)}$ as well $\theta_t^{(r)}$ for $r > 0$. Since we care about the stationary distribution and we start at $t = -\infty$, we have reached the steady state at $t = 0$. So, we compute $\mathbb{E}[\theta_0^{(r)} \otimes \theta_0^{(r)}]$.

Stationary covariance of the base recursion: We first start with $\theta_t^{(0)}$.

Proposition C.10. *We have that $\mathbb{E}[\theta_t^{(0)} \otimes \theta_t^{(0)}]$ is equal for all $t > -\infty$ and is bounded as*

$$\mathbb{E}[\theta_t^{(0)} \otimes \theta_t^{(0)}] \preceq \eta\sigma_{\text{sgd}}^2\mathbf{I} + \eta\sigma^2\mathbf{H}^{-1/2}\mathbf{P}_\beta\mathbf{H}^{-1/2},$$

where \mathbf{P}_β is defined in Equation (26) and we denote $\sigma^2 = G^2\gamma_\infty(\beta)^2/(2\rho)$.

Proof. The input $(\xi_t\mathbf{x}_t, \mathbf{w}_t)$ forms a white noise sequence, since for $t \neq \tau$, we have $\mathbb{E}[\xi_t\mathbf{x}_t\xi_\tau\mathbf{x}_\tau] = \mathbb{E}[\xi_t\mathbf{x}_t]\mathbb{E}[\xi_\tau\mathbf{x}_\tau] = \mathbf{0}$ (since $\xi_t\mathbf{x}_t$ for each t is i.i.d.) and $\mathbb{E}[\mathbf{w}_t\mathbf{w}_\tau] = \mathbf{0}$. The covariance of the input is

$$\mathbb{E}[(\xi_t\mathbf{x}_t, \mathbf{w}_t) \otimes (\xi_t\mathbf{x}_t, \mathbf{w}_t)] = \begin{bmatrix} \mathbb{E}[\xi_t^2\mathbf{x}_t\mathbf{x}_t] & \mathbf{0} \\ \mathbf{0} & \mathbb{E}[\mathbf{w}_t \otimes \mathbf{w}_t] \end{bmatrix} = \mathbb{E}[(\xi_\tau\mathbf{x}_\tau, \mathbf{w}_\tau) \otimes (\xi_\tau\mathbf{x}_\tau, \mathbf{w}_\tau)]$$

for all t, τ . This is further bounded by Assumption (B1) as

$$\mathbb{E}[(\xi_t\mathbf{x}_t, \mathbf{w}_t) \otimes (\xi_t\mathbf{x}_t, \mathbf{w}_t)] \preceq \begin{bmatrix} \sigma_{\text{sgd}}^2\mathbf{H} & \mathbf{0} \\ \mathbf{0} & \sigma^2\mathbf{I} \end{bmatrix}$$

The output covariance of the asymptotically stable LTI system (31) can be given in terms of the transfer function $\tilde{\mathbf{G}}(\omega) = [\mathbf{G}(\omega) \quad \mathbf{G}'(\omega)]$ characterized in Property C.9 using Theorem F.2. This gives that $\mathbb{E}[\theta_t^{(0)} \otimes \theta_t^{(0)}]$ is equal for each $t > -\infty$ and is bounded as

$$\mathbb{E}[\theta_t^{(0)} \otimes \theta_t^{(0)}] \preceq \frac{1}{2\pi} \int_{-\pi}^{\pi} (\eta^2\sigma_{\text{sgd}}^2\mathbf{M}_\omega\mathbf{H}\mathbf{M}_\omega^* + \eta^2\sigma^2|B(\omega)|^2\mathbf{M}_\omega\mathbf{M}_\omega^*) d\omega. \quad (32)$$

With the eigenvalue decomposition $\mathbf{H} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$, we get $\mathbf{M}_\omega = \mathbf{U}((1 - \exp(i\omega))\mathbf{I} - \eta\mathbf{\Lambda})^{-1}\mathbf{U}^\top$. This gives

$$\mathbf{M}_\omega \mathbf{H} \mathbf{M}_\omega^* = \mathbf{U} \operatorname{diag} \left((\lambda_j / |1 - \exp(i\omega) - \eta\lambda_j|^2)_{j=1}^d \right) \mathbf{U}^\top.$$

We invoke Corollary C.5 to say

$$\begin{aligned} \int_{-\pi}^{\pi} \mathbf{M}_\omega \mathbf{H} \mathbf{M}_\omega^* d\omega &= \mathbf{U} \operatorname{diag} \left(\left(\int_{-\pi}^{\pi} d\omega \lambda_j / |1 - \exp(i\omega) - \eta\lambda_j|^2 \right)_{j=1}^d \right) \mathbf{U}^\top \\ &\preceq \mathbf{U} \operatorname{diag} \left((2\pi/\eta)_{j=1}^d \right) \mathbf{U}^\top = \frac{2\pi}{\eta} \mathbf{I}. \end{aligned} \quad (33)$$

Similarly, we invoke Lemma C.4 to compute

$$\begin{aligned} \int_{-\pi}^{\pi} |B(\omega)|^2 \mathbf{M}_\omega \mathbf{M}_\omega^* d\omega &= \mathbf{U} \operatorname{diag} \left(\left(\int_{-\pi}^{\pi} d\omega |B(\omega)|^2 / |1 - \exp(i\omega) - \eta\lambda_j|^2 \right)_{j=1}^d \right) \mathbf{U}^\top \\ &\preceq \mathbf{U} \operatorname{diag} \left((2\pi \langle \boldsymbol{\beta}, \mathbf{T}_j \boldsymbol{\beta} \rangle / (\eta\lambda_j))_{j=1}^d \right) \mathbf{U}^\top \\ &= \frac{2\pi}{\eta} \mathbf{U} \mathbf{\Lambda}^{-1/2} \boldsymbol{\Sigma}_\beta \mathbf{\Lambda}^{-1/2} \mathbf{U}^\top = \frac{2\pi}{\eta} \mathbf{H}^{-1/2} \mathbf{P}_\beta \mathbf{H}^{-1/2}, \end{aligned} \quad (34)$$

where $\boldsymbol{\Sigma}_\beta$ and \mathbf{P}_β are defined in (26). Plugging in (33) and (33) into (32) completes the proof of the upper bound. \square

Stationary covariance of the higher-order recursion: Next, we turn to $\boldsymbol{\theta}_t^{(r)}$.

Proposition C.11. *For any $r \geq 1$, we have*

$$\mathbb{E} \left[\boldsymbol{\theta}_0^{(r)} \otimes \boldsymbol{\theta}_0^{(r)} \right] \preceq \eta (\eta R^2)^r \left(\sigma_{\text{sgd}}^2 + \frac{C_{\text{kurt}} \sigma^2}{R^2} \langle \boldsymbol{\beta}, \mathbf{T} \boldsymbol{\beta} \rangle \right).$$

Proof. Follows from combining Proposition C.10 with the more general Lemma C.12 below. \square

Lemma C.12. *For some $r \geq 1$, suppose that $\mathbb{E} \left[\boldsymbol{\theta}_t^{(r-1)} \otimes \boldsymbol{\theta}_t^{(r-1)} \right]$ is equal for each t and is bounded as $\mathbb{E} \left[\boldsymbol{\theta}_t^{(r-1)} \otimes \boldsymbol{\theta}_t^{(r-1)} \right] \preceq a\mathbf{I} + b\mathbf{H}^{-1/2} \mathbf{P}_\beta \mathbf{H}^{-1/2}$ for some scalars $a, b \geq 0$. Then, we have the following.*

(a) *We have that $\boldsymbol{\zeta}_t^{(r)} := (\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t) \boldsymbol{\theta}_t^{(r-1)}$ is a white-noise process with*

$$\mathbb{E} \left[\boldsymbol{\zeta}_t^{(r)} \otimes \boldsymbol{\zeta}_t^{(r)} \right] \preceq (aR^2 + bC_{\text{kurt}} \langle \boldsymbol{\beta}, \mathbf{T} \boldsymbol{\beta} \rangle) \mathbf{H}.$$

(b) *We have that $\mathbb{E} \left[\boldsymbol{\theta}_t^{(r)} \otimes \boldsymbol{\theta}_t^{(r)} \right]$ is equal for each t and is bounded as*

$$\mathbb{E} \left[\boldsymbol{\theta}_t^{(r)} \otimes \boldsymbol{\theta}_t^{(r)} \right] \preceq \eta (aR^2 + bC_{\text{kurt}} \langle \boldsymbol{\beta}, \mathbf{T} \boldsymbol{\beta} \rangle) \mathbf{I}.$$

Proof. Note that $\mathbb{E} \left[\boldsymbol{\zeta}_t^{(r)} \otimes \boldsymbol{\zeta}_\tau^{(r)} \right] = \mathbf{0}$ for $t \neq \tau$ since \mathbf{x}_t is independent of \mathbf{x}_τ and $\mathbb{E}[\mathbf{x}_t \otimes \mathbf{x}_t] = \mathbf{H}$. Since \mathbf{x}_t is independent of $\boldsymbol{\theta}_t^{(r-1)}$, we get from the tower rule of expectations that

$$\begin{aligned} \mathbb{E} \left[\boldsymbol{\zeta}_t^{(r)} \otimes \boldsymbol{\zeta}_t^{(r)} \right] &= \mathbb{E} \left[(\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t) \left(\boldsymbol{\theta}_t^{(r-1)} \otimes \boldsymbol{\theta}_t^{(r-1)} \right) (\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t) \right] \\ &= \mathbb{E} \left[(\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t) \mathbb{E} \left[\boldsymbol{\theta}_t^{(r-1)} \otimes \boldsymbol{\theta}_t^{(r-1)} \right] (\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t) \right], \end{aligned}$$

or that $(\boldsymbol{\zeta}_t^{(r)})$ is a white noise process. Its covariance can further be bounded as

$$\begin{aligned} \mathbb{E} \left[\boldsymbol{\zeta}_t^{(r)} \otimes \boldsymbol{\zeta}_t^{(r)} \right] &\preceq \mathbb{E} \left[(\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t) \left(a\mathbf{I} + b\mathbf{H}^{-1/2} \mathbf{P}_\beta \mathbf{H}^{-1/2} \right) (\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t) \right] \\ &\preceq a \mathbb{E} \left[\|\mathbf{x}_t\|_2^2 (\mathbf{x}_t \otimes \mathbf{x}_t) \right] + b \mathbb{E} \left[(\mathbf{x}_t \otimes \mathbf{x}_t) \mathbf{H}^{-1/2} \mathbf{P}_\beta \mathbf{H}^{-1/2} (\mathbf{x}_t \otimes \mathbf{x}_t) \right] \\ &\preceq aR^2 \mathbf{H} + bC_{\text{kurt}} \operatorname{Tr}[\mathbf{P}_\beta] \mathbf{H}, \end{aligned}$$

where the last inequality followed from Assumption **(B3)**. Further, note that $\text{Tr}[\mathbf{P}_\beta] = \langle \beta, \mathbf{T}\beta \rangle$ from (27).

The output covariance of the asymptotically stable LTI system (30) can be given in terms of the transfer function $\mathbf{G}(\omega) = -\eta \mathbf{M}_\omega$ using Theorem F.2 as

$$\mathbb{E} \left[\boldsymbol{\theta}_t^{(r)} \otimes \boldsymbol{\theta}_t^{(r)} \right] \preceq \frac{\eta^2 (aR^2 + bC_{\text{kurt}} \langle \beta, \mathbf{T}\beta \rangle)}{2\pi} \int_{-\pi}^{\pi} \mathbf{M}_\omega \mathbf{H} \mathbf{M}_\omega^* d\omega \stackrel{(33)}{\preceq} \eta (aR^2 + bC_{\text{kurt}} \langle \beta, \mathbf{T}\beta \rangle) \mathbf{I}.$$

□

Remainder Term: It remains to show that the remainder term δ_t can be neglected by taking $m \rightarrow \infty$.

Proposition C.13. *We have $\lim_{m \rightarrow \infty} \mathbb{E} \left[\delta_t^{(m)} \otimes \delta_t^{(m)} \right] = \mathbf{0}$.*

Proof. Let $\zeta_t^{(m+1)} := (\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t) \boldsymbol{\theta}_t^{(m)}$. By Lemma C.12 and Proposition C.11, we have ζ_t is a white-noise process with

$$\mathbb{E} \left[\zeta_t^{(m+1)} \otimes \zeta_t^{(m+1)} \right] \preceq (\eta R^2)^{m+1} \left(\sigma_{\text{sgd}}^2 + \frac{C_{\text{kurt}} \sigma^2}{R^2} \langle \beta, \mathbf{T}\beta \rangle \right) \mathbf{H} \rightarrow \mathbf{0}$$

as $m \rightarrow \infty$ since $\eta < 1/R^2$. Note that the update for $\delta_t^{(m)}$ exactly matches that of SGD (without added DP noise), and the noise covariance is $\mathbf{0}$. The statement of this result is equivalent to showing that the stationary covariance of SGD with zero residuals is zero. This observation is formalized in Lemma 4 of [29] (see also Theorem F.3 of Appendix F), which gives for any t that

$$\mathbf{0} \preceq \mathbb{E}[\delta_t^{(m)} \otimes \delta_t^{(m)}] \preceq \frac{\eta}{1 - \eta R^2} \left[(\eta R^2)^{m+1} \left(\sigma_{\text{sgd}}^2 + \frac{C_{\text{kurt}} \sigma^2}{R^2} \langle \beta, \mathbf{T}\beta \rangle \right) \right] \mathbf{I} \rightarrow \mathbf{0}$$

as $m \rightarrow \infty$. □

C.2.4 Part 4: Combining the Errors

Time-domain description: We now state and prove a time-domain description of the upper bound of Equation (13).

Theorem C.14. *Suppose Assumption C.2 holds. Consider the sequence $(\boldsymbol{\theta}_t)_{t=-\infty}^{\infty}$ produced by the Noisy-FTRL update in Equation (18) with some given weights $\beta \in \ell^2$ and noise variance $\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, G^2 \gamma_\infty(\beta)^2 / (2\rho) \mathbf{I})$. If the learning rate satisfies $\eta < 1/R^2$, we have*

$$F_\infty(\beta) \leq \left(1 + \left(1 - \sqrt{\eta R^2} \right)^{-2} \right) \eta R^2 \sigma_{\text{sgd}}^2 + \left(1 + C_{\text{kurt}} \left(1 - \sqrt{\eta R^2} \right)^{-2} \right) \frac{\eta G^2 \gamma_\infty(\beta)^2}{2\rho} \langle \beta, \mathbf{T}\beta \rangle.$$

Proof. We use shorthand $\sigma^2 = \frac{G^2 \gamma_\infty(\beta)^2}{2\rho}$. First, note that $\eta < 1/R^2$ also implies that $\eta \lambda_j < 1$ for each eigenvalue λ_j of \mathbf{H} . The right side is well-defined since Lemma F.17 gives

$$|\langle \beta, \mathbf{T}\beta \rangle| \leq \sum_{j=1}^d \left| \sum_{t=0}^{\infty} \sum_{\tau=0}^{\infty} \beta_t \beta_\tau (1 - \eta \lambda_j)^{|t-\tau|} \right| \leq \|\beta\|_2^2 \sum_{j=1}^d \frac{2}{\eta \lambda_j} < \infty \quad (35)$$

for $\beta \in \ell^2$. Next, using Proposition C.10, $\text{Tr}[\mathbf{H}] \leq R^2$, and $\text{Tr}[\mathbf{P}_\beta] = \langle \beta, \mathbf{T}\beta \rangle$, we get

$$\mathbb{E} \left\| \boldsymbol{\theta}_0^{(0)} \right\|_{\mathbf{H}}^2 = \text{Tr} \left[\mathbf{H} \mathbb{E} \left[\boldsymbol{\theta}_0^{(0)} \otimes \boldsymbol{\theta}_0^{(0)} \right] \right] \leq \eta R^2 \sigma_{\text{sgd}}^2 + \eta \sigma^2 \langle \beta, \mathbf{T}\beta \rangle. \quad (36)$$

Similarly, using Proposition C.11, we get for $r \geq 1$ that

$$\mathbb{E} \left\| \boldsymbol{\theta}_0^{(r)} \right\|_{\mathbf{H}}^2 \leq (\eta R^2)^{r+1} \left(\sigma_{\text{sgd}}^2 + \frac{C_{\text{kurt}} \sigma^2}{R^2} \langle \beta, \mathbf{T}\beta \rangle \right).$$

We can ignore the remainder term since $\mathbb{E} \left\| \boldsymbol{\delta}_t^{(m)} \right\|_{\mathbf{H}}^2 \rightarrow 0$ as $m \rightarrow \infty$, from Proposition C.13. Thus, we get using Property C.7 and the triangle inequality on the norm $\mathbf{u} \mapsto \sqrt{\mathbb{E} \langle \mathbf{u}, \mathbf{H} \mathbf{u} \rangle}$ of a random vector \mathbf{u} to get

$$\sqrt{\mathbb{E} \|\boldsymbol{\theta}'_0\|_{\mathbf{H}}^2} \leq \sum_{r=0}^{\infty} \sqrt{\mathbb{E} \|\boldsymbol{\theta}_0^{(r)}\|_{\mathbf{H}}^2}.$$

To complete the proof, we plug in Equations (35) and (36) and sum up the infinite series. We simplify the result using $\|\mathbf{x} + \mathbf{y}\|_{\mathbf{H}}^2 \leq 2\|\mathbf{x}\|_{\mathbf{H}}^2 + 2\|\mathbf{y}\|_{\mathbf{H}}^2$ and use $F(\boldsymbol{\theta}) - F(\boldsymbol{\theta}_*) = (1/2)\|\boldsymbol{\theta} - \boldsymbol{\theta}_*\|_{\mathbf{H}}^2$. \square

Frequency-domain description: We now state and prove the frequency domain description of the upper bound (13).

Theorem C.15. *Consider the setting of Theorem C.14. If $B \in L^2$, i.e., $\int_{-\pi}^{\pi} |B(\omega)|^2 d\omega < \infty$, we have*

$$\begin{aligned} F_{\infty}(B) &\leq \left(1 + \left(1 - \sqrt{\eta R^2}\right)^{-2}\right) \eta R^2 \sigma_{\text{sgd}}^2 \\ &\quad + \left(1 + C_{\text{kurt}} \left(1 - \sqrt{\eta R^2}\right)^{-2}\right) \frac{\eta^2 G^2 \gamma_{\infty}(B)^2}{2\pi\rho} \int_{-\pi}^{\pi} |B(\omega)|^2 h(\omega) d\omega. \end{aligned}$$

Proof. We again use the shorthand $\sigma^2 = \frac{G^2 \gamma_{\infty}(\boldsymbol{\beta})^2}{2\rho}$. First note that

$$h(\omega) \leq \sum_{j=1}^d \frac{\lambda_j}{1 + (1 - \eta\lambda_j)^2 - 2(1 - \eta\lambda_j)} = \sum_{j=1}^d \frac{1}{\eta^2 \lambda_j} = \frac{\text{Tr}[\mathbf{H}^{-1}]}{\eta^2}.$$

Thus, the right-side is well-defined since

$$\int_{-\pi}^{\pi} |B(\omega)|^2 h(\omega) d\omega \leq \frac{\text{Tr}[\mathbf{H}^{-1}]}{\eta^2} \int_{-\pi}^{\pi} |B(\omega)|^2 d\omega < \infty$$

by assumption. We use Lemma C.4 to get

$$\langle \boldsymbol{\beta}, \mathbf{T}\boldsymbol{\beta} \rangle = \sum_{j=1}^d \langle \boldsymbol{\beta}, \mathbf{T}_j \boldsymbol{\beta} \rangle \leq \sum_{j=1}^d \frac{\eta\lambda_j}{\pi} \int_{-\pi}^{\pi} \frac{|B(\omega)|^2 d\omega}{|1 - \exp(i\omega) - \eta\lambda_j|^2} = \frac{\eta}{\pi} \int_{-\pi}^{\pi} |B(\omega)|^2 h(\omega) d\omega. \quad \square$$

C.3 Proofs of Lower Bounds on the Asymptotic Suboptimality

We now state and prove the lower bound part of (13) on the asymptotic suboptimality.

Assumption C.16. *In addition to Assumption C.2, the data distribution \mathbb{P}_{data} satisfies the following:*

(A2') **Worst-Case Residuals:** *For $(\mathbf{x}, y) \sim \mathbb{P}_{\text{data}}$, the residual $\xi := y - \langle \boldsymbol{\theta}_*, \mathbf{x} \rangle$ has variance $\mathbb{E}[\xi^2] = \sigma_{\text{sgd}}^2$.*

Note that the variance of ξ^2 holds with equality under Assumption C.16.

Theorem C.17. *Suppose Assumption C.16 holds. Consider the sequence $(\boldsymbol{\theta}_t)_{t=-\infty}^{\infty}$ produced by the Noisy-FTRL update in Equation (18) with some given weights $\boldsymbol{\beta} \in \ell^1$. If the learning rate satisfies $\eta < 1/R^2$, we have*

$$F_{\infty}(\boldsymbol{\beta}) \geq \frac{\eta\sigma_{\text{sgd}}^2}{2} \text{Tr}[\mathbf{H}] + \frac{\eta^2 G^2 \gamma_{\infty}(B)^2}{4\pi\rho} \int_{-\pi}^{\pi} |B(\omega)|^2 h(\omega) d\omega \geq \frac{\eta\sigma_{\text{sgd}}^2}{2} \text{Tr}[\mathbf{H}] + \frac{\eta G^2 \gamma_{\infty}(\boldsymbol{\beta})^2}{4\rho} \langle \boldsymbol{\beta}, \mathbf{T}\boldsymbol{\beta} \rangle,$$

where $h(\omega)$ is defined in (14) and \mathbf{T} is defined in (24). Furthermore, the minimal stationary error over all choices of $\boldsymbol{\beta}$ is bounded as

$$\inf_{\boldsymbol{\beta}} F_{\infty}(\boldsymbol{\beta}) \geq \frac{1}{4} \left(2\eta\sigma_{\text{sgd}}^2 + \frac{\eta^2 G^2}{2\rho} \right) \text{Tr}[\mathbf{H}]$$

where the infimum is attained by $\boldsymbol{\beta}_*$ whose DTFT B_* verifies $|B_*(\omega)|^2 = 1/\sqrt{h(\omega)}$.

Note that we assume $\boldsymbol{\beta} \in \ell^1$, i.e., $\|\boldsymbol{\beta}\|_1 = \sum_{\tau=0}^{\infty} |\beta_{\tau}| < \infty$ for technical reasons. This implies that $\boldsymbol{\beta} \in \ell^2$, which we assumed for the upper bounds.

The key idea behind the proof is that the variance of $\boldsymbol{\theta}'_t$ is no smaller than that of an LTI system with $\boldsymbol{x}_t \otimes \boldsymbol{x}_t$ replaced by its expectation \mathbf{H} . We can quantify this latter covariance with equality under Assumption C.16. We set up some notation and develop some preliminary results before proving this theorem.

Formally, consider the sequences $(\boldsymbol{\theta}_t^{(0)})_{t=-\infty}^{\infty}$ and $(\boldsymbol{\delta}_t^{(0)})_{t=-\infty}^{\infty}$ as defined in (29) (cf. §C.2.1). They start at $t = -\infty$ from $\boldsymbol{\theta}_t^{(0)} = \boldsymbol{\theta}'_t$ and $\boldsymbol{\delta}_t^{(0)} = \mathbf{0}$. By Property C.7, we these satisfy $\boldsymbol{\theta}'_t = \boldsymbol{\theta}_t^{(0)} + \boldsymbol{\delta}_t^{(0)}$.

We use a technical result that $\boldsymbol{\theta}_t^{(0)}$ and $\boldsymbol{\delta}_t$ are uncorrelated. It is proved at the end of this section.

Proposition C.18. *Consider the setting of Theorem C.17. We have for all t that*

$$\mathbb{E} \left[\boldsymbol{\theta}_t^{(0)} \otimes \boldsymbol{\delta}_t^{(0)} \right] = \mathbf{0}.$$

We now give the proof of Theorem C.17.

Proof of Theorem C.17. We use shorthand $\sigma^2 = \frac{G^2 \gamma_{\infty}(\boldsymbol{\beta})^2}{2\rho}$. Since $\boldsymbol{\theta}'_t = \boldsymbol{\theta}_t^{(0)} + \boldsymbol{\delta}_t^{(0)}$, we have

$$\mathbb{E} [\boldsymbol{\theta}'_t \otimes \boldsymbol{\theta}'_t] = \mathbb{E} \left[\boldsymbol{\theta}_t^{(0)} \otimes \boldsymbol{\theta}_t^{(0)} \right] + \mathbb{E} \left[\boldsymbol{\delta}_t^{(0)} \otimes \boldsymbol{\delta}_t^{(0)} \right] \succeq \mathbb{E} \left[\boldsymbol{\theta}_t^{(0)} \otimes \boldsymbol{\theta}_t^{(0)} \right] \quad (37)$$

where the cross terms disappear from Proposition C.18 for the first equality. We can get an expression for this term by following the proof of Proposition C.10: under Assumption C.16, we have that Equation (32) holds with equality. Thus, we get for all $t > -\infty$ that

$$\begin{aligned} F_{\infty}(B) &= \text{Tr} [\mathbf{H} \mathbb{E} [\boldsymbol{\theta}'_t \otimes \boldsymbol{\theta}'_t]] \succeq \text{Tr} \left[\mathbf{H} \mathbb{E} \left[\boldsymbol{\theta}_t^{(0)} \otimes \boldsymbol{\theta}_t^{(0)} \right] \right] \\ &= \frac{1}{2\pi} \int_{-\pi}^{\pi} \left(\eta^2 \sigma_{\text{sgd}}^2 \text{Tr} \left[\mathbf{H}^{1/2} \mathbf{M}_{\omega} \mathbf{H} \mathbf{M}_{\omega}^* \mathbf{H}^{1/2} \right] + \eta^2 \sigma^2 |B(\omega)|^2 \text{Tr} \left[\mathbf{H}^{1/2} \mathbf{M}_{\omega} \mathbf{M}_{\omega}^* \mathbf{H}^{1/2} \right] \right) d\omega. \end{aligned} \quad (38)$$

We invoke Corollary C.5 to obtain

$$\begin{aligned} \int_{-\pi}^{\pi} \text{Tr} \left[\mathbf{H}^{1/2} \mathbf{M}_{\omega} \mathbf{H} \mathbf{M}_{\omega}^* \mathbf{H}^{1/2} \right] d\omega &= \sum_{j=1}^d \lambda_j^2 \int_{-\pi}^{\pi} \frac{d\omega}{|1 - \exp(i\omega) - \eta\lambda_j|^2} \\ &\geq \sum_{j=1}^d \frac{\pi\lambda_j}{\eta} = \frac{\pi}{\eta} \text{Tr} [\mathbf{H}]. \end{aligned}$$

Similarly, we invoke Lemma C.4 to compute

$$\begin{aligned} \int_{-\pi}^{\pi} |B(\omega)|^2 \text{Tr} \left[\mathbf{H}^{1/2} \mathbf{M}_{\omega} \mathbf{M}_{\omega}^* \mathbf{H}^{1/2} \right] d\omega &= \int_{-\pi}^{\pi} \left(\sum_{j=1}^d |B(\omega)|^2 \frac{\lambda_j}{|1 - \exp(i\omega) - \eta\lambda_j|^2} \right) d\omega \\ &= \int_{-\pi}^{\pi} |B(\omega)|^2 h(\omega) d\omega \geq \frac{\pi}{\eta} \langle \boldsymbol{\beta}, \mathbf{T} \boldsymbol{\beta} \rangle. \end{aligned}$$

This establishes the lower bound for specific choices of $\boldsymbol{\beta}$.

Now, we turn to the universal lower bound. Using the expression for $\gamma_{\infty}(B)$ from Property C.1, we get that the lower bound from the theorem statement is

$$F_{\infty}(B) \geq \frac{\eta\sigma_{\text{sgd}}^2}{2} \text{Tr} [\mathbf{H}] + \frac{\eta^2 G^2}{8\pi^2 \rho} \left(\int_{-\pi}^{\pi} \frac{d\omega}{|B(\omega)|^2} \right) \left(\int_{-\pi}^{\pi} |B(\omega)|^2 h(\omega) \right). \quad (39)$$

The Cauchy-Schwarz inequality gives us that

$$\left(\int_{-\pi}^{\pi} \frac{d\omega}{|B(\omega)|^2} \right) \left(\int_{-\pi}^{\pi} |B(\omega)|^2 h(\omega) \right) \geq \left(\int_{-\pi}^{\pi} \sqrt{h(\omega)} d\omega \right)^2,$$

with equality attained for $|B(\omega)|^2 = 1/\sqrt{h(\omega)}$. This gives the universal lower bound on (39) over all possible choices of B (or equivalently, all possible choices of β). To further lower bound this, we use $\cos(\omega) \geq -1$ to get

$$h(\omega) = \sum_{j=1}^d \frac{\lambda_j}{1 + (1 - \eta\lambda_j)^2 - 2(1 - \eta\lambda_j)\cos(\omega)} \geq \sum_{j=1}^d \frac{\lambda_j}{(2 - \eta\lambda_j)^2} \geq \frac{1}{4} \sum_{j=1}^d \lambda_j = \frac{\text{Tr}[\mathbf{H}]}{4}.$$

Thus, we get that (39) can be further lower bounded as

$$F_\infty(B) \geq \frac{\eta\sigma_{\text{sgd}}^2}{2} \text{Tr}[\mathbf{H}] + \frac{\eta^2 G^2}{8\pi^2 \rho} \left(\int_{-\pi}^{\pi} \frac{\sqrt{\text{Tr}[\mathbf{H}]}}{2} d\omega \right)^2 = \frac{\eta\sigma_{\text{sgd}}^2}{2} \text{Tr}[\mathbf{H}] + \frac{\eta^2 G^2}{8\rho} \text{Tr}[\mathbf{H}].$$

□

Missing technical proofs in the lower bound: We now give the proof of Proposition C.18, which first relies on the following intermediate result.

Proposition C.19. *Consider the setting of Theorem C.17. We have for all t, τ that*

$$\mathbb{E} \left[\mathbf{w}_\tau \otimes \delta_t^{(0)} \right] = \mathbf{0}.$$

Proof. For this proof, we start the sequences at $t = 0$ rather than $t = -\infty$. We drop the superscript to write $\delta_t^{(0)}$ as δ_t . Define shorthand $\mathbf{Q}_t := \mathbf{I} - \eta \mathbf{x}_t \otimes \mathbf{x}_t$ and $\mathbf{R}_t := \mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t$. We expand out the recursion to get

$$\begin{aligned} \delta_t &= \mathbf{Q}_{t-1} \delta_{t-1} + \eta \mathbf{R}_{t-1} \theta_{t-1}^{(0)} \\ &= \mathbf{Q}_{t-1} (\mathbf{Q}_{t-2} \delta_{t-2} + \eta \mathbf{R}_{t-2} \theta_{t-2}^{(0)}) + \eta \mathbf{R}_{t-1} \theta_{t-1}^{(0)} \\ &= \mathbf{Q}_{t-1} \mathbf{Q}_{t-2} \cdots \mathbf{Q}_0 \delta_0 + \eta \left(\mathbf{R}_{t-1} \theta_{t-1}^{(0)} + \mathbf{Q}_{t-1} \mathbf{R}_{t-2} \theta_{t-2}^{(0)} + \cdots + \mathbf{Q}_{t-1} \cdots \mathbf{Q}_1 \mathbf{R}_0 \theta_0^{(0)} \right). \end{aligned}$$

The first term is zero because $\delta_0 = \mathbf{0}$ at initialization. Since \mathbf{R}_τ is mean zero and independent of $\theta_\tau^{(0)}$ and \mathbf{R}_t for $t > \tau$, we have

$$\begin{aligned} \frac{1}{\eta} \mathbb{E}[\delta_t \otimes \mathbf{w}_\tau] &= \mathbb{E}[\mathbf{R}_{t-1}] \mathbb{E} \left[\theta_{t-1}^{(0)} \otimes \mathbf{w}_\tau \right] \\ &\quad + \mathbb{E}[\mathbf{Q}_{t-1}] \mathbb{E}[\mathbf{R}_{t-2}] \mathbb{E} \left[\theta_{t-2}^{(0)} \otimes \mathbf{w}_\tau \right] + \cdots + \mathbb{E}[\mathbf{Q}_{t-1} \cdots \mathbf{Q}_1] \mathbb{E}[\mathbf{R}_0] \mathbb{E} \left[\theta_0^{(0)} \otimes \mathbf{w}_\tau \right] \\ &= \mathbf{0}, \end{aligned}$$

giving us the desired result. □

Proof of Proposition C.18. We drop the superscript to write $\delta_t^{(0)}$ as δ_t . We prove the claim by induction. At initialization, we have $\delta_{-\infty} = \mathbf{0}$ so the hypothesis holds. Now assume that it holds at time t , i.e., $\mathbb{E} \left[\theta_t^{(0)} \otimes \delta_t \right] = \mathbf{0}$.

Next, we expand out $\mathbb{E} \left[\theta_{t+1}^{(0)} \otimes \delta_{t+1} \right]$ using their respective recursions. Note that $\mathbf{w}_t, \mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t$ and ξ_t are each zero mean and independent of all quantities appearing up to iteration t (formally, they are independent of the σ -algebra generated by $(\theta_t^{(0)}$ and δ_t). This gives

$$\mathbb{E} \left[\theta_{t+1}^{(0)} \otimes \delta_{t+1} \right] = (\mathbf{I} - \eta \mathbf{H}) \mathbb{E} \left[\theta_t^{(0)} \otimes \delta_t \right] (\mathbf{I} - \eta \mathbf{H}) - \eta \mathbb{E} \left[\sum_{\tau=0}^{\infty} \beta_\tau \left(\mathbf{w}_{t-\tau} \otimes \delta_t^{(0)} \right) \right] (\mathbf{I} - \eta \mathbf{H}). \quad (40)$$

The first term is zero by the induction hypothesis. For the second term, we can interchange the expectation and the infinite sum by the Fubini-Tonelli theorem since

$$\sum_{\tau=0}^{\infty} |\beta_\tau| \mathbb{E} \left| \left\langle \mathbf{w}_{t-\tau}, \delta_t^{(0)} \right\rangle \right| \leq \|\beta\|_1 \max_{\tau=0, \dots, \infty} \mathbb{E} \left| \left\langle \mathbf{w}_{t-\tau}, \delta_t^{(0)} \right\rangle \right| < \infty$$

since $\beta_1 \in \ell^1$ and $\mathbb{E} \left| \left\langle \mathbf{w}_{t-\tau}, \delta_t^{(0)} \right\rangle \right| < \infty$ because

$$\mathbb{E} \left\langle \mathbf{w}_{t-\tau}, \delta_t^{(0)} \right\rangle = \text{Tr} \left[\mathbb{E} \left[\mathbf{w}_{t-\tau} \otimes \delta_t^{(0)} \right] \right] = 0$$

by Proposition C.19. By Proposition C.19 again, we thus get

$$\mathbb{E} \left[\sum_{\tau=0}^{\infty} \beta_{\tau} \left(\mathbf{w}_{t-\tau} \otimes \delta_t^{(0)} \right) \right] = \sum_{\tau=0}^{\infty} \beta_{\tau} \mathbb{E} \left[\left(\mathbf{w}_{t-\tau} \otimes \delta_t^{(0)} \right) \right] = \mathbf{0}.$$

□

C.4 Asymptotics of ν -Noisy-FTRL

We now state and prove the upper bound for ν -Noisy-FTRL. Note that ν -Noisy-FTRL can be described in the frequency domain as $|\hat{B}^{\nu}(\omega)|^2 = |1 - \nu - \exp(i\omega)|^2$.

Proposition C.20. *Consider the setting of Theorem C.15 with $\sigma_{\text{sgd}}^2 = 0$. Then, ν -Noisy-FTRL with $\nu \leq \eta\mu$ satisfies*

$$F_{\infty}(\hat{\beta}^{\nu}) \leq C \max\{1, C_{\text{kurt}}\} \eta^2 G^2 \rho^{-1} \text{Tr}[\mathbf{H}] \log^2 \left(\frac{8}{\nu} \right) + \tilde{O}(\eta^3 R^2 \mu G^2 \rho^{-1}),$$

for a universal constant $C > 0$, and $\tilde{O}(\cdot)$ suppresses polylogarithmic terms in the problem parameters.

Proof. We use C to denote a universal constant that can change from line to line. Denote $\mathcal{I}(a, b)$ for $0 < a, b < 1$ as the integral

$$\mathcal{I}(a, b) := \int_{-\pi}^{\pi} \frac{|1 - a - \exp(i\omega)|}{|1 - b - \exp(i\omega)|^2} d\omega.$$

We can express the bound of Theorem C.15 with our specific choice of $B(\omega)$ as

$$F_{\infty}(\hat{B}^{\nu}) \leq C \max\{1, C_{\text{kurt}}\} \mathcal{I}(\nu, \nu) \sum_{j=1}^d \lambda_j \mathcal{I}(\nu, \eta\lambda_j). \quad (41)$$

The strategy is to reduce each \mathcal{I} term to standard elliptic integrals and leverage their well-studied properties to get the result. We start with the first term. We use Lemma F.15 to rewrite in terms of the elliptic integral of the first kind $K(k) = \int_0^{\pi/2} d\omega / \sqrt{1 - k^2 \sin^2(\omega)}$ (denoted as (a)). Then, we use Property F.10 which says that $K(k) = O(-\log \sqrt{1 - k^2})$ (denoted as (b)). This gives,

$$\mathcal{I}(\nu, \nu) \stackrel{(a)}{=} \frac{4}{2 - \nu} K \left(\frac{\sqrt{1 - \nu}}{1 - \nu/2} \right) \stackrel{(b)}{\leq} \frac{5}{2 - \nu} \log \left(\frac{4}{\nu} (2 - \nu) \right) \leq 5 \log \left(\frac{8}{\nu} \right). \quad (42)$$

Similarly, we can express the second integral in terms of the elliptic integral of the third kind $\Pi(\alpha^2, k)$, whose definition is given in (88). From Lemma F.16, we have for $a, b \in (0, 1)$ that

$$\mathcal{I}(a, b) = \frac{2a^2}{b^2(1 - a/2)} \Pi(\alpha^2, k) \quad \text{where} \quad \alpha^2 = \frac{b^2(1 - a) - a^2(1 - b)}{b^2(1 - a/2)^2}$$

and $k = \sqrt{1 - a}/(1 - a/2)$. We invoke Property F.11 to bound the behavior of $\Pi(\alpha^2, k)$ as $k \rightarrow 1^-$ (i.e. $a \rightarrow 0^+$) to get

$$\begin{aligned} \mathcal{I}(a, b) &\leq \frac{2a^2}{b^2(1 - a/2)} \frac{1}{\sqrt{1 - \alpha^2}} \log \frac{4}{\sqrt{1 - k^2}} (1 + O(a)) \\ &= \frac{2(1 - a/2)}{(1 - b/2)^2} \log \left(\frac{4}{a} (2 - a) \right) (1 + O(a)) \leq \frac{128}{49} \log(8/a) (1 + O(a)), \end{aligned}$$

where the last inequality holds for $a \leq b \leq 1/4$. We plug in $a = \nu$ and $b = \eta\lambda_j$ so these conditions are satisfied to get

$$\mathcal{I}(\eta\mu, \eta\lambda_j) \leq C \log \left(\frac{8}{\nu} \right) (1 + O(\nu)). \quad (43)$$

The last term is $O(\nu) \leq O(\eta\mu)$. Plugging in (42) and (43) into (41) and using $\text{Tr}[\mathbf{H}] = \sum_{j=1}^n \lambda_j \leq R^2$ completes the proof. □

C.5 Asymptotics of Anti-PGD

As we discussed in Table 3, anti-PGD [44] is a special case of Noisy-FTRL with $\beta = (1, -1, 0, \dots)$. Then, we have that $(\text{Toeplitz}(\beta))^{-1}$ is the lower triangular matrix of all ones, so we have $\gamma_T(\beta) = T$, or that its limiting sensitivity is infinite.

We can circumvent the infinity by damping $\beta = (1, -(1-\nu), 0, \dots)$ for some $0 < \nu < 1$ to be decided later. In this case, we have $B(\omega) = 1 - (1-\nu)\exp(-i\omega)$, so that $|B(\omega)|^2 = |1 - \nu - \exp(i\omega)|^2$, which is the analogue of ν -Noisy-FTRL with a square.

Proposition C.21. *Consider the setting of Theorem C.15 with $\sigma_{\text{sgd}}^2 = 0$ and $\beta = (1, -(1-\eta\lambda), 0, \dots)$ for some $\lambda \in (0, 1/\eta]$. Then, we have,*

$$F_\infty(\beta) = \Theta \left(\eta G^2 \rho^{-1} \left(\nu d + \frac{\eta \text{Tr}[\mathbf{H}]}{\nu} \right) \right).$$

Further, if the learning rate satisfies $\eta = c/\text{Tr}[\mathbf{H}]$ and we take $\beta = (1, -(1 - \sqrt{1/d}), \dots)$, we get

$$F_\infty(\beta) = \Theta \left((c^{1/2} + c^{-1/2}) \eta^{3/2} \sigma^2 \sqrt{d \text{Tr}[\mathbf{H}]} \right).$$

Proof. Let $\sigma^2 = G^2/(2\rho)$. From Theorems C.15 and C.17, we get that

$$F_\infty(\beta) = \Theta \left(\eta^2 \sigma^2 \left(\int_{-\pi}^{\pi} \frac{d\omega}{|1 - \nu - \exp(i\omega)|^2} \right) \left(\sum_{j=1}^d \lambda_j \int_{-\pi}^{\pi} \frac{|1 - \nu - \exp(i\omega)|^2}{|1 - \eta\lambda_j - \exp(i\omega)|^2} d\omega \right) \right). \quad (44)$$

Using Lemma F.12, we have

$$\int_{-\pi}^{\pi} \frac{d\omega}{|1 - \nu - \exp(i\omega)|^2} = \frac{2\pi}{\nu(2-\nu)} = \Theta \left(\frac{1}{\nu} \right).$$

For the second integral, we expand out the numerator and invoke Lemma F.12 again to get

$$\begin{aligned} \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{|1 - \nu - \exp(i\omega)|^2}{|1 - \eta\lambda_j - \exp(i\omega)|^2} d\omega &= \frac{1 + (1-\nu)^2}{\eta\lambda_j(2-\eta\lambda_j)} - 2(1-\nu) \frac{1 - \eta\lambda_j}{\eta\lambda_j(2-\eta\lambda_j)} \\ &= \Theta \left(\frac{\nu^2}{\eta\lambda_j} + 1 \right), \end{aligned}$$

where we use $1 \leq 2 - \nu \leq 2$ and the same for λ_j instead of λ . Plugging the two integrals back into (44) completes the proof. \square

C.6 Proofs of Technical Lemmas

We now prove Lemma C.4.

Proof of Lemma C.4. Denote

$$I = \int_{-\pi}^{\pi} \frac{|B(\omega)|^2 d\omega}{|1 - \eta\lambda_j - \exp(i\omega)|^2}.$$

The denominator is simply

$$|1 - \exp(i\omega) - \eta\lambda_j|^2 = 1 + (1 - \eta\lambda_j)^2 - 2(1 - \eta\lambda_j) \cos \omega. \quad (45)$$

We expand the numerator as

$$\begin{aligned} |B(\omega)|^2 &= \sum_{t=0}^{\infty} \beta_t^2 + \sum_{t=0}^{\infty} \sum_{\tau=0}^{t-1} \beta_t \beta_\tau (\exp(i\omega(t-\tau)) + \exp(-i\omega(t-\tau))) \\ &= \sum_{t=0}^{\infty} \beta_t^2 + 2 \sum_{t=0}^{\infty} \sum_{\tau=0}^{t-1} \beta_t \beta_\tau \cos(\omega(t-\tau)) \\ &= \sum_{t=0}^{\infty} \sum_{\tau=0}^{\infty} \beta_t \beta_\tau \cos(\omega(t-\tau)). \end{aligned} \quad (46)$$

This is bounded since the Cauchy-Schwarz inequality gives

$$|B(\omega)|^2 \leq \|\beta\|_2^2 < \infty.$$

Thus, we can apply Fubini's theorem to exchange the sum and integral to give

$$\begin{aligned} I &= \sum_{t=0}^{\infty} \sum_{\tau=0}^{\infty} \beta_t \beta_\tau \int_{-\pi}^{\pi} \frac{\cos(\omega(t-\tau)) d\omega}{1 + (1-\eta\lambda_j)^2 - 2(1-\eta\lambda_j)\cos(\omega)} \\ &= \sum_{t=0}^{\infty} \sum_{\tau=0}^{\infty} \frac{2\pi}{1 - (1-\eta\lambda_j)^2} (1-\eta\lambda_j)^{|t-\tau|} = \frac{2\pi \langle \beta, \mathbf{T}_j \beta \rangle}{\eta\lambda_j(2-\eta\lambda_j)}, \end{aligned}$$

where we evaluated the integral using Lemma F.12. We use $1 \leq 2 - \eta\lambda_j \leq 2$ to complete the proof. \square

D Finite-Time Privacy-Utility Tradeoffs for Linear Regression

The goal of this section is to establish the finite time convergence of DP-FTRL. The key idea of the proof is to establish high probability bounds on the ℓ_2 norm of the iterates of Noisy-FTRL and use that to deduce a clip norm that does not clip any gradients with high probability.

Outline: The outline of the rest of this section is as follows:

- **Appendix D.1:** Preliminaries, including setup, notation and assumptions.
- **Appendix D.2:** High probability bounds the iterates of Noisy-FTRL.
- **Appendix D.3:** Expected bounds on the iterates of Noisy-FTRL.
- **Appendix D.4:** Connecting DP-FTRL to Noisy-FTRL for the final bound privacy-utility bounds (Corollary D.14 for DP-SGD and Corollary D.15 for DP-FTRL).

D.1 Setup, Assumptions and Notation

In this section, we make precise the assumptions, notation, and give some preliminary results.

D.1.1 Assumptions

We make the following assumptions throughout this section.

Assumption D.1. *The data distribution \mathbb{P}_{data} satisfies the following:*

- (A1) **Input Distribution:** *The inputs have mean $\mathbb{E}[\mathbf{x}] = \mathbf{0}$ and covariance $\mathbb{E}[\mathbf{x} \otimes \mathbf{x}] =: \mathbf{H}$. We have $\mu \mathbf{I} \preceq \mathbf{H} \preceq L \mathbf{I}$ for $\mu, L > 0$. Further, $\mathbf{H}^{-1/2} \mathbf{x}$ is element-wise independent and sub-Gaussian with variance proxy 1, e.g. $\mathbf{H}^{-1/2} \mathbf{x} \sim \mathcal{N}(0, \mathbf{I})$.*
- (A2) **Noise Distribution:** *There exists a $\theta_* \in \mathbb{R}^d$ such that $y = \langle \theta_*, \mathbf{x} \rangle + \xi$, where ξ is independent of \mathbf{x} and is zero-mean sub-Gaussian with variance proxy σ_{sgd}^2 , e.g. $\xi \sim \mathcal{N}(0, \sigma_{\text{sgd}}^2)$.*

Not under Assumption (A2) that θ_* is the global minimizer of $F(\theta) = (1/2)\mathbb{E}(\langle \theta, \mathbf{x} \rangle - y)^2$.

D.1.2 Notation

- As in Assumption C.2, we denote R^2 as the smallest number such that the fourth moment of \mathbf{x} is bounded as

$$\mathbb{E} \left[\|\mathbf{x}\|_2^2 \mathbf{x} \otimes \mathbf{x} \right] \preceq R^2 \mathbf{H}. \quad (47)$$

Under Assumption (A1), we have $R^2 = \Theta(\text{Tr}[\mathbf{H}])$ always. While $\text{Tr}[\mathbf{H}] \leq R^2$ directly follows from (47) using Jensen's inequality, we show that $R^2 \leq 3\text{Tr}[\mathbf{H}]$ in Property C.3 in Appendix C.1.

- It is convenient to rewrite the Noisy-FTRL recursion (16) in terms of the difference $\theta'_t := \theta_t - \theta_*$ as

$$\theta'_{t+1} = (\mathbf{I} - \eta(\mathbf{x}_t \otimes \mathbf{x}_t))\theta'_t + \eta \xi_t \mathbf{x}_t - \eta \sum_{\tau=0}^t \beta_\tau \mathbf{w}_{t-\tau}. \quad (48)$$

We will show in the upcoming Property D.2 that $\theta'_t = \hat{\theta}_t + \tilde{\theta}^{\text{sgd}} + \tilde{\theta}^{\text{dp}}$, where $\hat{\theta}_t$ captures the effect of the initial iterate, $\tilde{\theta}^{\text{sgd}}$ captures the effect of the SGD noise, and $\tilde{\theta}^{\text{dp}}$ captures the effect of the additive DP noise. We will define these quantities now and state and prove Property D.2 later. Note that these recursions are defined for the same sequences of input realizations $(\mathbf{x}_0, \mathbf{x}_1, \dots)$ drawn from \mathbb{P}_{data} , linear model noise realizations (ξ_0, ξ_1, \dots) , and DP noise realizations $(\mathbf{w}_0, \mathbf{w}_1, \dots)$.

- We define the noise-free version of the DP-FTRL recursion as $\hat{\theta}_0 = \theta'_0$ and

$$\hat{\theta}_{t+1} = (\mathbf{I} - \eta(\mathbf{x}_t \otimes \mathbf{x}_t))\hat{\theta}_t. \quad (49)$$

- The effect of the SGD noise in the Noisy-FTRL process can be quantified by creating a process starting from $\tilde{\theta}_0^{\text{sgd}} = \mathbf{0}$ with no DP noise (i.e. $\mathbf{w}_\tau \equiv \mathbf{0}$):

$$\tilde{\theta}_{t+1}^{\text{sgd}} = (\mathbf{I} - \eta(\mathbf{x}_t \otimes \mathbf{x}_t))\tilde{\theta}_t^{\text{sgd}} + \eta \xi_t \mathbf{x}_t. \quad (50)$$

- The effect of the DP noise in the Noisy-FTRL process can be quantified by creating a process starting from $\tilde{\theta}_0^{\text{dp}} = \mathbf{0}$ with no SGD noise (i.e., $\xi_t \equiv 0$):

$$\tilde{\theta}_{t+1}^{\text{dp}} = (\mathbf{I} - \eta(\mathbf{x}_t \otimes \mathbf{x}_t))\tilde{\theta}_t^{\text{dp}} - \eta \sum_{\tau=0}^t \beta_\tau \mathbf{w}_{t-\tau}. \quad (51)$$

- For an input \mathbf{x}_t drawn from \mathbb{P}_{data} We define the matrix

$$\mathbf{Q}_t := \mathbf{I} - \eta \mathbf{x}_t \otimes \mathbf{x}_t. \quad (52)$$

Note that $\mathbb{E}[\mathbf{Q}_t] = \mathbf{I} - \eta \mathbf{H}$.

- Define the linear operator $\mathcal{P} : \mathbb{S}_+^d \rightarrow \mathbb{S}_+^d$ that operates on the cone of PSD matrices given by

$$\mathcal{P}M = \mathbb{E}[(\mathbf{I} - \eta \mathbf{x} \otimes \mathbf{x})M(\mathbf{I} - \eta \mathbf{x} \otimes \mathbf{x})], \quad (53)$$

where \mathbf{x} is an input drawn from \mathbb{P}_{data} . By definition, we have $\mathbb{E}[\mathbf{Q}_t M \mathbf{Q}_t] = \mathcal{P}M$ and by independence,

$$\mathbb{E}[\mathbf{Q}_t \mathbf{Q}_{t-1} M \mathbf{Q}_{t-1} \mathbf{Q}_t] = \mathcal{P}(\mathcal{P}M) = \mathcal{P}^2 M. \quad (54)$$

This extends to higher powers of \mathcal{P} as well. Finally, we will heavily use the fact that $\text{Tr}[\mathcal{P}M] \leq (1 - \eta\mu)\text{Tr}[M]$ for PSD matrices M (see Lemma F.18 for a proof).

- For each iteration t , we define the PSD matrix Σ_t^{sgd} as

$$\Sigma_t^{\text{sgd}} = \mathbf{x}_{t-1} \otimes \mathbf{x}_{t-1} + \mathbf{Q}_{t-1}(\mathbf{x}_{t-2} \otimes \mathbf{x}_{t-2})\mathbf{Q}_{t-1} + \dots + \mathbf{Q}_{t-1} \cdots \mathbf{Q}_1(\mathbf{x}_0 \otimes \mathbf{x}_0)\mathbf{Q}_1 \cdots \mathbf{Q}_{t-1}, \quad (55)$$

- For each iteration t , we define the PSD matrix Σ_t^{dp} as

$$\Sigma_t^{\text{dp}} = \sum_{\tau=0}^{t-1} \mathbf{V}_{t,\tau} \mathbf{V}_{t,\tau}^\top \quad \text{where} \quad (56)$$

$$\mathbf{V}_{t,\tau} = \begin{cases} \beta_\tau \mathbf{I} + \beta_{\tau-1} \mathbf{Q}_{t-1} + \dots + \beta_0 \mathbf{Q}_{t-1} \cdots \mathbf{Q}_{t-\tau}, & \text{if } 1 \leq \tau \leq t-1, \\ \beta_0 \mathbf{I}, & \text{if } \tau = 0. \end{cases}$$

D.1.3 Preliminary Results

The first result is a decomposition of the Noisy-FTRL process in to three processes: (a) gradient descent without additive noise, (b) a noise process with only noise from the linear model, and (c) a noise process with only the DP noise.

Property D.2. For the sequences $\theta'_t, \hat{\theta}_t, \tilde{\theta}_t^{\text{sgd}}, \tilde{\theta}_t^{\text{dp}}$ defined in Equations (48) to (51), we have the following:

$$\theta'_t = \hat{\theta}_t + \tilde{\theta}_t^{\text{sgd}} + \tilde{\theta}_t^{\text{dp}} \quad (57)$$

$$\hat{\theta}_t = \mathbf{Q}_t \cdots \mathbf{Q}_0 \theta'_0 \quad (58)$$

$$\tilde{\theta}_t^{\text{sgd}} = \eta (\mathbf{x}_t \xi_t + \mathbf{Q}_t \mathbf{x}_{t-1} \xi_{t-1} + \cdots + \mathbf{Q}_t \cdots \mathbf{Q}_1 \mathbf{x}_0 \xi_0) \quad (59)$$

$$\begin{aligned} \tilde{\theta}_t^{\text{dp}} &= -\eta \left(\sum_{\tau=0}^t \beta_\tau \mathbf{w}_{t-\tau} + \mathbf{Q}_t \sum_{\tau=0}^{t-1} \beta_\tau \mathbf{w}_{t-1-\tau} + \cdots + \mathbf{Q}_t \cdots \mathbf{Q}_1 (\beta_0 \mathbf{w}_0) \right) \\ &= -\eta \left(\beta_0 \mathbf{w}_{t-1} + (\beta_1 \mathbf{I} + \beta_0 \mathbf{Q}_{t-1}) \mathbf{w}_{t-2} + \cdots + (\beta_{t-1} \mathbf{I} + \beta_{t-2} \mathbf{Q}_{t-1} + \cdots + \beta_0 \mathbf{Q}_{t-1} \cdots \mathbf{Q}_1) \mathbf{w}_0 \right). \end{aligned} \quad (60)$$

Proof. The expressions follow from unrolling their respective updates. By unrolling the DP-FTRL update (48), we get,

$$\begin{aligned} \theta'_{t+1} &= \mathbf{Q}_t \theta'_t + \eta \mathbf{x}_t \xi_t - \eta \sum_{\tau=0}^t \beta_\tau \mathbf{w}_{t-\tau} \\ &= \mathbf{Q}_t \mathbf{Q}_{t-1} \theta'_{t-1} + \eta (\mathbf{x}_t \xi_t + \mathbf{Q}_t \mathbf{x}_{t-1} \xi_{t-1}) - \eta \left(\sum_{\tau=0}^t \beta_\tau \mathbf{w}_{t-\tau} + \mathbf{Q}_t \sum_{\tau=0}^{t-1} \beta_\tau \mathbf{w}_{t-1-\tau} \right) \\ &= \mathbf{Q}_t \cdots \mathbf{Q}_0 \theta'_0 + \eta (\mathbf{x}_t \xi_t + \mathbf{Q}_t \mathbf{x}_{t-1} \xi_{t-1} + \cdots + \mathbf{Q}_t \cdots \mathbf{Q}_1 \mathbf{x}_0 \xi_0) \\ &\quad - \eta \left(\sum_{\tau=0}^t \beta_\tau \mathbf{w}_{t-\tau} + \mathbf{Q}_t \sum_{\tau=0}^{t-1} \beta_\tau \mathbf{w}_{t-1-\tau} + \cdots + \mathbf{Q}_t \cdots \mathbf{Q}_1 (\beta_0 \mathbf{w}_0) \right). \end{aligned}$$

Unrolling Equations (49) to (51) respectively gives Equations (58) to (60), and comparing them with the expression above gives Equation (57). \square

D.2 High-Probability Bounds on Noisy-FTRL

The goal of this subsection is to prove a high probability bound on norms of the iterates of Noisy-FTRL. We require a technical convergence condition on the weights β .

Definition D.3. A sequence $\beta = (\beta_0, \beta_1, \dots)$ is said to satisfy *Half-Expo Decay* with parameter $\nu \in (0, 1)$ if for all nonnegative integers τ , we have

$$|\beta_0| (1-\nu)^{\tau/2} + |\beta_1| (1-\nu)^{(\tau-1)/2} + \cdots + |\beta_\tau| \leq C (1-\nu)^{\tau/2} \quad (61)$$

for a universal constant $C > 0$.

Theorem D.4. Fix a constant $0 < p < 1$ and suppose the Assumption D.1 holds. Consider the sequence $(\theta_t)_{t=0}^{T-1}$ of iterates and the sequence $(\mathbf{g}_t)_{t=0}^{T-1}$ of gradients when running Noisy-FTRL for T iterations with noise coefficients $\beta = (\beta_0, \dots, \beta_{T-1})$, DP noise $\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ of a given variance⁴ σ^2 , a learning rate $\eta \leq (cR^2 \log(T/p))$ for a universal constant $c \geq 1$. Further, suppose that β satisfies *Half-Expo Decay* with parameter ν for some $\nu \leq \eta\mu$. Then, with probability at least $1 - p$, we have

$$\begin{aligned} \|\theta'_t\|_2^2 &\leq C \left(\|\theta'_0\|_2^2 + \frac{\eta R^2 \sigma_{\text{sgd}}^2}{\mu} + \frac{\eta^2 \sigma^2 d \|\beta\|_1^2}{\nu} \right) \log^3 \left(\frac{T}{p} \right) \quad \text{and} \\ \|\mathbf{g}_t\|_2^2 &\leq CR^4 \left(\|\theta'_0\|_2^2 + \frac{\eta R^2 \sigma_{\text{sgd}}^2}{\mu} + \frac{\sigma_{\text{sgd}}^2}{R^2} + \frac{\eta^2 \sigma^2 d \|\beta\|_1^2}{\nu} \right) \log^5 \left(\frac{T}{p} \right). \end{aligned}$$

for a universal constant C .

We prove this theorem over a sequence of intermediate results.

⁴In the context of this paper, we have $\sigma^2 = G^2 \gamma(\beta)^2 / (2\rho)$.

D.2.1 Proof Setup: Definition of Events

The proof strategy relies on defining some events (that hold with high probability from concentration of measure) and prove the required boundedness under those events. Consider $0 < p < 1$ and a universal constant C from statement of Theorem D.4. We define the following events.

- Define the event where the inputs are bounded in norm as:

$$\mathcal{E}_1 := \bigcap_{t=0}^{T-1} \left\{ \|\mathbf{x}_t\|_2^2 \leq CR^2 \log\left(\frac{T}{p}\right) \right\}. \quad (62)$$

- Define an event where the noise in the linear model is bounded as:

$$\mathcal{E}_2 := \bigcap_{t=0}^{T-1} \left\{ |\xi_t|^2 \leq 2\sigma_{\text{sgd}}^2 \log\left(\frac{2T}{p}\right) \right\}. \quad (63)$$

- Define the event where the norm of $\tilde{\boldsymbol{\theta}}^{\text{sgd}}$ defined in (50) is bounded

$$\mathcal{E}_1^{\text{dp}} := \bigcap_{t=0}^{T-1} \left\{ \|\tilde{\boldsymbol{\theta}}^{\text{sgd}}\|_2^2 \leq C\eta^2 \sigma_{\text{sgd}}^2 \text{Tr}[\boldsymbol{\Sigma}_t^{\text{sgd}}] \log\left(\frac{T}{p}\right) \right\}, \quad (64)$$

where we define the random matrix $\boldsymbol{\Sigma}_t^{\text{sgd}} = \mathbf{x}_{t-1} \otimes \mathbf{x}_{t-1} + \mathbf{Q}_{t-1}(\mathbf{x}_{t-2} \otimes \mathbf{x}_{t-2})\mathbf{Q}_{t-1} + \dots + \mathbf{Q}_{t-1} \dots \mathbf{Q}_1(\mathbf{x}_0 \otimes \mathbf{x}_0)\mathbf{Q}_1 \dots \mathbf{Q}_{t-1}$ (see also (55)). When this event holds, we have that $\mathbf{0} \preceq \mathbf{Q}_t \preceq \mathbf{I}$ for $t = 0, \dots, T-1$ as long as $\eta \leq 1/(CR^2 \log(T/p))$. Indeed, in this case, we have

$$\mathbf{I} - \eta \mathbf{x}_t \otimes \mathbf{x}_t \succeq \left(1 - \eta \|\mathbf{x}_t\|_2^2\right) \mathbf{I} \succeq \mathbf{0}. \quad (65)$$

- The components of the sum defining $\boldsymbol{\Sigma}_t^{\text{sgd}}$ are the PSD matrices $\mathbf{W}_{t,\tau}$, defined for $\tau \leq t-1$ as

$$\mathbf{W}_{t,\tau} = \begin{cases} \mathbf{Q}_{t-1} \dots \mathbf{Q}_{\tau+1}(\mathbf{x}_\tau \otimes \mathbf{x}_\tau)\mathbf{Q}_{\tau+1} \dots \mathbf{Q}_{t-1}, & \text{if } \tau < t-1, \\ \mathbf{x}_{t-1} \otimes \mathbf{x}_{t-1}, & \text{if } \tau = t-1. \end{cases} \quad (66)$$

Define the event where these are bounded in trace as

$$\mathcal{E}_2^{\text{sgd}} := \bigcap_{t=0}^{T-1} \bigcap_{\tau=0}^{t-1} \left\{ \text{Tr}[\mathbf{W}_{t,\tau}] \leq \frac{T^2 R^2}{p} (1 - \eta\mu)^{t-1-\tau} \right\}. \quad (67)$$

- Define the event where the norm of $\tilde{\boldsymbol{\theta}}^{\text{dp}}$ defined in (51) is bounded as

$$\mathcal{E}_1^{\text{dp}} := \bigcap_{t=0}^{T-1} \left\{ \|\tilde{\boldsymbol{\theta}}^{\text{dp}}\|_2^2 \leq C\eta^2 \sigma^2 \text{Tr}[\boldsymbol{\Sigma}_t^{\text{dp}}] \log\left(\frac{T}{p}\right) \right\}, \quad (68)$$

where $\boldsymbol{\Sigma}_t^{\text{dp}}$ is defined in (56).

- Define the event where the matrix $\mathbf{V}_{t,\tau}$ defined in (56) is bounded in trace:

$$\mathcal{E}_2^{\text{dp}} := \bigcap_{t=0}^{T-1} \bigcap_{\tau=0}^{t-1} \left\{ \text{Tr}[\mathbf{V}_{t,\tau} \mathbf{V}_{t,\tau}^\top] \leq \frac{T^2 d}{p} \left(\sum_{k=0}^{\tau} |\beta_k| (1 - \eta\mu)^{(\tau-k)/2} \right) \right\}. \quad (69)$$

We show that all these events hold with high probability.

Proposition D.5. *Consider the setting of Theorem D.4. We have,*

$$\mathbb{P}\left(\mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_1^{\text{sgd}} \cap \mathcal{E}_2^{\text{sgd}} \cap \mathcal{E}_1^{\text{dp}} \cap \mathcal{E}_2^{\text{dp}}\right) \geq 1 - 6p.$$

Proof. We will show that each of the events holds with probability at least $1 - p$ and a union bound gives the desired result.

Event \mathcal{E}_1 : Since $\mathbf{z}_t = \mathbf{H}^{-1/2} \mathbf{x}_t$ is element-wise independent and 1-sub-Gaussian, we have from the Hanson-Wright inequality (Lemma F.6) that

$$\mathbb{P}(\|\mathbf{x}_t\|_2^2 > C \text{Tr}[\mathbf{H}] \log(1/p)) = \mathbb{P}(\langle \mathbf{z}_t, \mathbf{H} \mathbf{z}_t \rangle > C \text{Tr}[\mathbf{H}] \log(1/p)) \leq p.$$

Taking a union bound over $t = 0, 1, \dots, T-1$ gives that $\mathbb{P}(\mathcal{E}_1) \geq 1 - p$.

Event \mathcal{E}_2 : Since ξ_t is sub-Gaussian with mean zero and variance proxy σ_{sgd}^2 , we have,

$$\mathbb{P}(|\xi_t| > s) \leq 2 \exp\left(-\frac{s^2}{2\sigma_{\text{sgd}}^2}\right).$$

Setting the right side equal to p/T and taking a union bound over $t = 0, 1, \dots, T-1$ gives $\mathbb{P}(\mathcal{E}_2) \geq 1 - p$.

Event $\mathcal{E}_1^{\text{sgd}}$: From the expression for $\tilde{\boldsymbol{\theta}}_t^{\text{sgd}}$ from (59), we can say that $\tilde{\boldsymbol{\theta}}_t^{\text{sgd}}$ conditioned on $\mathbf{x}_0, \dots, \mathbf{x}_{t-1}$ is mean zero and satisfies

$$\tilde{\boldsymbol{\theta}}_t^{\text{sgd}} = \eta \underbrace{[\mathbf{x}_{t-1} \quad \mathbf{Q}_{t-1} \mathbf{x}_{t-1} \quad \cdots \quad (\mathbf{Q}_{t-1} \cdots \mathbf{Q}_1 \mathbf{x}_0)]}_{=: \mathbf{M}_t} \begin{bmatrix} \xi_{t-1} \\ \vdots \\ \xi_0 \end{bmatrix}.$$

Using the assumption that each ξ_τ is independent and sub-Gaussian with variance proxy σ_{sgd}^2 , we get from the Hanson-Wright inequality (Lemma F.6) again that

$$\mathbb{P}\left(\left\|\tilde{\boldsymbol{\theta}}_t^{\text{sgd}}\right\|_2^2 > C\eta^2 \sigma_{\text{sgd}}^2 \text{Tr}[\mathbf{M}_t \mathbf{M}_t^\top] \log(1/p)\right) = \mathbb{P}(\langle \boldsymbol{\xi}_{:t}, \mathbf{M}_t \mathbf{M}_t^\top \boldsymbol{\xi}_{:t} \rangle > C\eta^2 \sigma_{\text{sgd}}^2 \text{Tr}[\mathbf{M}_t \mathbf{M}_t^\top] \log(1/p)) \leq p.$$

Next, we confirm that

$$\text{Tr}[\mathbf{M}_t \mathbf{M}_t^\top] = \|\mathbf{x}_{t-1}\|_2^2 + \|\mathbf{Q}_{t-1} \mathbf{x}_{t-1}\|_2^2 + \cdots + \|\mathbf{Q}_{t-1} \cdots \mathbf{Q}_1 \mathbf{x}_0\|_2^2 = \text{Tr}[\boldsymbol{\Sigma}_t^{\text{sgd}}].$$

Finally, a union bound over $t = 0, 1, \dots, T-1$ gives that $\mathbb{P}(\mathcal{E}_1^{\text{sgd}}) \geq 1 - p$.

Event $\mathcal{E}_2^{\text{sgd}}$: Markov's inequality gives

$$\mathbb{P}(\text{Tr}[\mathbf{W}_{t,\tau}] > s) \leq \frac{1}{s} \mathbb{E}[\mathbf{W}_{t,\tau}] \leq (1 - \eta\mu)^{t-1-\tau} \frac{R^2}{s}$$

where the calculations for the expected bound are deferred to Lemma D.9. Taking a union bound over all $T(T+1)/2 \leq T^2$ choices of (t, τ) gives $\mathbb{P}(\mathcal{E}_2^{\text{sgd}}) \geq 1 - p$.

Event $\mathcal{E}_1^{\text{dp}}$: From the expression for $\tilde{\boldsymbol{\theta}}_t^{\text{dp}}$ from (60), we deduce that

$$\tilde{\boldsymbol{\theta}}_t^{\text{dp}} | \mathbf{x}_0, \dots, \mathbf{x}_{t-1} \sim \mathcal{N}(\mathbf{0}, \eta^2 \sigma^2 \boldsymbol{\Sigma}_t^{\text{dp}}).$$

Invoking the Hanson-Wright inequality (Lemma F.6) and union bounding over $t = 0, \dots, T-1$ gives $\mathbb{P}(\mathcal{E}_1^{\text{dp}}) \geq 1 - p$.

Event $\mathcal{E}_2^{\text{dp}}$: Markov's inequality gives

$$\mathbb{P}(\text{Tr}[\mathbf{V}_{t,\tau} \mathbf{V}_{t,\tau}^\top] > s) \leq \frac{1}{s} \mathbb{E}[\mathbf{V}_{t,\tau} \mathbf{V}_{t,\tau}^\top] \leq \left(\sum_{k=0}^{\tau} |\beta_k| (1 - \eta\mu)^{(\tau-k)/2}\right) \frac{d}{s}$$

where we defer the technical calculations involved in bounding the expectation above to Lemma D.10. Taking a union bound over all $T(T+1)/2 \leq T^2$ choices of (t, τ) gives $\mathbb{P}(\mathcal{E}_2^{\text{dp}}) \geq 1 - p$. \square

D.2.2 High Probability Bounds on Component Recursions

Bound on the noise-less iterates: We start with $\hat{\theta}_t$ from (49).

Proposition D.6. *Under event \mathcal{E}_1 and if $\eta \leq (CR^2 \log(T/p))^{-1}$, we have that $\|\hat{\theta}_t\|_2 \leq \|\theta'_0\|_2$.*

Proof. Using the fact that $\mathbf{0} \preceq \mathbf{Q}_t \preceq \mathbf{I}$ under \mathcal{E}_1 (cf. Equation (65)), we get

$$\|\hat{\theta}_t\|_2 = \|\mathbf{Q}_{t-1} \cdots \mathbf{Q}_0 \theta'_0\|_2 \leq \|\mathbf{Q}_{t-1}\|_2 \cdots \|\mathbf{Q}_0\|_2 \|\theta'_0\|_2 \leq \|\theta'_0\|_2.$$

□

Bound on $\tilde{\theta}_t^{\text{sgd}}$: We turn to $\tilde{\theta}_t^{\text{sgd}}$ from (50).

Proposition D.7. *Under events $\mathcal{E}_1, \mathcal{E}_1^{\text{sgd}}, \mathcal{E}_2^{\text{sgd}}$, and $\eta \leq (CR^2 \log(T/p))^{-1}$, we have*

$$\|\tilde{\theta}_t^{\text{sgd}}\|_2^2 \leq C \left(\frac{\eta R^2}{\mu} \right) \log^3 \left(\frac{T}{p} \right).$$

Proof. Under $\mathcal{E}_1^{\text{sgd}}$, we have

$$\|\tilde{\theta}^{\text{sgd}}\|_2^2 \leq C \eta^2 \sigma_{\text{sgd}}^2 \text{Tr} \left[\boldsymbol{\Sigma}_t^{\text{sgd}} \right] \log \left(\frac{T}{p} \right). \quad (70)$$

We bound $\text{Tr}[\boldsymbol{\Sigma}_t] = \sum_{\tau=0}^{t-1} \text{Tr}[\mathbf{W}_{t,\tau}]$ for $\mathbf{W}_{t,\tau}$ defined in (66). We have two bounds for $\text{Tr}[\mathbf{W}_{t,\tau}]$:

(a) Using $\mathbf{0} \preceq \mathbf{Q}_t \preceq \mathbf{I}$ under \mathcal{E}_1 (cf. Equation (65)), we bound

$$\text{Tr}[\mathbf{W}_{t,\tau}] = \|\mathbf{Q}_{t-1} \cdots \mathbf{Q}_{\tau+1} \mathbf{x}_\tau\|_2^2 \leq \|\mathbf{Q}_{t-1}\|_2^2 \cdots \|\mathbf{Q}_{\tau+1}\|_2^2 \|\mathbf{x}_\tau\|_2^2 \leq CR^2 \log(T/p).$$

(b) Under event $\mathcal{E}_2^{\text{sgd}}$, we have the bound

$$\text{Tr}[\mathbf{W}_{t,\tau}] \leq \frac{T^2 R^2}{p} (1 - \eta\mu)^{t-1-\tau}.$$

Using the first bound for the last $\tau \leq t-1$ iterations and the second bound for the rest, we get

$$\begin{aligned} \text{Tr} \left[\boldsymbol{\Sigma}_t^{\text{sgd}} \right] &\leq \sum_{k=0}^{t-\tau-1} \frac{T^2 R^2}{p} (1 - \eta\mu)^{t-1-\tau} \mathbb{1}(\tau < t-1) + \tau (CR^2 \log(T/p)) \\ &\leq \frac{T^2 R^2}{p} (1 - \eta\mu)^\tau \sum_{k=0}^{t-\tau-1} (1 - \eta\mu)^k \mathbb{1}(\tau < t-1) + \tau (CR^2 \log(T/p)) \\ &\leq \frac{T^2 R^2}{p} \frac{\exp(-\eta\mu\tau)}{\eta\mu} \mathbb{1}(\tau < t-1) + \tau (CR^2 \log(T/p)). \end{aligned}$$

Choosing $\tau = \min \left\{ t-1, \frac{1}{\eta\mu} \log \left(\frac{T^2}{Cp \log(T/p)} \right) \right\}$ as per Lemma F.20 gives

$$\text{Tr} \left[\boldsymbol{\Sigma}_t^{\text{sgd}} \right] \leq \frac{CR^2 \log(T/p)}{\eta\mu} \left(1 + \log \left(\frac{T^2}{p \log(T/p)} \right) \right) \leq \frac{C'R^2}{\eta\mu} \log^2(T/p)$$

for some absolute constants C, C' . Plugging this back into (70) completes the proof. □

Bound on $\tilde{\theta}_t^{\text{dp}}$: We turn to $\tilde{\theta}_t^{\text{dp}}$ from (51).

Proposition D.8. *Consider the setting of Theorem D.4. Under events $\mathcal{E}_1, \mathcal{E}_1^{\text{dp}}, \mathcal{E}_2^{\text{dp}}$, and $\eta \leq (CR^2 \log(T/p))^{-1}$, we have*

$$\|\tilde{\theta}_t^{\text{sgd}}\|_2^2 \leq C \left(\frac{\eta R^2}{\mu} \right) \log^3 \left(\frac{T}{p} \right).$$

Proof. Based on the bound on $\|\tilde{\boldsymbol{\theta}}_t^{\text{dp}}\|_2$ from $\mathcal{E}_1^{\text{dp}}$, we bound $\text{Tr} [\boldsymbol{\Sigma}_t^{\text{dp}}] = \sum_{\tau=0}^{t-1} \text{Tr} [\mathbf{V}_{t,\tau} \mathbf{V}_{t,\tau}^\top]$. We bound each trace on the right side in two ways:

- (a) We have $\text{Tr} [\mathbf{V}_{t,\tau} \mathbf{V}_{t,\tau}^\top] \leq \|\boldsymbol{\beta}\|_1^2 d$ from Lemma D.10.
- (b) Under $\mathcal{E}_2^{\text{dp}}$ and the assumption (*) of Half-Expo Decay of β with parameter $\nu \leq \eta\mu$, we also have

$$\begin{aligned} \text{Tr} [\mathbf{V}_{t,\tau} \mathbf{V}_{t,\tau}^\top] &\leq \frac{T^2 d}{p} \left(\sum_{k=0}^{\tau} |\beta_k| (1 - \eta\mu)^{(\tau-k)/2} \right)^2 \\ &\leq \frac{T^2 d}{p} \left(\sum_{k=0}^{\tau} |\beta_k| (1 - \nu)^{(\tau-k)/2} \right)^2 \\ &\stackrel{(*)}{\leq} \frac{CT^2 d}{p} (1 - \nu)^\tau. \end{aligned}$$

Using the first bound for the first τ iterations and the second bound for the rest, we get

$$\begin{aligned} \text{Tr} [\boldsymbol{\Sigma}_t^{\text{dp}}] &\leq \tau \left(\|\boldsymbol{\beta}\|_1^2 d \right) + \sum_{k=\tau}^{t-1} \frac{CT^2 d}{p} (1 - \nu)^k \mathbb{1}(\tau > t - 1) \\ &\leq \tau \left(\|\boldsymbol{\beta}\|_1^2 d \right) + \frac{CT^2 d}{p} (1 - \nu)^\tau \sum_{k=0}^{\infty} (1 - \nu)^k \mathbb{1}(\tau > t - 1) \\ &\leq \tau \left(\|\boldsymbol{\beta}\|_1^2 d \right) + \frac{CT^2 d \exp(-\nu\tau)}{p\nu} \mathbb{1}(\tau > t - 1). \end{aligned}$$

Choosing $\tau \leq \left\{ t - 1, \frac{1}{\nu} \log(CT^2/p\|\boldsymbol{\beta}\|_1^2) \right\}$ as per Lemma F.20, we get,

$$\text{Tr} [\boldsymbol{\Sigma}_t^{\text{dp}}] \leq \frac{\|\boldsymbol{\beta}\|_1^2 d}{\nu} \left(1 + \log \left(\frac{CT^2}{p\|\boldsymbol{\beta}\|_1^2} \right) \right) \leq C' \frac{\|\boldsymbol{\beta}\|_1^2 d}{\nu} \log \left(\frac{T}{p} \right),$$

where we used $\|\boldsymbol{\beta}\|_1 \geq |\beta_0| = 1$ and C, C' are some universal constants. Combining this with the bound on $\|\tilde{\boldsymbol{\theta}}_t^{\text{dp}}\|_2$ asserted by $\mathcal{E}_1^{\text{dp}}$ completes the proof. \square

D.2.3 Completing the Proof of the High Probability Bounds

We are now ready to prove Theorem D.4.

Proof of Theorem D.4. Under events $\mathcal{E}_1, \mathcal{E}_1^{\text{sgd}}, \mathcal{E}_2^{\text{sgd}}, \mathcal{E}_1^{\text{dp}}, \mathcal{E}_2^{\text{dp}}$, we have bounds on the norms of $\hat{\boldsymbol{\theta}}_t, \tilde{\boldsymbol{\theta}}_t^{\text{sgd}}, \tilde{\boldsymbol{\theta}}_t^{\text{dp}}$ respectively from Propositions D.6 to D.8. We combine them together with the triangle inequality and Equation (57) of Property D.2 to the claimed bound on $\|\boldsymbol{\theta}'_t\|_2$.

Next, for the gradients, we use the triangle and Cauchy-Schwarz inequalities on the definition $\mathbf{g}_t = \mathbf{x}_t \langle \mathbf{x}_t, \boldsymbol{\theta}'_t \rangle - \mathbf{x}_t \xi_t$ to get

$$\|\mathbf{g}_t\|_2^2 \leq 2 \|\mathbf{x}_t\|_2^4 \|\boldsymbol{\theta}'_t\|_2^2 + 2 \|\mathbf{x}_t\|_2^2 |\xi_t|_2^2.$$

Plugging in the bounds on $\|\mathbf{x}_t\|_2$ and $|\xi_t|_2$ from \mathcal{E}_1 and \mathcal{E}_2 respectively gives the claimed bound on $\|\mathbf{g}_t\|_2^2$.

Finally, all the events above hold with probability at least $1 - 6p$ from Proposition D.5. Substituting $p/6$ for p and adjusting the constants completes the proof. \square

D.2.4 Helper Lemmas

Lemma D.9. Consider the setting of Theorem D.4 and consider the PSD matrices $\mathbf{W}_{t,\tau}$, defined for $\tau \leq t - 1$ as

$$\mathbf{W}_{t,\tau} = \begin{cases} \mathbf{Q}_{t-1} \cdots \mathbf{Q}_{\tau+1} (\mathbf{x}_\tau \otimes \mathbf{x}_\tau) \mathbf{Q}_{\tau+1} \cdots \mathbf{Q}_{t-1}, & \text{if } \tau < t - 1, \\ \mathbf{x}_{t-1} \otimes \mathbf{x}_{t-1}, & \text{if } \tau = t - 1. \end{cases}$$

We have that $\mathbb{E}[\text{Tr} [\mathbf{W}_{t,\tau}]] \leq R^2 (1 - \eta\mu)^{t-1-\tau}$.

Proof. For $\tau = t - 1$, we have $\mathbb{E}[\mathbf{W}_{t,t-1}] = \text{Tr}[\mathbf{H}] \leq R^2$. For $\tau < t - 1$, we have by independence of each \mathbf{x}_t that

$$\begin{aligned} \text{Tr}[\mathbb{E}[\mathbf{W}_{t,\tau}]] &= \text{Tr}[\mathbb{E}[\mathbf{Q}_{t-1} \cdots \mathbf{Q}_{\tau+1} \mathbf{H} \mathbf{Q}_{\tau+1} \cdots \mathbf{Q}_{t-1}]] = \text{Tr}[\mathbb{E}[\mathbf{Q}_{t-1} \cdots \mathbf{Q}_{\tau} (\mathcal{P}\mathbf{H}) \mathbf{Q}_{\tau} \cdots \mathbf{Q}_{t-1}]] = \cdots \\ &= \text{Tr}[\mathcal{P}^{t-1-\tau} \mathbf{H}]. \end{aligned}$$

Recursively bounding $\text{Tr}[\mathcal{P}^\tau \mathbf{H}] = \text{Tr}[\mathcal{P}(\mathcal{P}^{\tau-1} \mathbf{H})] \leq (1 - \eta\mu) \text{Tr}[\mathcal{P}^{\tau-1} \mathbf{H}]$ from Lemma F.18 completes the proof. \square

Lemma D.10. Consider $\mathbf{V}_{t,\tau}$ as defined in (56). We have that

$$\mathbb{E}[\text{Tr}[\mathbf{V}_{t,\tau} \mathbf{V}_{t,\tau}^\top]] \leq d \left(\sum_{k=0}^{\tau} |\beta_k| (1 - \eta\mu)^{(\tau-k)/2} \right).$$

Further, if the event $\mathcal{E} = \cap_{\tau=1}^t \{\mathbf{Q}_t \succeq \mathbf{0}\}$ holds, then we also have

$$\text{Tr}[\mathbf{V}_{t,\tau} \mathbf{V}_{t,\tau}^\top] \leq d \left(\sum_{k=0}^{\tau} |\beta_k| \right)^2.$$

Proof. Since t is fixed throughout, we simply write $\mathbf{V}_{t,\tau}$ as \mathbf{V}_τ . We define a sequence of matrices $\mathbf{A}_0, \dots, \mathbf{A}_\tau$ as $\mathbf{A}_0 = \beta_0 \mathbf{I}$ and

$$\mathbf{A}_{k+1} = \beta_{k+1} \mathbf{I} + \mathbf{Q}_{t-\tau+k} \mathbf{A}_k$$

for $k = 0, \dots, \tau - 1$. We first prove the expected bound followed by the absolute bound.

Expected bound: Then, we successively deduce the following.

- (a) We have $\mathbf{A}_k = \beta_k \mathbf{I} + \beta_{k-1} \mathbf{Q}_{t-\tau+k-1} + \cdots + \beta_0 \mathbf{Q}_{t-\tau+k-1} \cdots \mathbf{Q}_{t-\tau}$ by simply unrolling the recursions.
- (b) We immediately recognize that $\mathbf{V}_\tau = \mathbf{A}_\tau$.
- (c) By independence of each \mathbf{Q}_t , taking an expectation of the expression in (a) gives

$$\mathbb{E}[\mathbf{A}_k] = \sum_{l=0}^k \beta_l (\mathbf{I} - \eta\mathbf{H})^{k-l}.$$

- (d) We establish a recursion

$$\mathbb{E} \text{Tr}[\mathbf{A}_{k+1} \mathbf{A}_{k+1}^\top] \leq d\beta_{k+1}^2 + 2d|\beta_{k+1}| \sum_{l=0}^k |\beta_l| (1 - \eta\mu)^{k-l+1} + (1 - \eta\mu) \mathbb{E} \text{Tr}[\mathbf{A}_k \mathbf{A}_k^\top].$$

Indeed, by expanding out the square of the recursion and using the independence of the \mathbf{x}_t 's, we get

$$\begin{aligned} \mathbb{E} \text{Tr}[\mathbf{A}_{k+1} \mathbf{A}_{k+1}^\top] &= \beta_{k+1}^2 \text{Tr}[\mathbf{I}] + 2\beta_{k+1} \text{Tr}[(\mathbf{I} - \eta\mathbf{H}) \mathbb{E}[\mathbf{A}_k]] + \text{Tr}[\mathcal{P}(\mathbb{E}[\mathbf{A}_k \mathbf{A}_k^\top])] \\ &\leq d\beta_{k+1}^2 + 2|\beta_{k+1}| \sum_{l=0}^k |\beta_l| \text{Tr}[(\mathbf{I} - \eta\mathbf{H})^{k-l+1}] + (1 - \eta\mu) \mathbb{E} \text{Tr}[\mathbf{A}_k \mathbf{A}_k^\top], \end{aligned}$$

where we plugged in the expression for $\mathbb{E}[\mathbf{A}_k]$ from item (c) and used Lemma F.18 to bound the last term. Using $\mathbf{0} \preceq \mathbf{I} - \eta\mathbf{H} \preceq (1 - \eta\mu)\mathbf{I}$ gives the claimed expression.

- (e) Using induction and the recursion from part (d), we prove that

$$\mathbb{E} \text{Tr}[\mathbf{A}_k \mathbf{A}_k^\top] \leq d \left(\sum_{l=0}^k |\beta_l| (1 - \eta\mu)^{(k-l)/2} \right)^2.$$

Together with part (b), this gives the desired result.

Indeed, the base case holds because $\mathbb{E}\text{Tr}[\mathbf{A}_0\mathbf{A}_0^\top] = \beta_0^2 d$. Supposing the induction hypothesis holds for some $k < \tau - 1$, we use the recursion of item (d) to get

$$\begin{aligned} \frac{1}{d} \mathbb{E}\text{Tr}[\mathbf{A}_{k+1}\mathbf{A}_{k+1}^\top] &\leq \beta_{k+1}^2 + 2|\beta_{k+1}| \sum_{l=0}^k |\beta_l| (1 - \eta\mu)^{k-l+1} + \left(\sum_{l=0}^k |\beta_l| (1 - \eta\mu)^{\frac{k-l+1}{2}} \right)^2 \\ &\leq \beta_{k+1}^2 + 2|\beta_{k+1}| \sum_{l=0}^k |\beta_l| (1 - \eta\mu)^{\frac{k-l+1}{2}} + \left(\sum_{l=0}^k |\beta_l| (1 - \eta\mu)^{\frac{k-l+1}{2}} \right)^2 \\ &= \left(\sum_{l=0}^{k+1} |\beta_l| (1 - \eta\mu)^{\frac{k-l+1}{2}} \right)^2, \end{aligned}$$

where the second inequality used $1 - \eta\mu \leq 1$.

Absolute bound: Next, we prove the absolute bound, assuming that \mathcal{E} holds. Again, we successively deduce:

- (a) We starting with $\mathbf{A}_k = \beta_k \mathbf{I} + \beta_{k-1} \mathbf{Q}_{t-\tau+k-1} + \cdots + \beta_0 \mathbf{Q}_{t-\tau+k-1} \cdots \mathbf{Q}_{t-\tau}$.
- (b) Then, we get

$$|\text{Tr}[\mathbf{A}_k]| \leq |\beta_k|d + |\beta_{k-1}| |\text{Tr}[\mathbf{Q}_{t-\tau+k-1}]| + \cdots + |\beta_0| |\text{Tr}[\mathbf{Q}_{t-\tau+k-1} \cdots \mathbf{Q}_{t-\tau}]| \leq d \sum_{l=0}^k |\beta_l|,$$

where we bound each of the traces by d using Lemma F.19 (since we have $\mathbf{Q}_k \preceq \mathbf{I}$ under \mathcal{E}).

- (c) By a similar logic, we get

$$\begin{aligned} &\left| \text{Tr}[\mathbf{Q}_{t-\tau+k} \mathbf{A}_k + \mathbf{A}_k^\top \mathbf{Q}_{t-\tau+k}] \right| \\ &\leq 2|\beta_k| |\text{Tr}[\mathbf{Q}_{t-\tau+k}]| + 2|\beta_1| |\text{Tr}[\mathbf{Q}_{t-\tau+k} \mathbf{Q}_{t-\tau+k-1}]| + \cdots + 2|\beta_0| |\text{Tr}[\mathbf{Q}_{t-\tau+k} \cdots \mathbf{Q}_{t-\tau}]| \\ &\leq 2d \sum_{l=0}^k |\beta_l|. \end{aligned}$$

- (d) We prove by induction that $\text{Tr}[\mathbf{A}_k \mathbf{A}_k^\top] \leq d \left(\sum_{l=0}^k |\beta_l| \right)^2$.

The base case holds since $\text{Tr}[\mathbf{A}_0 \mathbf{A}_0^\top] = d\beta_0^2$. Supposing the induction hypothesis holds for some integer $1 \leq k < t - 1$, we use the recursion of \mathbf{A}_{k+1} to calculate

$$\begin{aligned} \text{Tr}[\mathbf{A}_{k+1} \mathbf{A}_{k+1}^\top] &= d\beta_{k+1}^2 + \beta_{k+1} \text{Tr}[\mathbf{Q}_{t-\tau+k} \mathbf{A}_k + \mathbf{A}_k^\top \mathbf{Q}_{t-\tau+k}] + \text{Tr}[\mathbf{Q}_{t-\tau+k} \mathbf{A}_k \mathbf{A}_k^\top \mathbf{Q}_{t-\tau+k}] \\ &\leq d\beta_{k+1}^2 + 2d|\beta_{k+1}| \sum_{l=0}^k |\beta_l| + \text{Tr}[\mathbf{A}_k \mathbf{A}_k^\top] \leq d \left(\sum_{l=0}^{k+1} |\beta_l| \right)^2. \end{aligned}$$

Finally, item (d) together with $\mathbf{A}_\tau = \mathbf{V}_{t,\tau}$ completes the proof. \square

D.3 Expected Bounds on Noisy-FTRL

Our goal of this section is to prove the following finite-time convergence guarantee of Noisy-FTRL in terms of the asymptotic suboptimality.

Theorem D.11. *Consider problem (15) and suppose Assumption C.2 holds. For a given a starting iterate $\boldsymbol{\theta}_0 \in \mathbb{R}^d$, weights $\boldsymbol{\beta} \in \ell^2$, learning rate $\eta < 1/R^2$, consider the sequence $(\boldsymbol{\theta}_t)_{t=0}^\infty$ produced by the iteration (16) where $\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ with $\sigma^2 = G^2 \gamma_\infty(\boldsymbol{\beta})^2 / (2\rho)$. Then, for any $t \geq 0$, we have,*

$$\mathbb{E}[F(\boldsymbol{\theta}_t) - F(\boldsymbol{\theta}_\star)] \leq \left(\sqrt{\frac{L}{\mu} \exp(-\eta\mu t)} (F(\boldsymbol{\theta}_0) - F(\boldsymbol{\theta}_\star)) + \sqrt{F_\infty(\boldsymbol{\beta})} \right)^2.$$

We start with some preliminary lemmas. The first lemma is about the covariance of the noise process and is a generalization of [29, Lemma 3] to linearly correlated additive noise.

Lemma D.12. Consider the sequence $(\tilde{\boldsymbol{\theta}}_t)_{t=0}^\infty$ generated by Noisy-FTRL starting from $\tilde{\boldsymbol{\theta}}_t = \boldsymbol{\theta}_*$ with noise correlations $\beta \in \ell^2$ and learning rate $\eta \leq 1/R^2$. Under Assumption C.2, we have that its covariance

$$\mathbf{S}_t := \mathbb{E} \left[\left(\tilde{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_* \right) \otimes \left(\tilde{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_* \right) \right]$$

satisfies: (a) $\mathbf{S}_t \preceq \mathbf{S}_{t+1}$ for all $t \geq 0$, and (b) the sequence $(\mathbf{S}_t)_{t=0}^\infty$ converges element-wise as $t \rightarrow \infty$.

Proof. Recall the notation $\mathbf{Q}_t = \mathbf{I} - \eta \mathbf{x}_t \otimes \mathbf{x}_t$ and $\mathcal{PM} = \mathbb{E}[\mathbf{Q}_t \mathbf{M} \mathbf{Q}_t]$. We use the shorthand $\tilde{\boldsymbol{\theta}}'_t := \tilde{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_*$. We first prove that the covariance is increasing in a PSD sense and argue that its limit exists.

Part 1: Non-decreasing noise: By unrolling the update equation and using $\tilde{\boldsymbol{\theta}}'_t = \mathbf{0}$, we get

$$\begin{aligned} \tilde{\boldsymbol{\theta}}'_t &= \eta (\mathbf{x}_{t-1} \xi_{t-1} + \mathbf{Q}_{t-1} \mathbf{x}_{t-2} \xi_{t-2} + \cdots + \mathbf{Q}_{t-1} \cdots \mathbf{Q}_1 \mathbf{x}_0 \xi_0) \\ &\quad - \eta \left(\beta_0 \mathbf{w}_{t-1} + (\beta_1 \mathbf{I} + \beta_0 \mathbf{Q}_{t-1}) \mathbf{w}_{t-2} + \cdots + (\beta_{t-1} \mathbf{I} + \beta_{t-2} \mathbf{Q}_{t-1} + \cdots + \beta_0 \mathbf{Q}_{t-1} \cdots \mathbf{Q}_1) \mathbf{w}_0 \right). \end{aligned} \quad (71)$$

Next, we calculate $\mathbb{E} \left[\tilde{\boldsymbol{\theta}}'_t \otimes \tilde{\boldsymbol{\theta}}'_t \right]$. By independence, all the cross terms cancel out, so it suffices to write out the second moment of each of the terms above. For the SGD noise terms that contain $\mathbf{x}_\tau \xi_\tau$, we get for $\tau = 0, \dots, t-1$ that

$$\mathbb{E} \left[(\mathbf{Q}_{t-1} \cdots \mathbf{Q}_{t-\tau+1} \mathbf{x}_{t-\tau} \xi_{t-\tau}) \otimes (\mathbf{Q}_{t-1} \cdots \mathbf{Q}_{t-\tau+1} \mathbf{x}_{t-\tau} \xi_{t-\tau}) \right] = \mathcal{P}^\tau (\mathbb{E}[\xi^2 \mathbf{x} \otimes \mathbf{x}]) =: \mathcal{T}_\tau. \quad (72)$$

Since it is a second moment term, we have $\mathcal{T}_\tau \succeq \mathbf{0}$. For the DP noise terms, denote $\mathbf{x}^{\otimes 2} = \mathbf{x} \otimes \mathbf{x} = \mathbf{x} \mathbf{x}^\top$. Then, we have for $\tau = 0$ to $t-1$ that

$$\begin{aligned} &\frac{1}{\sigma^2} \mathbb{E} \left((\beta_\tau \mathbf{I} + \beta_{\tau-1} \mathbf{Q}_{t-1} + \beta_{\tau-2} \mathbf{Q}_{t-1} \mathbf{Q}_{t-2} + \cdots + \beta_0 \mathbf{Q}_{t-1} \cdots \mathbf{Q}_{t-\tau}) \mathbf{w}_{t-\tau-1} \right)^{\otimes 2} \\ &= \mathbb{E} (\beta_\tau \mathbf{I} + \beta_{\tau-1} \mathbf{Q}_{t-1} + \beta_{\tau-2} \mathbf{Q}_{t-1} \mathbf{Q}_{t-2} + \cdots + \beta_0 \mathbf{Q}_{t-1} \cdots \mathbf{Q}_{t-\tau})^{\otimes 2} \\ &= \beta_\tau^2 \mathbf{I} + 2\beta_\tau \sum_{k=0}^{\tau-1} \beta_k (\mathbf{I} - \eta \mathbf{H})^{\tau-k} + \sum_{k=0}^{\tau-1} \sum_{l=0}^{\tau-1-k} \beta_k \beta_l \mathbb{E} [\mathbf{Q}_{t-1} \cdots \mathbf{Q}_{t-\tau+k} \mathbf{Q}_{t-\tau+l} \cdots \mathbf{Q}_{t-1}] \\ &= \beta_\tau^2 \mathbf{I} + 2\beta_\tau \sum_{k=0}^{\tau-1} \beta_k (\mathbf{I} - \eta \mathbf{H})^{\tau-k} + 2 \sum_{k=0}^{\tau-1} \sum_{l=0}^k \beta_k \beta_l \mathbb{E} [\mathbf{Q}_{t-1} \cdots \mathbf{Q}_{t-\tau+l} (\mathbf{I} - \eta \mathbf{H})^{k-l} \mathbf{Q}_{t-\tau+l} \cdots \mathbf{Q}_{t-1}] \\ &= \beta_\tau^2 \mathbf{I} + 2\beta_\tau \sum_{k=0}^{\tau-1} \beta_k (\mathbf{I} - \eta \mathbf{H})^{\tau-k} + 2 \sum_{k=0}^{\tau-1} \sum_{l=0}^k \beta_k \beta_l \mathcal{P}^{\tau-k} ((\mathbf{I} - \eta \mathbf{H})^{k-l}) =: \mathcal{T}'_\tau. \end{aligned} \quad (73)$$

By virtue of this being a second moment, we have that $\mathcal{T}'_\tau \succeq \mathbf{0}$. Plugging in (72) and (73) into the second moment of (71), we get,

$$\begin{aligned} \mathbb{E} \left[\tilde{\boldsymbol{\theta}}'_{t+1} \otimes \tilde{\boldsymbol{\theta}}'_{t+1} \right] &= \eta^2 \sum_{\tau=0}^t (\mathcal{T}_\tau + \sigma^2 \mathcal{T}'_\tau) \\ &= \mathbb{E} \left[\tilde{\boldsymbol{\theta}}'_t \otimes \tilde{\boldsymbol{\theta}}'_t \right] + \eta^2 (\mathcal{T}_t + \sigma^2 \mathcal{T}'_t) \succeq \mathbb{E} \left[\tilde{\boldsymbol{\theta}}'_t \otimes \tilde{\boldsymbol{\theta}}'_t \right]. \end{aligned}$$

This shows that the noise is non-decreasing in a PSD sense.

Part 2: Convergence of the covariance: Next, we show that the noise sequence converges. From the update equation $\tilde{\boldsymbol{\theta}}'_{t+1} = \mathbf{Q}_t \tilde{\boldsymbol{\theta}}'_t + \eta \mathbf{x}_t \xi_t - \eta \sum_{\tau=0}^t \beta_\tau \mathbf{w}_{t-\tau}$, we get

$$\begin{aligned} \mathbf{S}_{t+1} &= \mathcal{P} \mathbf{S}_t + \eta^2 \mathbb{E}[\xi^2 \mathbf{x} \otimes \mathbf{x}] + \eta^2 \sigma^2 \sum_{\tau=0}^t \beta_\tau^2 \mathbf{I} \\ &\quad - \eta (\mathbf{I} - \eta \mathbf{H}) \sum_{\tau=0}^t \beta_\tau \mathbb{E} \left[\tilde{\boldsymbol{\theta}}'_t \otimes \mathbf{w}_{t-\tau} \right] - \eta \sum_{\tau=0}^t \beta_\tau \mathbb{E} \left[\mathbf{w}_{t-\tau} \otimes \tilde{\boldsymbol{\theta}}'_t \right] (\mathbf{I} - \eta \mathbf{H}). \end{aligned}$$

For $\tau = 0$, the term $\mathbb{E}[\tilde{\boldsymbol{\theta}}'_t \otimes \mathbf{w}_{t-\tau}]$ and its transpose are both $\mathbf{0}$. For $\tau > 0$, we have from (71) that

$$\begin{aligned} -\mathbb{E} \left[\tilde{\boldsymbol{\theta}}'_t \otimes \mathbf{w}_{t-\tau} \right] &= \eta \mathbb{E} \left[\beta_{\tau-1} \mathbf{I} + \beta_{\tau-2} \mathbf{Q}_{t-1} + \cdots + \beta_0 \mathbf{Q}_{t-1} \cdots \mathbf{Q}_{t-\tau+1} \right] \mathbb{E} [\mathbf{w}_{t-\tau} \otimes \mathbf{w}_{t-\tau}] \\ &= \eta \sigma^2 \left(\beta_{\tau-1} \mathbf{I} + \beta_{\tau-2} (\mathbf{I} - \eta \mathbf{H}) + \cdots + \beta_0 (\mathbf{I} - \eta \mathbf{H})^{\tau-1} \right). \end{aligned}$$

Plugging this back in gives

$$\begin{aligned} \mathbf{S}_{t+1} &= \mathcal{P} \mathbf{S}_t + \eta^2 \mathbb{E} [\xi^2 \mathbf{x} \otimes \mathbf{x}] + \eta^2 \sigma^2 \sum_{\tau=0}^t \beta_\tau^2 \mathbf{I} + 2\eta^2 \sigma^2 \sum_{\tau=1}^t \sum_{k=0}^{\tau-1} \beta_\tau \beta_k (\mathbf{I} - \eta \mathbf{H})^{\tau-k} \\ &= \mathcal{P} \mathbf{S}_t + \eta^2 \mathbb{E} [\xi^2 \mathbf{x} \otimes \mathbf{x}] + \eta^2 \sigma^2 \sum_{\tau=0}^t \sum_{k=0}^{\tau-1} \beta_\tau \beta_k (\mathbf{I} - \eta \mathbf{H})^{|\tau-k|}. \end{aligned} \quad (74)$$

Next, we take a trace of (74). For the first term, we get

$$\begin{aligned} \text{Tr} [\mathcal{P} \mathbf{S}_t] &= \text{Tr} [\mathbf{S}_t] - 2\eta \text{Tr} [\mathbf{H} \mathbf{S}_t] + \eta^2 \text{Tr} \left[\mathbf{S}_t \mathbb{E} [\|\mathbf{x}_t\|_2^2 \mathbf{x}_t \otimes \mathbf{x}_t] \right] \\ &\leq \text{Tr} [\mathbf{S}_t] - \eta \text{Tr} [\mathbf{H} \mathbf{S}_t] (2 - \eta R^2) \\ &\leq (1 - \eta \mu) \text{Tr} [\mathbf{S}_t], \end{aligned}$$

where we use (a) $\mathbb{E} [\|\mathbf{x}_t\|_2^2 \mathbf{x}_t \otimes \mathbf{x}_t] \preceq R^2 \mathbf{H}$, (b) $\eta \leq 1/R^2$, and (c) $\mathbf{H} \succeq \mu \mathbf{I}$. By assumption, we also get that $\text{Tr} [\mathbb{E} [\xi^2 \mathbf{x} \otimes \mathbf{x}]] \leq \sigma_{\text{sgd}}^2 \text{Tr} [\mathbf{H}] \leq \sigma_{\text{sgd}}^2 R^2$. Finally, we have using Lemma F.17 that

$$\sum_{\tau=0}^t \sum_{k=0}^{\tau-1} \beta_\tau \beta_k \sum_{j=1}^d (1 - \eta \lambda_j)^{|\tau-k|} \leq \|\boldsymbol{\beta}\|_2^2 \sum_{j=1}^d \left(\frac{2 - \eta \lambda_j}{\eta \lambda_j} \right) \leq \frac{2 \|\boldsymbol{\beta}\|_2^2 \text{Tr} [\mathbf{H}^{-1}]}{\eta}.$$

Thus, we get

$$\text{Tr} [\mathbf{S}_{t+1}] \leq (1 - \eta \mu) \text{Tr} [\mathbf{S}_t] + 2\eta \sigma^2 \|\boldsymbol{\beta}\|_2^2 \text{Tr} [\mathbf{H}^{-1}] + \eta^2 R^2 \sigma_{\text{sgd}}^2.$$

By unrolling this out, we get a uniform bound for all t :

$$\text{Tr} [\mathbf{S}_t] \leq \frac{1}{\mu} \left(2\sigma^2 \|\boldsymbol{\beta}\|_2^2 \text{Tr} [\mathbf{H}^{-1}] + \eta R^2 \sigma_{\text{sgd}}^2 \right) < \infty$$

since $\boldsymbol{\beta} \in \ell^2$. For any fixed vector \mathbf{v} , $\langle \mathbf{v}, \mathbf{S}_t \mathbf{v} \rangle$ thus has a limit from the monotone convergence theorem. From this, it follows that every diagonal entry of \mathbf{S}_t converges (take \mathbf{v} as a standard basis vector) and then every off-diagonal entry of \mathbf{S}_t also converges (take \mathbf{v} as the sum of two standard basis vectors). This shows that \mathbf{S}_t converges element-wise. \square

We are now ready to prove Theorem D.11.

Proof of Theorem D.11. Define $F_\infty^*(\boldsymbol{\beta})$ as the asymptotic suboptimality of a process that starts from $\boldsymbol{\theta}_0 = \boldsymbol{\theta}_*$. We will prove the desired result with $F_\infty^*(\boldsymbol{\beta})$ in the place of $F_\infty(\boldsymbol{\beta})$. Finally, we will show that $F_\infty(\boldsymbol{\beta})$ is independent of its starting iterate so $F_\infty(\boldsymbol{\beta}) = F_\infty^*(\boldsymbol{\beta})$.

We first separate out the effects of the noise and the initial iterate using Property D.2. We invoke Lemma D.12 for the former and directly bound the latter. Lastly, we combine them both with a triangle inequality. Recall that use the shorthand $\boldsymbol{\theta}'_t := \boldsymbol{\theta}_t - \boldsymbol{\theta}_*$ and $\mathbf{Q}_t := \mathbf{I} - \eta \mathbf{x}_t \otimes \mathbf{x}_t$.

Effect of the initialization: We first calculate

$$\mathbb{E} [\mathbf{Q}_t^2] = \mathbf{I} - 2\eta \mathbf{H} + \eta^2 \mathbb{E} \left[\|\mathbf{x}_t\|_2^2 \mathbf{x}_t \otimes \mathbf{x}_t \right] \preceq \mathbf{I} - 2\eta \mathbf{H} + \eta^2 R^2 \mathbf{H} \preceq \mathbf{I} - \eta \mathbf{H} \preceq (1 - \eta \mu) \mathbf{I},$$

where the first inequality follows from (47), the second since $\eta \leq 1/R^2$, and the third since $\mathbf{H} \succeq \mu \mathbf{I}$. Letting \mathcal{F}_t denote the sigma algebra generated by $\mathbf{x}_0, \dots, \mathbf{x}_{t-1}$, we get

$$\mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_{t+1} \right\|_2^2 \middle| \mathcal{F}_t \right] = \left\langle \hat{\boldsymbol{\theta}}_t, \mathbb{E} [\mathbf{Q}_t^2] \hat{\boldsymbol{\theta}}_t \right\rangle \leq (1 - \eta \mu) \left\| \hat{\boldsymbol{\theta}}_t \right\|_2^2 \leq \exp(-\eta \mu) \left\| \hat{\boldsymbol{\theta}}_t \right\|_2^2.$$

Taking an unconditional expectation and unrolling this and using $\mu \mathbf{I} \preceq \mathbf{H} \preceq L \mathbf{I}$ (Assumption **(A1)**) gives

$$\mathbb{E} \left\| \hat{\boldsymbol{\theta}}_t \right\|_{\mathbf{H}}^2 \leq L \mathbb{E} \left\| \hat{\boldsymbol{\theta}}_t \right\|_2^2 \leq L \exp(-\eta \mu t) \|\boldsymbol{\theta}'_0\|_2^2 \leq \frac{L}{\mu} \exp(-\eta \mu t) \|\boldsymbol{\theta}'_0\|_{\mathbf{H}}^2. \quad (75)$$

Effect of the noise: Define $\tilde{\boldsymbol{\theta}}'_t := \tilde{\boldsymbol{\theta}}_t^{\text{sgd}} + \tilde{\boldsymbol{\theta}}_t^{\text{dp}}$. We get from Lemma D.12 that there exists a PSD matrix \mathbf{S}_∞ such that

$$\mathbf{0} = \mathbb{E} \left[\tilde{\boldsymbol{\theta}}'_0 \otimes \tilde{\boldsymbol{\theta}}'_0 \right] \preceq \mathbb{E} \left[\tilde{\boldsymbol{\theta}}'_1 \otimes \tilde{\boldsymbol{\theta}}'_1 \right] \preceq \cdots \preceq \lim_{t \rightarrow \infty} \mathbb{E} \left[\tilde{\boldsymbol{\theta}}'_t \otimes \tilde{\boldsymbol{\theta}}'_t \right] =: \mathbf{S}_\infty.$$

Multiplying by \mathbf{H} and taking a trace, we get,

$$0 \leq \mathbb{E} \left\| \tilde{\boldsymbol{\theta}}'_0 \right\|_{\mathbf{H}}^2 \leq \mathbb{E} \left\| \tilde{\boldsymbol{\theta}}'_1 \right\|_{\mathbf{H}}^2 \leq \cdots \leq \lim_{t \rightarrow \infty} \mathbb{E} \left\| \tilde{\boldsymbol{\theta}}'_t \right\|_{\mathbf{H}}^2 = \text{Tr} [\mathbf{H} \mathbf{S}_\infty]. \quad (76)$$

Thus, $\tilde{\boldsymbol{\theta}}_t = \tilde{\boldsymbol{\theta}}'_t + \boldsymbol{\theta}_*$ is a process that starts from $\tilde{\boldsymbol{\theta}}_0 = \boldsymbol{\theta}_*$ and satisfies the conditions of Lemma D.12. This in turn gives

$$0 \leq \mathbb{E} \left[F(\tilde{\boldsymbol{\theta}}_0) - F(\boldsymbol{\theta}_*) \right] \leq \mathbb{E} \left[F(\tilde{\boldsymbol{\theta}}_1) - F(\boldsymbol{\theta}_*) \right] \leq \cdots \leq \lim_{t \rightarrow \infty} \mathbb{E} \left[F(\tilde{\boldsymbol{\theta}}_t) - F(\boldsymbol{\theta}_*) \right] = \frac{1}{2} \text{Tr} [\mathbf{H} \mathbf{S}_\infty], \quad (77)$$

which equals $F_\infty^*(\boldsymbol{\beta})$ by definition.

Combining both processes: From the triangle inequality of the norm $\mathbf{u} \mapsto \sqrt{\mathbb{E} \|\mathbf{u}\|_{\mathbf{H}}^2}$, we get

$$\sqrt{\mathbb{E} \|\boldsymbol{\theta}'_t\|_{\mathbf{H}}^2} \leq \sqrt{\mathbb{E} \|\hat{\boldsymbol{\theta}}_t\|_{\mathbf{H}}^2} + \sqrt{\mathbb{E} \|\tilde{\boldsymbol{\theta}}'_t\|_{\mathbf{H}}^2}.$$

Plugging in (75) and (76) gives

$$\begin{aligned} \sqrt{\mathbb{E} [F(\boldsymbol{\theta}_t) - F(\boldsymbol{\theta}_*)]} &\leq \sqrt{\frac{L}{2\mu} \exp(-\eta \mu t) \|\hat{\boldsymbol{\theta}}_0\|_{\mathbf{H}}^2} + \sqrt{\frac{1}{2} \text{Tr} [\mathbf{H} \mathbf{S}_\infty]} \\ &= \sqrt{\frac{L}{\mu} \exp(-\eta \mu t) (F(\boldsymbol{\theta}_0) - F(\boldsymbol{\theta}_*))} + \sqrt{F_\infty^*(\boldsymbol{\beta})}, \end{aligned}$$

where the last equality followed from (77). This establishes the required statement with F_∞^* in place of F^∞ . Taking $t \rightarrow \infty$, we see that

$$\sqrt{F_\infty(\boldsymbol{\beta})} = \lim_{t \rightarrow \infty} \sqrt{\mathbb{E} [F(\boldsymbol{\theta}_t) - F(\boldsymbol{\theta}_*)]} = \sqrt{F_\infty^*(\boldsymbol{\beta})},$$

for any fixed η or that $F_\infty = F_\infty^*$ irrespective of $\boldsymbol{\theta}_0$. \square

D.4 Privacy-Utility Guarantees of DP-FTRL

We now state a general privacy-utility bound for DP-FTRL in terms of the asymptotics of Noisy-FTRL run with the same parameters.

Theorem D.13. *Fix a constant $0 < p < 1$ and suppose the Assumption D.1 holds. Fix some noise coefficients $\boldsymbol{\beta} = (\beta_0, \dots, \beta_{T-1})$ that satisfy Half-Expo Decay with parameter $\eta \tilde{\nu}$ for some $\tilde{\nu} \leq \mu$. Consider the sequence $(\boldsymbol{\theta}_t)_{t=0}^{T-1}$ of iterates and the sequence $(\mathbf{g}_t)_{t=0}^{T-1}$ of gradients when running DP-FTRL for T iterations with noise coefficients $\boldsymbol{\beta}$, gradient clip norm $G = cR^2 \max \left\{ \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}_*\|_2, \sqrt{\eta R^2 \sigma_{\text{sgd}}^2 / \mu}, \sigma_{\text{sgd}} / R \right\} \log^{5/2} \left(\frac{T}{p} \right)$, and a learning rate*

$$\eta \leq \min \left\{ \frac{1}{CR^2 \log(T/p)}, \frac{\tilde{\nu} \rho}{8C^2 R^4 d \gamma_\infty(\boldsymbol{\beta})^2 \|\boldsymbol{\beta}\|_1^2 \log^5(T/p)} \right\},$$

and DP noise $\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \sigma_{\text{dp}}^2 G^2 \mathbf{I})$ with squared noise multiplier $\sigma_{\text{dp}}^2 = \gamma(\boldsymbol{\beta})^2 / (2\rho)$. Then, we have the following:

- (a) $(\boldsymbol{\theta}_t)_{t=0}^T$ is ρ -zCDP.
(b) Let \mathcal{E} denote the event where no gradients are clipped, i.e, $\mathcal{E} = \cap_{t=0}^{T-1} \{\|\mathbf{g}_t\|_2 \leq G\}$. We have, $\mathbb{P}(\mathcal{E}) \geq 1 - p$.
(c) We have,

$$\mathbb{E}[(F(\boldsymbol{\theta}_t) - F(\boldsymbol{\theta}_*)) \cdot \mathbb{1}(\mathcal{E})] \leq \frac{2L}{\mu} \exp(-\eta\mu t) (F(\boldsymbol{\theta}_0) - F(\boldsymbol{\theta}_*)) + 2\hat{F}_\infty(\boldsymbol{\beta}),$$

where $\hat{F}_\infty(\boldsymbol{\beta})$ is the asymptotic suboptimality of Noisy-FTRL run with the same parameters.

Proof. Part (a) follows from Theorem 2.1. For part (b), we bound the gradient norms from Theorem D.4 as

$$\begin{aligned} \|\mathbf{g}_t\|_2 &\leq CR^2 \left(\|\boldsymbol{\theta}'_0\|_2 + \sqrt{\frac{\eta R^2 \sigma_{\text{sgd}}^2}{\mu}} + \frac{\sigma_{\text{sgd}}}{R} + G \sqrt{\frac{\eta \sigma^2 d \|\boldsymbol{\beta}\|_1^2}{\tilde{\nu}}} \right) \log^{5/2} \left(\frac{T}{p} \right) \\ &\leq CR^2 \left(\|\boldsymbol{\theta}'_0\|_2 + \sqrt{\frac{\eta R^2 \sigma_{\text{sgd}}^2}{\mu}} + \frac{\sigma_{\text{sgd}}}{R} \right) \log^{5/2} \left(\frac{T}{p} \right) + \frac{G}{4} \\ &\leq 4 \max \left\{ CR^2 \max \left\{ \|\boldsymbol{\theta}'_0\|_2, \sqrt{\frac{\eta R^2 \sigma_{\text{sgd}}^2}{\mu}}, \frac{\sigma_{\text{sgd}}}{R} \right\} \log^{5/2} \left(\frac{T}{p} \right), \frac{G}{4} \right\} \leq G \end{aligned}$$

where the second inequality followed from the condition on the learning rate and we take $c = 4C$ in the definition of G for the last inequality. Thus, \mathcal{E} holds whenever the bound of Theorem D.4 holds, so we have $\mathbb{P}(\mathcal{E}) \geq 1 - p$.

For part (c), consider the sequence $(\phi_t)_{t=0}^T$ produced by running Noisy-FTRL with $\phi_0 = \boldsymbol{\theta}_0$ and the same realizations $(\mathbf{x}_t, \xi_t, \mathbf{w}_t)$ of random inputs, linear model noise, and DP noise. On \mathcal{E} , we have that $\phi_t = \boldsymbol{\theta}_t$ for all t . Thus, we have,

$$\mathbb{E}[(F(\boldsymbol{\theta}_t) - F(\boldsymbol{\theta}_*)) \cdot \mathbb{1}(\mathcal{E})] = \mathbb{E}[(F(\phi_t) - F(\boldsymbol{\theta}_*)) \cdot \mathbb{1}(\mathcal{E})] \leq \mathbb{E}[F(\phi_t) - F(\boldsymbol{\theta}_*)],$$

since $\mathbb{1}(\mathcal{E}) \leq 1$. This can now be bounded using Theorem D.11 to complete the proof. \square

We can instantiate these rates for DP-SGD and DP-FTRL. Recall that we have $\kappa = L/\mu$, $d_{\text{eff}} = \text{Tr}[\mathbf{H}]/L$, and $R^2 = \Theta(\text{Tr}[\mathbf{H}])$.

Corollary D.14. *Consider the setting of Theorem D.13 with T large enough that $T/\log^5(T/p) \geq c\kappa^2 d_{\text{eff}}^2 d/\rho$. The the final suboptimality of DP-SGD at an appropriate choice of the learning rate is (ignoring absolute constants),*

$$\begin{aligned} \mathbb{E}[(F(\boldsymbol{\theta}_T) - F(\boldsymbol{\theta}_*)) \cdot \mathbb{1}(\mathcal{E})] &\leq \frac{L}{\mu} \exp \left(-\frac{\rho T}{c\kappa^2 d_{\text{eff}}^2 d \log^5(T/p)} \right) \\ &\quad + \kappa d_{\text{eff}} \left(\frac{d \text{Tr}[\mathbf{H}] \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}_*\|_2^2}{\rho T} + \frac{d\sigma_{\text{sgd}}^2}{\rho T} + \frac{\sigma_{\text{sgd}}^2}{T} \right) \text{polylog}(T). \end{aligned}$$

Proof. We plug in the asymptotic suboptimality bound of Noisy-SGD into the bound of Theorem D.13. We get two terms depending on the learning rate η : the first $\exp(-\mu\eta T)$ term and the second $O(\eta)$ term coming from the asymptotic suboptimality. We balance both the terms subject to the maximum bound on η using Lemma F.21 to get

$$\begin{aligned} \mathbb{E}[(F(\boldsymbol{\theta}_T) - F(\boldsymbol{\theta}_*)) \cdot \mathbb{1}(\mathcal{E})] &\leq \frac{L}{\mu} \exp \left(-\frac{\rho\mu^2 T}{cR^4 d \log^5(T/p)} \right) \\ &\quad + \frac{\text{polylog}(T)}{\mu T} \left(\frac{dR^4 \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}_*\|_2^2}{\rho} + \frac{d\sigma_{\text{sgd}}^2 R^2}{\rho} + \sigma_{\text{sgd}}^2 R^2 \right). \end{aligned}$$

Rearranging the constants completes the proof. \square

Corollary D.15. Consider the setting of Theorem D.13 with T large enough that $T/\log^7(T/p) \geq \frac{c\kappa^2 d_{\text{eff}}^2 d}{\rho} \log\left(\frac{c\kappa^2 d_{\text{eff}}^2 d}{\rho}\right)$. For ν -DP-FTRL with an appropriate choice of the parameter ν and learning rate η , we have (ignoring absolute constants),

$$\begin{aligned} \mathbb{E}[(F(\boldsymbol{\theta}_T) - F(\boldsymbol{\theta}_*)) \cdot \mathbb{1}(\mathcal{E})] &\leq \frac{L}{\mu} \exp\left(-\frac{\rho T}{c\kappa^2 d_{\text{eff}}^2 d \log^7(T/p) \log(\kappa^2 d_{\text{eff}}^2 d/\rho)}\right) \\ &\quad + \kappa d_{\text{eff}} \left(\frac{\kappa d_{\text{eff}} \text{Tr}[\mathbf{H}] \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}_*\|_2^2}{\rho T^2} + \frac{\kappa d_{\text{eff}} \sigma_{\text{sgd}}^2}{\rho T^2} + \frac{\sigma_{\text{sgd}}^2}{T} \right) \text{polylog}(T). \end{aligned}$$

Proof. We plug in the asymptotic error for ν -Noisy-FTRL from Proposition C.20 into Theorem D.13 to get that

$$\mathbb{E}[(F(\boldsymbol{\theta}_T) - F(\boldsymbol{\theta}_*)) \cdot \mathbb{1}(\mathcal{E})] \leq \frac{L}{\mu} \exp(-\mu\eta T) + \eta \sigma_{\text{sgd}}^2 R^2 + \eta^2 \frac{R^2 G^2}{\rho} \log^2 \frac{1}{\eta\mu}, \quad (78)$$

where G^2 is as given in the statement of Theorem D.13. For our choice of β , we have $\|\beta\|_1^2 \leq 4$ always and $\gamma(\beta)^2 \leq 5 \log(1/\eta\mu)$ from Equation (42) (from the proof of Proposition C.20). Thus, the largest learning rate permitted must satisfy

$$\eta \log^2 \frac{1}{\eta\mu} \leq \frac{\eta\rho}{cR^2 d \log^5(T/p)}.$$

From Lemma F.22, we can ensure with a more stringent condition

$$\eta \leq \frac{\mu\rho}{cR^4 d \log^5(T/p) \log^2(cR^4 d \log(T/p)/(\mu^2\rho))}.$$

Finally, this is implied by imposing the requirement

$$\eta \leq \frac{\mu\rho}{cR^4 d \log^7(T/p) \log\left(\frac{R^4 d}{\mu^2\rho}\right)} =: \eta_{\max}.$$

We now tune η to minimize the bound (78) subject to $\eta \leq \eta_{\max}$ using Lemma F.21. Thus gives,

$$\begin{aligned} \mathbb{E}[(F(\boldsymbol{\theta}_T) - F(\boldsymbol{\theta}_*)) \cdot \mathbb{1}(\mathcal{E})] &\leq \frac{L}{\mu} \exp\left(-\frac{\rho\mu^2 T}{cR^4 d \log^7(T/p) \log\frac{R^4 d}{\rho\mu^2}}\right) \\ &\quad + \frac{\text{polylog}(T)}{\mu T} \left(\frac{R^6 \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}_*\|_2^2}{\rho\mu T} + \frac{R^4 \sigma_{\text{sgd}}^2}{\rho\mu^2 T^2} + \sigma_{\text{sgd}}^2 R^2 \right). \end{aligned}$$

Rewriting the constants completes the proof. \square

E Proofs for General Strongly Convex Functions

We now generalize §4 to general strongly convex problems. Here, we bound the asymptotic suboptimality of DP-FTRL and DP-SGD by the value of a convex program.

Theorem E.1. Suppose $f(\cdot; \mathbf{z})$ is G -Lipschitz, and the stochastic gradients are uniformly bounded as $\|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}; \mathbf{z}) - \mathbb{E}_{\mathbf{z}' \sim \mathbb{P}_{\text{data}}}[\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}; \mathbf{z}')]\|_2 \leq \sigma_{\text{sgd}}$. Then, if F is μ -strongly convex and L -smooth, we can bound F_{∞} for any noise correlation $B(\omega)$ in the frequency domain as:

$$\min_{\substack{\psi: [-\pi, \pi] \rightarrow \mathbb{R}_+ \\ \psi \in \mathcal{C}(\eta, L, \mu)}} \frac{Ld}{2\pi} \int_{-\pi}^{\pi} (G^2 \rho^{-1} |B(\omega)|^2 \gamma_{\infty}(B)^2 + \sigma_{\text{sgd}}^2) \psi(\omega) d\omega, \quad (79)$$

where $\gamma_{\infty}(B)$ is the limiting sensitivity in the Fourier domain, and $\mathcal{C}(\eta, \mu, L)$ is a convex set (details in Appendix E).

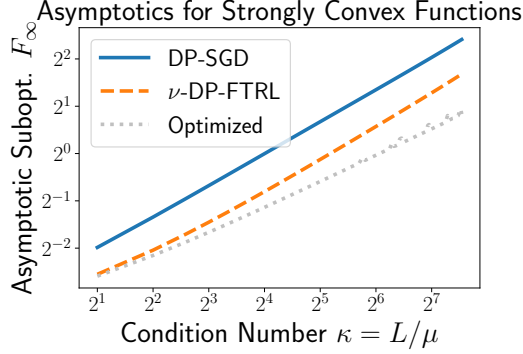


Figure 4: **DP-FTRL attains a tighter bound on F_∞** with the growing condition number. Here, “Optimized” approximately minimizes (79). The plots holds for smooth and strongly convex functions ($L = 1 = G, \sigma_{\text{sgd}} = 0$).

While technically an infinite dimensional optimization problem over the function ψ , we can approximate the solution by discretizing ψ into k points uniformly over $[-\pi, \pi]$. Further, if we discretize B similarly, we can obtain a **second-order cone program** with k conic constraints and $O(k)$ decision variables. As $k \rightarrow \infty$, the solution approaches the solution to (79). Empirically, we observe that the values stabilize quickly as k increases. We stop the computation when the change in bound as a function of k drops below a threshold. We use $k = 1000$ for.

Further, given the optimal $\psi = \psi^*$, we can run an alternating minimization where we minimize the objective of (79) with respect to ψ for fixed B and with respect to B for fixed ψ . This leads to an iteratively improving choice of B . We find empirically that this iterative procedure converges quickly and leads to a provable theoretical gap between the upper bounds on F_∞ achievable by DP-SGD and DP-FTRL.

We numerically compare the bound (79) for DP-SGD and ν -DP-FTRL with weights from (7). Figure 4 shows that the gap between DP-SGD and ν -DP-FTRL is multiplicative, i.e., the absolute gap grows with the increasing condition number $\kappa = L/\mu$ (which reflects practical scenarios). The suboptimality of “Optimized” DP-FTRL (optimized as described above) grows even more slowly with κ .

Overall, ν -DP-FTRL significantly improves upon DP-SGD and has only a single tunable parameter ν and no expensive computation to generate the noise correlations. We focus on ν -DP-FTRL for experiments in this paper, but leave the possibility of improving results further based on Optimized DP-FTRL for future work.

E.1 Proofs

We prove the results from Theorem E.1. Under the assumptions of the theorem, clipping does not occur in DP-FTRL so the updates can be written as

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta((B\mathbf{w})_t + (\mathbf{g}_t + \hat{\mathbf{w}}_t)) \quad (80)$$

where

$$\mathbf{g}_t = \nabla F(\boldsymbol{\theta}_t), \hat{\mathbf{w}}_t = \nabla f(\boldsymbol{\theta}_t; \mathbf{z}_t) - \mathbb{E}_{\mathbf{z} \sim \mathbb{P}_{\text{data}}}[\nabla f(\boldsymbol{\theta}_t; \mathbf{z})]$$

where $\hat{\mathbf{w}}_t$ is a random variable that, conditioned on $\boldsymbol{\theta}_t$, is bounded by σ_{sgd} with probability 1.

Let \mathbf{I}_d denote the $d \times d$ identity matrix.

Theorem E.2. Let $\boldsymbol{\lambda} = \{\lambda_t\}_{t=-\infty}^{\infty}$ be such that $\lambda_t \geq 0 \quad \forall t \in \mathbb{Z}$,

$$\sum_{t=-\infty}^{\infty} \lambda_t \leq 2\lambda_0$$

and let Λ denote the Discrete-time Fourier transform (DTFT) of λ . Let

$$M_\lambda(\omega) = A(\omega)^* \widetilde{M}_\lambda(\omega) A(\omega) \quad (81a)$$

$$A(\omega) = \begin{pmatrix} \eta \mathbf{I}_d & 0 \\ (1 - \exp(i\omega)) \mathbf{I}_d & -\eta \mathbf{I}_d \end{pmatrix} \quad (81b)$$

$$\widetilde{M}_\lambda(\omega) = \begin{pmatrix} -\mu L (\Lambda(\omega) + \Lambda(\omega)^*) \mathbf{I}_d & \mu \Lambda(\omega) \mathbf{I}_d + L \Lambda(\omega)^* \mathbf{I}_d \\ \mu \Lambda^*(\omega) \mathbf{I}_d + L \Lambda(\omega) \mathbf{I}_d & -(\Lambda(\omega) + \Lambda(\omega)^*) \mathbf{I}_d \end{pmatrix} \quad (81c)$$

Then, for any non-negative valued function $\psi : [-\pi, \pi] \mapsto \mathbb{R}_+$ such that

$$M_\lambda(\omega) \preceq \begin{pmatrix} -\eta^2 \mathbf{I}_d & 0 \\ 0 & \psi(\omega) \mathbf{I}_d \end{pmatrix} \quad \forall \omega \in [-\pi, \pi] \quad (82)$$

We have that

$$\lim_{t \rightarrow \infty} \mathbb{E} \left[\frac{\sum_{t=-T}^T \|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_2^2}{2T+1} \right] \leq \frac{2d}{2\pi\eta^2} \int_{-\pi}^{\pi} \left(|B(\omega)|^2 G^2 \rho^{-1} \gamma_\infty(B)^2 + \sigma_{\text{sgd}}^2 \right) \psi(\omega) d\omega$$

where S_{sgd} is the power spectral density of $\widetilde{\boldsymbol{w}}$. In particular, if the density of $\boldsymbol{\theta}_t$ converges to a stationary distribution, the expected value of

$$\lim_{t \rightarrow \infty} \mathbb{E} \left[\|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_2^2 \right]$$

under the stationary distribution is bounded as above.

Proof. We assume without loss of generality that $\nabla F(0) = 0$ so that the origin is the global optimum of F (else we can translate the origin to achieve this).

Since $\mathbf{g} = \nabla F(\boldsymbol{\theta})$ satisfies

$$\langle \mathbf{g} - L\boldsymbol{\theta}, \mu\boldsymbol{\theta} - \mathbf{g} \rangle \geq 0 \quad \forall \boldsymbol{\theta}, \mathbf{g}$$

Then, we can write down the following family of integral quadratic constraints relating $\mathbf{g} = (\dots, \mathbf{g}_0, \mathbf{g}_1, \mathbf{g}_2, \dots)$ and $\boldsymbol{\theta} = (\dots, \boldsymbol{\theta}_0, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots)$ in terms of their Fourier transforms $\Theta(\omega), G(\omega)$ ([23] equations 27-29):

$$\int_{-\pi}^{\pi} \begin{pmatrix} \Theta(\omega) \\ G(\omega) \end{pmatrix}^* \begin{pmatrix} -\mu L (\Lambda(\omega) + \Lambda(\omega)^*) \mathbf{I}_d & \mu (\Lambda(\omega)) \mathbf{I}_d + L (\Lambda(\omega)^*) \mathbf{I}_d \\ \mu (\Lambda^*(\omega)) \mathbf{I}_d + L (\Lambda(\omega)) \mathbf{I}_d & -(\Lambda(\omega) + \Lambda(\omega)^*) \mathbf{I}_d \end{pmatrix} \begin{pmatrix} \Theta(\omega) \\ G(\omega) \end{pmatrix} d\omega \geq 0 \quad (83)$$

Noting that from (80), we have that

$$\Theta(\omega) (\exp(i\omega) - 1) = -\eta (G(\omega) + Z(\omega)) \implies G(\omega) = \left(\frac{1 - \exp(i\omega)}{\eta} \right) \Theta(\omega) - Z(\omega)$$

where Z denotes the DTFT of $\boldsymbol{\zeta} = \mathbf{B}\mathbf{w} + \widehat{\mathbf{w}}$. Plugging this into the above quadratic constraint and multiplying by η^2 , we obtain

$$\int_{-\pi}^{\pi} \begin{pmatrix} \Theta(\omega) \\ Z(\omega) \end{pmatrix}^* M_\lambda(\omega) \begin{pmatrix} \Theta(\omega) \\ Z(\omega) \end{pmatrix} d\omega \geq 0 \quad (84)$$

Since $M_\lambda(\omega) \preceq \begin{pmatrix} -\eta^2 \mathbf{I}_d & 0 \\ 0 & \psi(\omega) \mathbf{I}_d \end{pmatrix}$ we obtain that

$$\begin{aligned} \int_{-\pi}^{\pi} \begin{pmatrix} \Theta(\omega) \\ Z(\omega) \end{pmatrix}^* \begin{pmatrix} -\eta^2 \mathbf{I}_d & 0 \\ 0 & \psi(\omega) \mathbf{I}_d \end{pmatrix} \begin{pmatrix} \Theta(\omega) \\ Z(\omega) \end{pmatrix} d\omega \geq 0 &\implies \frac{\mathbb{E} \left[\int_{-\pi}^{\pi} \|\Theta(\omega)\|^2 \right]}{\mathbb{E} \left[\int_{-\pi}^{\pi} \|\sqrt{\psi(\omega)} Z(\omega)\|^2 \right]} \leq 1 \\ \implies \frac{\lim_{T \rightarrow \infty} \mathbb{E} \left[\frac{\sum_{t=-T}^T \|\boldsymbol{\theta}_t\|^2}{2T+1} \right]}{\lim_{T \rightarrow \infty} \mathbb{E} \left[\frac{\sum_{t=-T}^T \|\sqrt{\psi}[\boldsymbol{\zeta}](t)\|^2}{2T+1} \right]} &\leq \frac{1}{\eta^2} \end{aligned}$$

where $\sqrt{\zeta}[z]$ denotes the *LTI* operator with transfer function $\sqrt{\zeta}(\omega)$ applied to the signal ζ .

The denominator of the final line above is the power spectral density of $\sqrt{\kappa}[\zeta]$ (since $\sqrt{\kappa}[\zeta]$ is a wide-sense stationary stochastic process). By the Cauchy-Schwartz inequality for random variables, this is bounded above by

$$2d \left(|B(\omega)|^2 \rho^{-1} \gamma_\infty(B)^2 + \sigma_{\text{sgd}}^2 \right) \psi(\omega)$$

where the first term in brackets is the power spectral density of the Gaussian random process $B\hat{w}$ and the second term is an upper bound on the power spectral density of \hat{w} . Hence, by Theorem F.2, we have the result. \square

Full Proof of Theorem E.1: We are now ready to prove the theorem.

Proof of Theorem E.1. Given the above theorem and smooth convexity parameter L , we know that the asymptotic suboptimality F_∞ is bounded above by

$$\frac{2Ld}{2\pi\eta^2} \int_{-\pi}^{\pi} \left(|B(\omega)|^2 \rho^{-1} \gamma_\infty(B)^2 G^2 + \sigma_{\text{sgd}}^2 \right) \psi(\omega) d\omega$$

Now, the constraint (82) can be rewritten as

$$\begin{pmatrix} -\eta^2 & 0 \\ 0 & \psi(\omega) \end{pmatrix} - \begin{pmatrix} \eta & 0 \\ 1 - \exp(i\omega) & -\eta \end{pmatrix}^{*\top} \begin{pmatrix} -\mu L (\Lambda(\omega) + \Lambda(\omega)^*) & \mu \Lambda(\omega) + L \Lambda(\omega)^* \\ \mu \Lambda^*(\omega) + L \Lambda(\omega) & -(\Lambda(\omega) + \Lambda(\omega)^*) \end{pmatrix} \begin{pmatrix} \eta & 0 \\ 1 - \exp(i\omega) & -\eta \end{pmatrix} \succeq 0 \quad (85)$$

since all the matrices involved are Hadamard products of the 2×2 matrices above and the identity matrix.

Thus, for each ω , $\psi(\omega)$ must satisfy a 2×2 PSD constraint which can be rewritten as a Second Order Cone Program (SOCP) constraint. Furthermore, the constraint on λ from theorem E.2 is a linear constraint. Since the projection of a convex set in ψ, λ to ψ is convex, ψ belongs to a convex set. Furthermore, if we take λ such that $\lambda_\tau = 0$ for $|\tau| > T_{\max}$ for some $T_{\max} > 0$, the constraint on λ can be written as

$$2\lambda_0 \geq \sum_{\tau=-T_{\max}}^{T_{\max}} \lambda_\tau$$

Further, if we discretize ω to a uniform grid on $[-\pi, \pi]$, the constraints (85) can be written as a finite collection of SOCP constraints linking $\psi(\omega)$ and λ . \square

F Technical Definitions and Lemmas

We review a number of relevant technical definitions and lemmas here:

- **Appendix F.1:** Fourier Analysis of Linear Time-Invariant Systems.
- **Appendix F.2:** Stationary covariance of SGD.
- **Appendix F.3:** Concentration of Measure.
- **Appendix F.4:** Review of definitions and useful properties of elliptic integrals.

F.1 Linear Time-Invariant (LTI) Systems

We first review the definition and some useful properties of discrete-time Linear Time-Invariant (LTI) systems. We refer to the textbook [43] for a more detailed description.

Definition F.1. An input-output system $\mathbf{y}_t = \mathcal{A}_t(\mathbf{x})$ with an input sequence $\mathbf{x} = (\mathbf{x}_t)_{t=-\infty}^{\infty}$ in some input space \mathcal{X} and an output sequence $(\mathbf{y}_t)_{t=-\infty}^{\infty}$ in an output space \mathcal{Y} is said to be *LTI* if satisfies two properties:

- **Linearity:** For any \mathcal{X} -valued sequences $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$ and scalars $\alpha_1, \alpha_2, \dots$, we have

$$\mathcal{A}_t \left(\sum_{j=1}^{\infty} \alpha_j \mathbf{x}^{(j)} \right) = \sum_{j=1}^{\infty} \alpha_j \mathcal{A}_t(\mathbf{x}^{(j)}).$$

- **Time-Invariance:** For any $t_0 \in \mathbb{Z}$, the sequence \mathbf{x}' defined as $\mathbf{x}'_t := \mathbf{x}_{t-t_0}$ satisfies $\mathcal{A}_t(\mathbf{x}') = \mathcal{A}_{t-t_0}(\mathbf{x})$.

LTI systems can be viewed as linear operators defined on the Hilbert space of *signals*

$$\ell_{2e}^d = \left\{ (\mathbf{x}_t)_{t=-\infty}^{\infty} : \sum_{\tau=-t}^t \|\mathbf{x}_\tau\|^2 < \infty \quad \forall t \in \mathbb{Z} \right\}$$

with the inner product $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{t=-\infty}^{\infty} \mathbf{x}_t^\top \mathbf{y}_t$ with the convention $\ell_{2e} = \ell_{2e}^1$. LTI systems can be described in linear algebraic notation by the action of an infinite Toeplitz matrix on an element of ℓ_{2e} :

$$\mathbf{y} = \mathbf{H}\mathbf{x} \implies \mathbf{y}_t = \sum_{\tau=-\infty}^{\infty} \mathbf{H}_{t,\tau} \mathbf{x}_\tau = \sum_{\tau=-\infty}^{\infty} \mathbf{h}_{t-\tau} \mathbf{x}_\tau = (\mathbf{h} \star \mathbf{x})_t \quad \forall t \in \mathbb{Z}.$$

This property is represented more elegantly in the Fourier domain. Consider the discrete-time Fourier transform (DTFT)

$$\mathbf{X}(\omega) = \sum_{t=-\infty}^{\infty} \mathbf{x}_t \exp(-i\omega t),$$

of \mathbf{x}_t and the corresponding DTFTs $\mathbf{Y}(\omega), \mathbf{G}(\omega)$ of \mathbf{y}, \mathbf{h} respectively. Then, we have $\mathbf{Y}(\omega) = \mathbf{G}(\omega)\mathbf{X}(\omega)$. Here, \mathbf{h} is known as the **impulse response** and $\mathbf{G}(\omega)$ is known as the **transfer function**.

An LTI system is said to be asymptotically stable if its output decays to zero for any input sequence that is bounded, ie, for which there exists $T > 0$ such that $x_t = 0 \quad \forall t > T$.

Variance of LTI systems driven by white noise: The Fourier-domain analysis of an LTI system (particularly its transfer function) helps us characterize the covariance of the output \mathbf{y}_t as a function of the covariance of the input \mathbf{x}_t . For simplicity, we assume that $\mathbf{x}_t \in \mathbb{R}^d$ and $\mathbf{y}_t \in \mathbb{R}^p$ so that \mathbf{H}_t and $\mathbf{G}(\omega)$ can be identified with matrices in $\mathbb{R}^{p \times d}$.

Theorem F.2. Consider an asymptotically-stable LTI system with \mathbb{R}^d -valued inputs $(\mathbf{x}_t)_{t=-\infty}^{\infty}$ and \mathbb{R}^p -valued outputs $(\mathbf{y}_t)_{t=-\infty}^{\infty}$ and a transfer function $\mathbf{G}(\omega) \in \mathbb{R}^{p \times d}$. Suppose that \mathbf{x}_t is a stationary white noise sequence with covariance Σ , i.e., $\mathbb{E}[\mathbf{x}_t] = \mathbf{0}$ and $\mathbb{E}[\mathbf{x}_t \otimes \mathbf{x}_\tau] = \Sigma$ if $t = \tau$ and $\mathbf{0}$ otherwise for all t, τ . Then, we have for all $t > -\infty$ that

$$\mathbb{E}[\mathbf{y}_t \otimes \mathbf{y}_t] = \frac{1}{2\pi} \int_{-\pi}^{\pi} \mathbf{G}(\omega) \Sigma \mathbf{G}(\omega)^* d\omega.$$

F.2 Stationary Covariance of Stochastic Gradient Descent for Linear Regression

We now give a result characterizing the stationary covariance of SGD for linear regression [7, 18, 29, 30].

Theorem F.3 (Lemma 5 of [29]). Given a fixed non-random vector $\delta_0 \in \mathbb{R}^d$, consider the recursion

$$\delta_{t+1} = (\mathbf{I} - \eta \mathbf{x}_t \otimes \mathbf{x}_t) \delta_t + \eta \zeta_t,$$

for all $t \geq 0$ where

- \mathbf{x}_t are i.i.d. with mean $\mathbf{0}$, covariance \mathbf{H} , and
- ζ_t are i.i.d. with mean $\mathbf{0}$, covariance $\mathbb{E}[\zeta_t \otimes \zeta_t] \preceq \nu^2 \mathbf{H}$.

Further, if $\mathbb{E} \left[\|\mathbf{x}_t\|_2^2 (\mathbf{x}_t \otimes \mathbf{x}_t) \right] \preceq R^2 \mathbf{H}$ and $\eta < 1/R^2$, then we have for all $t \geq 0$.

$$\mathbb{E}[\delta_t \otimes \delta_t] \preceq \frac{\eta \nu^2}{1 - \eta R^2} \mathbf{I}.$$

F.3 Concentration of Measure

We recall the definition of sub-Gaussian random variables and list some useful concentration inequalities.

Definition F.4. A real-valued random variable X is said to be sub-Gaussian with variance proxy σ^2 if for all $\lambda \in \mathbb{R}$, we have

$$\mathbb{E}[\exp(\lambda(X - \mu))] \leq \exp(\lambda^2 \sigma^2 / 2),$$

where $\mu = \mathbb{E}[X]$. If in addition, the variance of X exactly equals σ^2 , it is said to be strictly sub-Gaussian.

The cumulants of strict sub-Gaussian random variables are closely related to those of a Gaussian [6, Prop. 3.2].

Property F.5. If X is strictly sub-Gaussian with mean zero and variance σ^2 , we have $\mathbb{E}[X^3] = 0$ and $\mathbb{E}[X^4] \leq 3\sigma^4 = \mathbb{E}[Y^4]$ for $Y \sim \mathcal{N}(0, \sigma^2)$.

Next, we state the **Hanson-Wright inequality** for the concentration of quadratic forms; see e.g. [46].

Lemma F.6. Let $\xi = (\xi_1, \dots, \xi_d)$ be such that each ξ_j is independent and sub-Gaussian with mean zero and variance proxy σ^2 . Then, we have for any matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$,

$$\mathbb{P}(\langle \xi, \mathbf{A}\xi \rangle - \mathbb{E}[\langle \xi, \mathbf{A}\xi \rangle] > t) \leq \exp\left(-c \min\left\{\frac{t^2}{\sigma^4 \|\mathbf{A}\|_F^2}, \frac{t}{\sigma^2 \|\mathbf{A}\|_2}\right\}\right),$$

for a universal constant c . Consequently, for any $\rho < 1/3$ and symmetric PSD matrix \mathbf{A} , we have with probability $1 - \rho$ that

$$\langle \xi, \mathbf{A}\xi \rangle \leq C\sigma^2 \left(\text{Tr}[\mathbf{A}] \sqrt{\log \frac{1}{\rho}} + \|\mathbf{A}\|_2 \log \frac{1}{\rho} \right) \leq C'\sigma^2 \text{Tr}[\mathbf{A}] \log \frac{1}{\rho},$$

for universal constants C, C' .

The second part follows from the first one under the simplifications $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F \leq \text{Tr}[\mathbf{A}]$ and $\mathbb{E}[\langle \xi, \mathbf{A}\xi \rangle] \leq \sigma^2 \text{Tr}[\mathbf{A}]$ for \mathbf{A} PSD.

Remark F.7. Explicit values for the constant c in Lemma F.6 (and thus for C, C') are known for the case when $\xi_1, \dots, \xi_d \sim \mathcal{N}(0, \sigma^2)$: $c \approx 0.1457 \geq 1/8$, $C \leq 8$, $C' \leq 16$ [42].

F.4 Review of Elliptic Integrals

We recall some definitions and useful properties of elliptic integrals. We refer to [1, §19] and [11] for details.

The three canonical elliptic integral forms are:

- (i) The complete elliptic integral of the first kind $K : (0, 1) \rightarrow [0, \infty)$ is

$$K(k) := \int_0^{\pi/2} \frac{d\omega}{\sqrt{1 - k^2 \sin^2(\omega)}}. \quad (86)$$

- (ii) The complete elliptic integral of the second kind $E : (0, 1) \rightarrow [0, \infty)$ is

$$E(k) := \int_0^{\pi/2} \sqrt{1 - k^2 \sin^2(\omega)} \, d\omega. \quad (87)$$

- (iii) The complete elliptic integral of the third kind $\Pi : (\mathbb{R} \setminus \{\pm 1\}) \times (0, 1) \rightarrow \mathbb{R}$ is denoted conventionally as $\Pi(\alpha^2, k)$ where α^2 is allowed to take negative values. It is defined as

$$\Pi(\alpha^2, k) := \int_0^{\pi/2} \frac{d\omega}{(1 - \alpha^2 \sin^2(\omega))\sqrt{1 - k^2 \sin^2(\omega)}}. \quad (88)$$

The corresponding integrals where $1 - k^2 \sin^2(\omega)$ is replaced with $1 + k^2 \sin^2(\omega)$ can also be expressed using the elliptic integrals [1, Eq. (19.7.2), (19.7.5)].

Property F.8. For any $m \in (0, 1)$, we have

$$\int_0^{\pi/2} \frac{d\omega}{\sqrt{1 + m \sin^2(\omega)}} = \frac{1}{\sqrt{1+m}} K\left(\sqrt{\frac{m}{1+m}}\right). \quad (89)$$

Property F.9. For any $m \in (0, 1)$ and any $\alpha^2 \in \mathbb{R} \setminus \{\pm 1\}$ such that $\alpha^2 + m \neq 0$, we have

$$\begin{aligned} & \int_0^{\pi/2} \frac{d\omega}{(1 - \alpha^2 \sin^2(\omega))\sqrt{1 + m \sin^2(\omega)}} \\ &= \frac{m}{(m + \alpha^2)\sqrt{1+m}} K\left(\sqrt{\frac{m}{1+m}}\right) + \frac{\alpha^2}{(m + \alpha^2)\sqrt{1+m}} \Pi\left(\frac{m + \alpha^2}{1+m}, \sqrt{\frac{m}{1+m}}\right). \end{aligned} \quad (90)$$

The next few properties are about the asymptotics of the elliptic integrals; see e.g. [1, Eq. (19.9.1)] for $K(\cdot)$ and [1, Eq. (19.12.4)] for Π .

Property F.10. For all $k \in (0, 1)$, we have

$$\log\left(\frac{4}{\sqrt{1-k^2}}\right) \leq K(k) \leq \left(1 + \frac{1-k^2}{4}\right) \log\left(\frac{4}{\sqrt{1-k^2}}\right) \leq \frac{5}{4} \log\left(\frac{4}{\sqrt{1-k^2}}\right).$$

Property F.11. For all $k, \alpha^2 \in (0, 1)$, we have

$$\Pi(\alpha^2, k) \leq \frac{1}{1-\alpha^2} \log\left(\frac{4}{\sqrt{1-k^2}}\right) \left(1 + O\left(\sqrt{1-k^2}\right)\right).$$

F.5 Useful Integrals

We list a number of useful definite integrals in this section.

Direct Evaluation: The first one is a cosine integral divided by a quadratic form.⁵

Lemma F.12. For reals $0 < |b| < a$ and an integer l , we have

$$\int_{-\pi}^{\pi} \frac{\cos(l\omega) d\omega}{a^2 + b^2 - 2ab \cos \omega} = \frac{2\pi}{a^2 - b^2} \left(\frac{b}{a}\right)^{|l|}.$$

The next lemma is also about rational cosine functions.⁶

Lemma F.13. For scalar a , we have

$$\int_{-\pi}^{\pi} \frac{d\omega}{1 + a \cos(\omega)} = \begin{cases} \frac{2\pi}{1-a^2}, & \text{if } |a| < 1, \\ +\infty, & \text{if } |a| = 1. \end{cases}$$

The next one is similar to the previous one.

Lemma F.14. We have that

$$\int_{-\pi}^{\pi} \frac{d\omega}{\sqrt{1 - \cos(\omega)}} = +\infty.$$

Proof. We successively deduce

$$\int_{-\pi}^{\pi} \frac{d\omega}{\sqrt{1 - \cos(\omega)}} = \frac{1}{\sqrt{2}} \int_{-\pi}^{\pi} \frac{d\omega}{|\sin(\omega/2)|} = 2\sqrt{2} \int_0^{\pi/2} \frac{d\omega}{\sin(\omega)} = +\infty,$$

where we used that $\int d\omega / \sin(\omega) = -\log |\csc(\omega) + \cot(\omega)| + C$. □

⁵See <https://math.stackexchange.com/a/816253>.

⁶See <https://math.stackexchange.com/a/1235309>.

Reductions to Elliptic Integrals: We now list several cosine integrals that can be reduced to elliptic integrals (see Appendix F.4 for their definitions).

Lemma F.15. For any $a \in (0, 1)$, we have

$$\int_{-\pi}^{\pi} \frac{d\omega}{|1 - a - \exp(i\omega)|} = \frac{4}{2 - a} K\left(\frac{\sqrt{1 - a}}{1 - a/2}\right), \quad (91)$$

where $K(\cdot)$ is the complete elliptic integral of the first kind, cf. (86).

Proof. Using $\cos(\omega) = 1 - 2 \sin^2(\omega/2)$ and the substitution $\omega' = \omega/2$, we successively deduce

$$\begin{aligned} \int_{-\pi}^{\pi} \frac{d\omega}{|1 - a - \exp(i\omega)|} &= 2 \int_0^{\pi} \frac{d\omega}{\sqrt{1 + (1 - a)^2 - 2(1 - a) \cos(\omega)}} \\ &= 2 \int_0^{\pi} \frac{d\omega}{\sqrt{a^2 + 4(1 - a) \sin^2(\omega/2)}} \\ &= \frac{4}{a} \int_0^{\pi/2} \frac{d\omega'}{\sqrt{1 + 4\left(\frac{1-a}{a^2}\right) \sin^2(\omega')}}. \end{aligned}$$

Applying Property F.8 to reduce this to the standard elliptic integral completes the proof. \square

The next lemma handles a more general case. Note that it recovers Lemma F.15 when $a = b$ since $\Pi(0, k) = K(k)$ by definition.

Lemma F.16. For any $a, b \in (0, 1)$, we have

$$\int_{-\pi}^{\pi} \frac{|1 - a - \exp(i\omega)|}{|1 - b - \exp(i\omega)|^2} d\omega = \frac{2a^2}{b^2(1 - a/2)} \Pi\left(\frac{b^2(1 - a) - a^2(1 - b)}{b^2(1 - a/2)^2}, \frac{\sqrt{1 - a}}{1 - a/2}\right), \quad (92)$$

where Π is the complete elliptic integral of the third kind, cf. (88).

Proof. We assume that $a \neq b$ to begin and handle the case of $a = b$ by continuity. Denote $h(a, \omega) = \sqrt{1 + (1 - a)^2 - 2(1 - a) \cos(\omega)}$

$$\begin{aligned} \int_{-\pi}^{\pi} \frac{|1 - a - \exp(i\omega)|}{|1 - b - \exp(i\omega)|^2} d\omega &= \int_{-\pi}^{\pi} \frac{|1 - a - \exp(i\omega)|^2}{|1 - a - \exp(i\omega)| |1 - b - \exp(i\omega)|^2} d\omega \\ &= \frac{1 + (1 - a)^2}{h(a, \omega) h(b, \omega)^2} - 2(1 - a) \frac{\cos(\omega)}{h(a, \omega) h(b, \omega)^2}. \end{aligned}$$

We next add and subtract terms to make the numerator of the second term read $h(b, \omega)^2$ to give

$$\int_{-\pi}^{\pi} \frac{|1 - a - \exp(i\omega)|}{|1 - b - \exp(i\omega)|^2} d\omega = \int_{-\pi}^{\pi} \frac{1 + (1 - a)^2 - \frac{1-a}{1-b} (1 + (1 - b)^2)}{h(a, \omega) h(b, \omega)^2} d\omega + \frac{1 - a}{1 - b} \int_{-\pi}^{\pi} \frac{d\omega}{h(a, \omega)}. \quad (93)$$

From Lemma F.15, the second term above can be written as

$$\frac{1 - a}{1 - b} \int_{-\pi}^{\pi} \frac{d\omega}{h(a, \omega)} = \frac{4(1 - a)}{(1 - b)(2 - a)} K\left(\frac{\sqrt{1 - a}}{1 - a/2}\right). \quad (94)$$

The first term of (93) can similarly be reduced to the elliptic integral form with $\cos(\omega) = 1 - 2 \sin^2(\omega/2)$ and the substitution $\omega' = \omega/2$ as

$$\begin{aligned} \int_{-\pi}^{\pi} \frac{d\omega}{h(a, \omega) h(b, \omega)^2} &= \frac{2}{ab^2} \int_0^{\pi} \frac{d\omega}{\sqrt{1 + \frac{4(1-a)}{a^2} \sin^2(\omega/2)} \left(1 + \frac{4(1-b)}{b^2} \sin^2(\omega/2)\right)} \\ &= \frac{4}{ab^2} \int_0^{\pi/2} \frac{d\omega'}{\sqrt{1 + \frac{4(1-a)}{a^2} \sin^2(\omega')} \left(1 + \frac{4(1-b)}{b^2} \sin^2(\omega')\right)}. \end{aligned}$$

This can be written in terms of elliptic integrals using Property F.9 as

$$\begin{aligned} & \int_0^{\pi/2} \frac{d\omega'}{\sqrt{1 + \frac{4(1-a)}{a^2} \sin^2(\omega') \left(1 + \frac{4(1-b)}{b^2} \sin^2(\omega')\right)}} \\ &= \frac{a}{2-a} \left(\frac{b^2(1-a)}{b^2(1-a) - a^2(1-b)} \right) K(k) - \frac{a^3(1-b)}{(2-a)(b^2(1-a) - a^2(1-b))} \Pi(\alpha^2, k), \end{aligned} \quad (95)$$

with $k = \sqrt{1-a}/(1-a/2)$ and

$$\alpha^2 = \frac{b^2(1-a) - a^2(1-b)}{b^2(1-a/2)^2}.$$

Plugging in (94) and (95) into (93), we find that the $K(\cdot)$ term cancels out, completing the proof. \square

F.6 Other Helper Results

We list a number of other miscellaneous useful results.

Lemma F.17. For a sequence $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots) \in \ell^2$ and a constant $0 \leq c < 1$, we have

$$\sum_{t=0}^{\infty} \sum_{\tau=0}^{\infty} \beta_t \beta_{\tau} c^{|t-\tau|} \leq \left(\frac{1+c}{1-c} \right) \|\boldsymbol{\beta}\|_2^2.$$

Proof. We break the sum into powers of c and use the Cauchy-Schwarz inequality (*) to get

$$\begin{aligned} \sum_{t=0}^{\infty} \sum_{\tau=0}^{\infty} \beta_t \beta_{\tau} c^{|t-\tau|} &= \|\boldsymbol{\beta}\|_2^2 + 2 \sum_{k=1}^{\infty} c^k \left(\sum_{t=0}^{\infty} \beta_t \beta_{t+k} \right) \\ &\stackrel{(*)}{\leq} \|\boldsymbol{\beta}\|_2^2 + 2 \sum_{k=1}^{\infty} c^k \|\boldsymbol{\beta}\|_2^2. \end{aligned}$$

Summing up the geometric series with a multiplier $0 \leq c < 1$ completes the proof. \square

Lemma F.18. Consider a random vector \mathbf{x} that satisfies $\mathbb{E}[\mathbf{x}] = 0$, $\mathbb{E}[\mathbf{x} \otimes \mathbf{x}] = \mathbf{H} \succeq \mu \mathbf{I}$ for some $\mu > 0$ and $\mathbb{E}[\|\mathbf{x}\|_2^2 \mathbf{x} \otimes \mathbf{x}] \preceq R^2 \mathbf{H}$. Then, we have for all $\eta \leq 1/R^2$ and all PSD matrices \mathbf{M} that

$$\text{Tr}[(\mathbf{I} - \eta \mathbf{x} \otimes \mathbf{x}) \mathbf{M} (\mathbf{I} - \eta \mathbf{x} \otimes \mathbf{x})] \leq (1 - \eta \mu) \text{Tr}[\mathbf{M}].$$

Proof. The left side above (call it ‘‘LHS’’) is bounded by

$$\begin{aligned} \text{LHS} &= \text{Tr}[\mathbf{M}] - 2\eta \text{Tr}[\mathbf{M} \mathbf{x} \otimes \mathbf{x}] + \eta^2 \text{Tr} \left[\mathbb{E} \left[\|\mathbf{x}\|_2^2 \mathbf{x} \otimes \mathbf{x} \right] \mathbf{M} \right] \\ &\leq \text{Tr}[\mathbf{M}] - 2\eta \text{Tr}[\mathbf{H} \mathbf{M}] + \eta^2 R^2 \text{Tr}[\mathbf{H} \mathbf{M}] \\ &\leq \text{Tr}[\mathbf{M}] - \eta \text{Tr}[\mathbf{H} \mathbf{M}] \\ &\leq (1 - \eta \mu) \text{Tr}[\mathbf{M}], \end{aligned}$$

where we used (a) $\mathbb{E}[\|\mathbf{x}\|_2^2 \mathbf{x} \otimes \mathbf{x}] \preceq R^2 \mathbf{H}$, (b) $\eta \leq 1/R^2$, and (c) $\mathbf{H} \succeq \mu \mathbf{I}$. \square

Lemma F.19. For PSD matrices $\mathbf{0} \preceq \mathbf{A}_1, \dots, \mathbf{A}_k \preceq \mathbf{I}$ of shape $d \times d$, we have $|\text{Tr}[\mathbf{A}_1 \cdots \mathbf{A}_k]| \leq d$.

Proof. Recall the inner product $\langle \mathbf{A}, \mathbf{B} \rangle = \text{Tr}[\mathbf{A} \mathbf{B}^{\top}]$ on the space of real $d \times d$ matrices. Using Hölder’s inequality on the Schatten p -norms, we get

$$|\text{Tr}[\mathbf{A}_1 \cdots \mathbf{A}_k]| = |\langle \mathbf{A}_1, \mathbf{A}_k \cdots \mathbf{A}_2 \rangle| \leq \|\mathbf{A}_1\|_{S_1} \|\mathbf{A}_k \cdots \mathbf{A}_2\|_{S_{\infty}}.$$

DP-FTRL Variant	Citation	Corr. matrix B	Anytime?	Computation Cost	
				Generation	Training (per step)
DP-SGD	[2]	Identity	✓	$O(1)$	$O(1)$
Honaker/TreeAgg	[33]	Lower-Triangular (LT)	✓	$O(1)$	$O(\log T)$
Optimal CC	[21]	Toeplitz & LT	✓	$O(1)$	$O(T)$
ν -DP-FTRL	Ours	Toeplitz & LT	✓	$O(1)$	$O(T)$
FFT	[17]	Toeplitz	-	$O(1)$	$O(T \log^2 T)$
Full Honaker	[25]	Arbitrary	-	$O(T^2)$	$O(T^2)$
Multi-Epoch (ME)	[17]	Arbitrary	-	$O(T^3)$	$O(T^2)$

Table 4: **Variants of DP-FTRL**: the noise correlation matrix B and whether the correlation matrix B can be created/optimized agnostic to the time horizon T (denoted as “Anytime”), and the computation cost.

Here, the Schatten 1-norm $\|\cdot\|_{S_1}$ is the ℓ_1 norm of the singular values (i.e. the nuclear norm); this is just the trace for a PSD matrix. Thus,

$$\|\mathbf{A}_1\|_{S_1} = \text{Tr}[\mathbf{A}_1] \leq \text{Tr}[\mathbf{I}] = 1.$$

The $\|\cdot\|_{S_\infty}$ is the ℓ_∞ norm of the singular values, i.e. the operator norm $\|\cdot\|_2$. We get,

$$\|\mathbf{A}_k \cdots \mathbf{A}_2\|_2 \leq \|\mathbf{A}_k\|_2 \cdots \|\mathbf{A}_2\|_2 \leq 1.$$

□

Lemma F.20. For some fixed integer $t \geq 1$ and constants $a > 0$, $\rho \in (0, 1)$, define the function

$$f(\tau) = \tau + \frac{1}{\rho a} \exp(-a\tau) \mathbb{1}(\tau < t - 1).$$

For $\hat{\tau} = \min\{t - 1, a^{-1} \log(1/\rho)\}$, we have,

$$f(\hat{\tau}) = \min\left\{t - 1, \frac{1}{a}(1 + \log(1/\rho))\right\} \leq \frac{1}{a}(1 + \log(1/\rho)).$$

Proof. The convex function $\tau \mapsto \tau + \frac{1}{\rho a} \exp(-a\tau)$ is minimized at $\tau_\star = a^{-1} \log(1/\rho) > 0$ with a minimum value of $a^{-1}(1 + \log(1/\rho))$. If $t - 1 \leq \hat{\tau}_\star$, we take $\hat{\tau} = t - 1$ and $f(\hat{\tau}) = t - 1 \leq \hat{\tau} \leq a^{-1}(1 + \log(1/\rho))$. □

The next lemma is from [45, Lemma 13].

Lemma F.21. Consider a function $\varphi : [0, \eta_{\max}] \rightarrow \mathbb{R}_+$ given by

$$\varphi(\eta) = A \exp(-\mu\eta T) + B\eta + C\eta^2 \log^2\left(\frac{1}{\eta\mu}\right),$$

given some constants $\eta_{\max}, \mu, A, B, C > 0$. If $T \geq (\mu\eta_{\max})^{-1}$, then we have

$$\varphi(\eta_\star) \leq A \exp(-\mu\eta_{\max}T) + \frac{3B}{\mu T} \left(1 \vee \log \frac{A\mu T}{B}\right) + \frac{3C}{\mu^2 T^2} \left(1 \vee \log \frac{A\mu^2 T^2}{C}\right)^2 \log^2(T),$$

for some $\eta_\star \leq \eta_{\max}$ depending on A, B, C, μ, T .

Lemma F.22. For $0 < c < 1/4$, we have,

$$0 < x \leq \frac{c}{9 \log^2(9/c)} \implies x \log^2(1/x) \leq c.$$

G Empirical Details

We train image-classification models using the CIFAR-10 dataset and language models using the Stack Overflow Next Word Prediction (SONWP) dataset available on `tensorflow-datasets`.

G.1 Image classification

Image classification has long been studied in DP ML. For example, the original DP-SGD work of Abadi et al. [2] focused on this task. We use CIFAR10 which has 50,000 training and 10,000 test examples. We evaluate and compute test accuracies on the entire test set, following the open-sourced code of Kairouz et al. [33]. We reuse the network architecture, dataset processing and initialization strategies presented in Kairouz et al. [33]; in particular, the architecture we use can be found in their Table 2 (b).

Setup and Tuning: We train all mechanisms for 2000 steps using a batch size of 500 and a clip norm of 1. This leads to a ML training dynamics of 20 epochs and 100 steps per epoch. We performed some initial small grid searches which showed nearly ubiquitously that momentum of 0.95 (searched over the grid 0, 0.85, 0.9, 0.95) and a linear learning rate cooldown $0.05 \times$ the initial learning rate over the last 500 steps of training improved model utility for all privacy levels. Thus, we fix these settings for all mechanisms except DP-SGD, for which no momentum performed best. For each mechanism, we then run a tuning grid search for the learning rate on coefficients in $\{1, 2, 5\}$ on powers in $[-2, 3]$, selecting the best mechanism for each privacy level from this interval. Final experiments are repeated 12 times in each setting and show 95% bootstrapped confidence intervals.

Some mechanisms include additional hyper parameters that specify the exact mechanism’s structure. For example, ME is specified by both the number of steps n and the max number of participations k . We include such parameters in the mechanism name. For all mechanisms, $n = 2000$.

G.2 Language modeling

Language modeling has been prominently studied in user-level DP contexts, starting with McMahan et al. [41]. DP training is important for real-world applications of language models trained on user data as these models can memorize their training data if appropriate mitigations are not applied [5, 13–15, 27, 37]. Indeed, DP already plays an important role in this application, as evidenced by Google’s use of DP for training on-device language models McMahan & Thakurta [40], Xu et al. [50]. StackOverflow Next Word Prediction (SONWP) contains over 10^8 examples non-iid over 342,477 users. The goal of this task is to predict the next word given a sequence of words. We use the same setup as Choquette-Choo et al. [17].

Setup and Tuning: All mechanisms use a clip norm of 1, a server momentum of 0.95, and a client learning rate of 1.0. They also use a server learning rate cooldown over the last 25% rounds. Initial tuning showed these to favorable settings. We train all mechanisms for 2052 steps and report the final evaluation accuracy of the model as reported on a held-out set of 10,000 examples. We zero out large updates of $\ell_\infty \geq 100$. We use the tuned server learning rates from Choquette-Choo et al. [17] for all existing mechanisms. For our new ν -DP-FTRL mechanisms, we do not perform extensive tuning due to computational costs and instead choose a near-optimal value of the server learning rate based on its correlation with normalized loss of the mechanism (i.e., normalized noise added due to B) and the results of Choquette-Choo et al. [16, Figure 11].