

On The Concurrence of Layer-wise Preconditioning Methods and Provable Feature Learning

Anonymous Authors¹

Abstract

Layer-wise preconditioning methods are a family of memory-efficient optimization algorithms that introduce preconditioners per axis of each layer’s weight tensors. These methods have seen a recent resurgence, demonstrating impressive performance relative to entry-wise (“diagonal”) preconditioning methods such as Adam(W) on a wide range of neural network optimization tasks. Complementary to their practical performance, we demonstrate that layer-wise preconditioning methods are provably necessary from a statistical perspective. To showcase this, we consider two prototypical models, *linear representation learning* and *single-index learning*, which are widely used to study how typical algorithms efficiently learn useful *features* to enable generalization. In these problems, we show SGD is a suboptimal feature learner when extending beyond ideal isotropic inputs $\mathbf{x} \sim N(\mathbf{0}, \mathbf{I})$ and well-conditioned settings typically assumed in prior work. We demonstrate theoretically and numerically that this suboptimality is fundamental, and that layer-wise preconditioning emerges naturally as the solution. We further show that standard tools like Adam preconditioning and batch-norm only mildly mitigate these issues, supporting the unique benefits of layer-wise preconditioning.

1. Introduction

Well-designed optimization algorithms have been an enabler to the staggering growth and success of machine learning. For the broader ML community, the Adam (Kingma & Ba, 2015) optimizer is likely the go-to scalable and performant choice for most tasks. However, despite its popularity in practice, it has been notoriously challenging to understand Adam-like optimizers theoretically, especially from a *statistical* (e.g. generalization) perspective.¹ In fact, there exist many theoretical settings where Adam and similar methods underperform in convergence or generalization relative to,

e.g., well-tuned SGD (see e.g. Wilson et al. (2017); Keskar & Socher (2017); Reddi et al. (2018); Gupta et al. (2021); Xie et al. (2022); Dereich et al. (2024)), further complicating a principled understanding the role of Adam-like optimizers in deep learning. Given these challenges, is there an alternative algorithmic paradigm that is comparable to the Adam family in practice that is also well-motivated from a statistical learning perspective? Encouragingly, in a recent large-scale deep learning optimization competition, AlgoPerf (MLCommons, 2024), Adam and its variants were outperformed in various “hold-out error per unit-compute”² metrics by a method known as Shampoo (Gupta et al., 2018), a member of a layer-wise “Kronecker-Factored” family of preconditioners, formally described in Section 2, contrasted with “diagonal” preconditioning methods like Adam.

Notable members of the Kronecker-Factored preconditioning family include Shampoo and KFAC (Martens & Grosse, 2015), as well as their many variants and descendants. These algorithms are motivated from an approximation-theoretic perspective, aiming to approximate some *ideal* curvature matrix (e.g. the Hessian or Fisher Information) in a way that mitigates the computational and memory challenges associated with second-order algorithms such as Newton’s Method (NM) or Natural Gradient Descent (NGD). However, towards establishing the benefit of these preconditioners, an approximation viewpoint is bottlenecked by our limited understanding of how the *idealized* second-order methods perform on neural-network learning tasks, even disregarding the computational considerations. It in fact remains unclear whether these second-order methods are inherently superior to the approximations designed to emulate them. For example, recent work has shown that, surprisingly, KFAC generally *outperforms* its ideal counterpart NGD in convergence rate and generalization error on typical deep learning tasks (Benzing, 2022). Thus, a key question remains:

*How do we explain the performance of
Kronecker-Factored preconditioned optimizers?*

In a seemingly distant area, the learning theory community has been interested in studying the solutions learned by abstractions of “typical” deep learning set-ups, where the overall goal is to theoretically demonstrate how neural networks

¹To be contrasted with an “optimization” perspective, e.g. guarantees of convergence to a critical point of the *training* objective.

²See Dahl et al. (2023, Section 4.2) for details.

learn features from data to perform better than classical “fixed features” methods, e.g. kernel machines. Much of this line of work focuses on analyzing the performance of SGD on simplified models of deep learning (see e.g. Collins et al. (2021); Damian et al. (2022); Ba et al. (2022); Barak et al. (2022); Abbe et al. (2023); Dandi et al. (2024a); Berthier et al. (2024); Nichani et al. (2024b); Collins et al. (2024)). Almost invariably, certain innocuous-looking assumptions are made, such as isotropic covariates $\mathbf{x} \sim N(\mathbf{0}, \mathbf{I})$. Under these conditions, SGD has been shown to exhibit desirable generalization properties. However, some works deviate from these assumptions in specific settings (Amari et al., 2020; Zhang et al., 2024b), and suggest that SGD can exhibit severely suboptimal generalization. Thus, toward extending our understanding of feature learning, it seems beneficial to consider a broader family of optimizers. This raises the following question:

What is a practical family of optimization algorithms that overcomes the deficiencies of SGD for standard feature learning tasks?

We answer the above two questions by focusing on two prototypical problems used to theoretically study feature learning: *linear representation learning* and *single-index learning*. In both problems, we show that SGD is clearly sub-optimal outside ideal settings, such as when the ubiquitous isotropic data $N(\mathbf{0}, \mathbf{I})$ assumption is violated. By inspecting the root cause behind these suboptimalities, we show that Kronecker-Factored preconditioners arise naturally as a *first-principles* solution to these issues. We provide novel non-approximation-theoretic motivations for this class of algorithms, while establishing new and improved learning-theoretic guarantees. We hope that this serves as strong evidence of an untapped synergy between deep learning optimization and feature learning theory.

Contributions.

- We study the linear representation learning problem under general anisotropic $\mathbf{x} \sim N(\mathbf{0}, \Sigma_{\mathbf{x}})$ covariates and show that the convergence of SGD can be drastically slow, even under mild anisotropy. Also, the convergence rate suffers an undesirable dependence on the “conditioning” of the instance even for ideal step-sizes. We arrive at a variant KFAC as the natural solution to these deficiencies of SGD, giving rise to the first *condition-number-free convergence rate* for the problem (Section 3.1).
- Next, we consider the problem of learning a single-index model using a two-layer neural network in the high-dimensional proportional limit. We show that for anisotropic covariates $\mathbf{x} \sim N(\mathbf{0}, \Sigma_{\mathbf{x}})$, SGD fails to learn useful features, whereas it is known that it learns suitable features in the isotropic setting. Furthermore, we show that KFAC is a natural fix to SGD, greatly enhancing the learned features in anisotropic settings (Section 3.2).
- Lastly, we carefully numerically verify our theoretical predictions. Notably, we confirm the findings in Benzing (2022) that full second-order methods heavily underperform KFAC in convergence rate and stability. We also show standard tools like Adam-like preconditioning and batch-norm (Ioffe & Szegedy, 2015) do not fix the issues we identify, even for our simple models, and may even *hurt* generalization in the latter’s case.

In addition to the works discussed earlier, we provide extensive related work and background in Appendix A.

Notation. We denote vector quantities by **bold** lower-case, and matrix quantities by **bold** upper-case. We use \odot to denote element-wise (Hadamard) product, \otimes for Kronecker product, and $\text{vec}(\cdot)$ the *column-major* vectorization operator. Positive (semi-)definite matrices are denoted by $\mathbf{Q} \succ (\succeq) \mathbf{0}$, and the corresponding partial order $\mathbf{P} \preceq \mathbf{Q} \implies \mathbf{Q} - \mathbf{P} \succeq \mathbf{0}$. We use $\|\cdot\|_{\text{op}}$, $\|\cdot\|_F$ to denote the operator (spectral) and Frobenius norms. We use $E[f(\mathbf{x})]$ to denote the expectation of $f(\mathbf{x})$, and $P[A(\mathbf{x})]$ to denote the probability of event $A(\mathbf{x})$. Given a batch $\{\mathbf{x}_i\}_{i=1}^n$, we denote the *empirical* expectation $\widehat{E}[f(\mathbf{x})] = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i)$. Given an indexed set of vectors, we use the upper case to denote the (row-wise) stacked matrix, e.g. $\mathbf{X} \triangleq [\mathbf{x}_1 \ \dots \ \mathbf{x}_n]^T \in \mathbb{R}^{n \times d_x}$. We reserve Σ ($\widehat{\Sigma}$) for (sample) covariance matrices, e.g. $\Sigma_{\mathbf{x}} = E[\mathbf{x}\mathbf{x}^T]$, $\widehat{\Sigma}_{\mathbf{x}} = \widehat{E}[\mathbf{x}\mathbf{x}^T] = \frac{1}{n} \mathbf{X}^T \mathbf{X}$. We use $\lesssim, \gtrsim, \approx$ to omit universal numerical constants, and standard asymptotic notation $o(\cdot), \mathcal{O}(\cdot), \Omega(\cdot), \Theta(\cdot)$. Lastly, we use the index shorthand $[n] = \{1, \dots, n\}$, and subscript $+$ to denote the “next iterate”, e.g. $\mathbf{G}_+ = \text{Next}(\mathbf{G})$.

2. Kronecker-Factored Approximation

One of the longest-standing research efforts in optimization literature is dedicated to understanding the role of (local) curvature toward accelerating convergence rates of optimization methods. An example is Newton’s method, where the curvature matrix (Hessian) serves as a preconditioner of the gradient, enabling one-shot convergence in quadratic optimization, in which gradient descent enjoys at best a linear convergence rate dictated by the conditioning of the problem. However, for high-dimensional variables, computing and storing the full curvature matrix is often infeasible. Thus enter Quasi-Newton and (preconditioned) Conjugate Gradient methods, where the goal is to reap the benefits of curvature under computational or structural specifications, such as {block-diagonal, low-rank, sparsity, etc.} constraints (e.g. BFGS family (Goldfarb, 1970; Liu & Nocedal, 1989; Nocedal & Wright, 1999)), or accessing the curvature matrix only through matrix-vector products (see e.g. Pearlmutter (1994); Schraudolph (2002); Martens (2010)).

Nevertheless, the use of these methods for neural network optimization introduces new considerations. Consider an L -layer fully-connected neural network (omitting

biases) $\mathbf{f}_\theta(\mathbf{x}) \triangleq \mathbf{W}_L \sigma(\mathbf{W}_{L-1} \cdots \sigma(\mathbf{W}_1 \mathbf{x}) \cdots)$, where $\mathbf{W}_\ell \in \mathbb{R}^{d \times d}$, $\ell \in [L]$ and $\theta \in \mathbb{R}^{Ld^2}$ is the concatenation of $\theta_\ell \triangleq \text{vec}(\mathbf{W}_\ell)$, $\ell \in [L]$. Firstly, establishing convergence of SGD (Arora et al., 2019), NGD (Zhang et al., 2019), or Gauss-Newton (Cai et al., 2019) (or their corresponding gradient flows) to global minima of the *training* objective is non-trivial, as optimization over θ is non-convex. Moreover, these results do not directly characterize the structure of the resulting features learned by the algorithms. Secondly, on the practical front, full preconditioners on θ require memory $\mathcal{O}(L^2 d^4)$, which grows prohibitively with depth and width. Block-diagonal approximations (where one curvature block $\mathbf{M}_\ell \in \mathbb{R}^{d^2 \times d^2}$ corresponds to a layer θ_ℓ) still require $\mathcal{O}(L d^4)$. Thus, entry-wise preconditioning as in Adam, with footprint $\mathcal{O}(L d^2) \approx \dim(\theta)$, is usually considered the only scalable class of preconditioners.

However, a distinct notion of ‘‘Kronecker-Factored’’ preconditioning emerged approximately concurrently with Adam, with representative examples such as KFAC and Shampoo. As its name suggests, since full block-diagonal approximations are too expensive, a Kronecker-Factored approximation is made instead, where $\mathbf{M}_\ell^{-1} \nabla_{\theta_\ell} \mathcal{L}(\theta) = (\mathbf{Q}_\ell \otimes \mathbf{P}_\ell)^{-1} \nabla_{\theta_\ell} \mathcal{L}(\theta)$, $\mathbf{P}_\ell, \mathbf{Q}_\ell \succeq \mathbf{0}$. Using properties of the Kronecker product (see Lemma E.1), this has the convenient interpretation of pre- and post-multiplying the weights in their *matrix form*:

$$(\mathbf{Q}_\ell \otimes \mathbf{P}_\ell)^{-1} \nabla_{\theta_\ell} \mathcal{L}(\theta) \iff \mathbf{P}_\ell^{-1} \nabla_{\mathbf{W}_\ell} \mathcal{L}(\theta) \mathbf{Q}_\ell^{-1}. \quad (1)$$

As such, the memory requirement of Kronecker-Factored layer-wise preconditioning is $\mathcal{O}(L d^2)$, matching that of entry-wise preconditioning. The notion of curvature differs from case to case, e.g., for KFAC, this is the Fisher Information matrix corresponding to the distribution parameterized by $f_\theta(\mathbf{x})$, whereas for Shampoo this is the full-matrix Adagrad preconditioner, in turn closely related to the Gauss-Newton matrix.³ We provide some sample derivations and background in Appendix D. However, as aforementioned, an approximation viewpoint falls short of explaining the practical performance of Kronecker-Factored methods, as they typically converge *faster* than their corresponding second-order method (Benzing, 2022) on deep learning tasks. This motivates understanding the unique benefits of layer-wise preconditioning methods from first principles, which brings us to the following section.

3. Feature Learning via Kronecker-Factored Preconditioning

We present two prototypical models of feature learning, *linear representation learning* and *single-index learning*, and demonstrate how typical guarantees for the features learned by SGD break down outside of idealized settings. We then

show how to rectify these issues by deriving a modified algorithm from first principles, and demonstrate that both cases in fact coincide with a particular Kronecker-factored preconditioning method. We now set-up the model architecture and algorithm primitive considered in both problems. We consider two-layer feedforward neural network predictors:

$$f_{\mathbf{F}, \mathbf{G}}(\mathbf{x}) = \mathbf{F} \sigma(\mathbf{G} \mathbf{x}), \quad (2)$$

where $\mathbf{F} \in \mathbb{R}^{d_y \times d_h}$, $\mathbf{G} \in \mathbb{R}^{d_h \times d_x}$ denote the weight matrices and $\sigma(\cdot)$ is a predetermined activation function. For scalar outputs $d_y = 1$, we use $f_{\mathbf{f}, \mathbf{G}}(\mathbf{x}) = \mathbf{f}^\top \sigma(\mathbf{G} \mathbf{x})$. For our purposes, we omit the bias vectors from both layers. We further denote the intermediate covariate pre- and post-activation $\mathbf{h} \triangleq \mathbf{G} \mathbf{x}$, $\mathbf{z} \triangleq \sigma(\mathbf{G} \mathbf{x})$. We consider a standard mean-squared-error (MSE) regression objective and its (batch) empirical counterpart:

$$\begin{aligned} \mathcal{L}(\mathbf{F}, \mathbf{G}) &\triangleq \frac{1}{2} \mathbb{E}_{(\mathbf{x}, \mathbf{y})} [\|\mathbf{y} - f_{\mathbf{F}, \mathbf{G}}(\mathbf{x})\|^2] \\ \widehat{\mathcal{L}}(\mathbf{F}, \mathbf{G}) &\triangleq \frac{1}{2} \widehat{\mathbb{E}} [\|\mathbf{y} - f_{\mathbf{F}, \mathbf{G}}(\mathbf{x})\|^2] \end{aligned} \quad (3)$$

Given a batch of inputs $\{\mathbf{x}_i\}_{i=1}^n$, we define the left and right preconditioners of the two layers (recall equation (1)):

$$\begin{aligned} \mathbf{Q}_G &= \widehat{\Sigma}_{\mathbf{x}} = \widehat{\mathbb{E}}[\mathbf{x} \mathbf{x}^\top], \quad \mathbf{Q}_F = \widehat{\Sigma}_{\mathbf{z}} = \widehat{\mathbb{E}}[\mathbf{z} \mathbf{z}^\top], \\ \mathbf{P}_G &= \widehat{\mathbb{E}} \left[\left(\frac{\partial f_{\mathbf{F}, \mathbf{G}}}{\partial \mathbf{h}} \right)^\top \frac{\partial f_{\mathbf{F}, \mathbf{G}}}{\partial \mathbf{h}} \right] \text{(or } \mathbf{I}_{d_h} \text{)}, \quad \mathbf{P}_F = \mathbf{I}_{d_y}. \end{aligned} \quad (4)$$

We introduce the flexibility of $\mathbf{P}_G = \mathbf{I}_{d_h}$ for when \mathbf{P}_G does not play a significant role; notably, this recovers certain Kronecker-Factored preconditioners that avoid extra backwards passes (see Appendix D). We consider a stylized alternating descent primitive, where we iteratively perform

$$\begin{aligned} \mathbf{G}_+ &= \mathbf{G} - \eta_G \mathbf{P}_G^{-1} \nabla_{\mathbf{G}} \widehat{\mathcal{L}}(\mathbf{F}, \mathbf{G}) (\mathbf{Q}_G + \lambda_G \mathbf{I}_{d_X})^{-1} \\ \mathbf{F}_+ &= \mathbf{F} - \eta_F \mathbf{P}_F^{-1} \nabla_{\mathbf{F}} \widehat{\mathcal{L}}(\mathbf{F}, \mathbf{G}_+) \mathbf{Q}_F^{-1}, \end{aligned} \quad (5)$$

where $\eta_G, \eta_F > 0$ are layer-wise learning rates, and $\lambda_G \in \mathbb{R}$ is a regularization parameter. In line with most prior work, we consider an alternating scheme for analytical convenience. We also assume that $\mathbf{G}_+, \mathbf{F}_+$ are computed on *independent batches* of data, equivalent to sample-splitting strategies in prior analysis (Collins et al. (2021); Zhang et al. (2024b); Ba et al. (2022); Moniri et al. (2024) etc).

The preconditioners (4) and update (5) bear a striking resemblance to KFAC (cf. Appendix D). In fact, the preconditioners align exactly with KFAC if we view $\mathcal{L}(\mathbf{F}, \mathbf{G})$ as a negative log-likelihood of a conditionally Gaussian model with fixed variance: $\hat{\mathbf{y}}(\mathbf{x}) \sim N(f_{\mathbf{F}, \mathbf{G}}(\mathbf{x}), \sigma^2 \mathbf{I})$. This is in some sense a coincidence (and a testament to the prevalence of KFAC’s design): rather than deriving the above preconditioners via approximating the Fisher Information matrix, we will show shortly how they arise as a natural

³This is itself a positive-definite approximation of the Hessian.

adjustment to SGD in our featured problems. We note that Kronecker-Factored preconditioning methods often involve further moving parts such as damping exponents $\mathbf{P}^{-\rho}$, additional ridge parameters on various preconditioners, and momentum. Though of great importance in practice, they are beyond the scope of this paper,⁴ and we only feature parameters that play a role in our analysis.

A convenient observation that unifies various stylized algorithms on two-layer networks is that the \mathbf{F} -update in (5) can be interpreted as an exponential moving average (EMA) over least-squares estimators conditioned on \mathbf{z} .

Lemma 3.1. *Given \mathbf{G} , define the least-squares estimator:*

$$\hat{\mathbf{F}}_{\text{ls}} \triangleq \underset{\hat{\mathbf{F}}}{\operatorname{argmin}} \frac{1}{2} \hat{\mathbb{E}} \left[\|\mathbf{y} - \hat{\mathbf{F}} \underbrace{\sigma(\mathbf{G}\mathbf{x})}_{\mathbf{z}}\|^2 \right] = \mathbf{Y}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1}$$

Given $\eta_F \in (0, 1]$, then the \mathbf{F} -update in (5) can be re-written as an EMA of $\hat{\mathbf{F}}_{\text{ls}}$; i.e., $\mathbf{F}_+ = (1 - \eta_F)\mathbf{F} + \eta_F \hat{\mathbf{F}}_{\text{ls}}$.

In particular, many prior works (e.g. Collins et al. (2021); Nayer & Vaswani (2022); Thekumparampil et al. (2021); Zhang et al. (2024b)) consider an alternating “minimization-descent” approach, where out of analytical convenience \mathbf{F} is updated by performing least-squares regression holding the hidden layer \mathbf{G} fixed. In light of Lemma 3.1, this corresponds to the case where $\eta_F = 1$.

3.1. Linear Representation Learning

Assume we have data generated by the following process

$$\mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Sigma_x), \quad \mathbf{y}_i = \mathbf{F}_* \mathbf{G}_* \mathbf{x}_i + \varepsilon_i, \quad (6)$$

where Σ_x is the input covariance, $\mathbf{F}_* \in \mathbb{R}^{d_y \times k}$, $\mathbf{G}_* \in \mathbb{R}^{k \times d_x}$ are (unknown) rank- k matrices, and $\varepsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Sigma_\varepsilon)$ is additive label noise independent of all other randomness. We consider Gaussian data throughout the paper for conciseness; all results in this section can be extended to subgaussian \mathbf{x}, ε via standard tools, affecting only the universal constants (see Appendix E). Let us define $\sigma_\varepsilon^2 \triangleq \lambda_{\max}(\Sigma_\varepsilon)$. Accordingly, our predictor model is a two-layer linear net (2) with $\mathbf{F} \in \mathbb{R}^{d_y \times k}$, $\mathbf{G} \in \mathbb{R}^{k \times d_x}$.

The goal of linear representation learning is to learn the low-dimensional feature space that \mathbf{G}_* maps to, which is equivalent to determining its row-space $\text{rowsp}(\mathbf{G}_*)$. Recovering $\mathbf{F}_*, \mathbf{G}_*$ is an ill-posed problem, as for any invertible $\mathbf{L} \in \mathbb{R}^{k \times k}$, the matrices $\mathbf{F}_* \mathbf{L}, \mathbf{L}^{-1} \mathbf{G}_*$ remain optimal. Therefore, we measure recovery of $\text{rowsp}(\mathbf{G}_*)$ via a *subspace distance*.

Definition 3.2 (Subspace Distance (Stewart & Sun, 1990)). Let $\mathbf{G}, \mathbf{G}_* \in \mathbb{R}^{k \times d_x}$ be matrices whose rows are orthonormal. Let $\mathcal{P}_*^\perp \in \mathbb{R}^{d_x \times d_x}$ be the projection matrix onto

⁴We refer the interested reader to Ishikawa & Karakida (2023) for discussion of these settings in the “maximal-update parameterization” framework (Yang & Hu, 2021b).

$\text{rowsp}(\mathbf{G}_*)^\perp$. Define the distance between the subspaces spanned by the rows of \mathbf{G} and \mathbf{G}_* by

$$\text{dist}(\mathbf{G}, \mathbf{G}_*) \triangleq \|\mathbf{G} \mathcal{P}_*^\perp\|_{\text{op}} \quad (7)$$

The subspace distance quantitatively captures the alignment between two subspaces, ranging between 0 (occurring iff $\text{rowsp}(\mathbf{G}_*) = \text{rowsp}(\mathbf{G})$) and 1 (occurring iff $\text{rowsp}(\mathbf{G}_*) \perp \text{rowsp}(\mathbf{G})$). We further make the following non-degeneracy assumptions.

Assumption 3.3. We assume \mathbf{G}_* is row-orthonormal, and \mathbf{F}_* is full-rank, $\text{rank}(\mathbf{F}_*) = k \leq d_y$. This is without loss of generality: if $k > d_y$, then recovering a k -dimensional row-space from \mathbf{y}_i is underdetermined. If $\text{rank}(\mathbf{F}_*) = k' < k$, then it suffices to consider $\mathbf{G}_* \in \mathbb{R}^{k' \times d_x}$.

The linear representation learning problem has often been studied in the context of multi-task learning (Du et al., 2021; Tripuraneni et al., 2020; Collins et al., 2021; Thekumparampil et al., 2021; Zhang et al., 2024b).

Remark 3.4 (Multi-task Learning). Multi-task learning considers data generated as $\mathbf{y}_i^{(t)} = \mathbf{F}_*^{(t)} \mathbf{G}_* \mathbf{x}_i^{(t)} + \varepsilon_i^{(t)}$ for distinct tasks $t = 1, \dots, T$, with the same goal of recovering the *shared* representation \mathbf{G}_* . Our algorithm and guarantees naturally extend here, see Appendix B.3 for full details. In particular, by embedding $\mathbf{F}_* = [\mathbf{F}_*^{(1)\top} \dots \mathbf{F}_*^{(T)\top}]^\top$, Assumption 3.3 is equivalent to the “task-diversity” conditions in the above works: $\text{rank}(\mathbf{F}_*) = \text{rank} \left(\sum_{t=1}^T \mathbf{F}_*^{(t)\top} \mathbf{F}_*^{(t)} \right) = k$.

We maintain the “single-task” setting in this section for concise bookkeeping while preserving the essential features of the representation learning problem. Various algorithms have been proposed toward provably recovering the representation \mathbf{G}_* . A prominent example is an alternating minimization-SGD scheme (Collins et al., 2021; Vaswani, 2024). In the cited works, a local convergence result⁵ is established for isotropic data $\Sigma_x = \mathbf{I}_{d_x}$. In Zhang et al. (2024b), it is shown that using SGD can drastically slow convergence even under mild anisotropy; their proposed algorithmic adjustment equates to applying the right-preconditioner $\mathbf{Q}_G = \hat{\Sigma}_x$. However, their local convergence result suffers a dependence on the condition number of \mathbf{F}_* , slowing the linear convergence rate for ill-conditioned \mathbf{F}_* . Let us now specify the algorithm template used in this section, that also encompasses the above work:

$$\bar{\mathbf{G}}_+ \text{ via (5), } \mathbf{G}_+ = \text{Ortho}(\bar{\mathbf{G}}_+), \quad \mathbf{F}_+ \text{ via (5).} \quad (8)$$

Notably, we row-orthonormalize the representation after each update. Besides ease of analysis, we have observed this numerically mitigates the elements of \mathbf{G} from blow-up

⁵By local convergence here we mean $\text{dist}(\mathbf{G}, \mathbf{G}_*)$ is sufficiently (but non-vanishingly) small.

when running variants of SGD. The alternating min-SGD algorithms in Collins et al. (2021); Vaswani (2024) are equivalent to iterating (8) setting $\mathbf{P}_G = \mathbf{Q}_G = \mathbf{I}$, $\eta_F = 1$ in (4), whereas Zhang et al. (2024b) use $\mathbf{P}_G = \mathbf{I}$, $\mathbf{Q}_G = \widehat{\Sigma}_x$, $\eta_F = 1$. Let us now write out the full-batch gradient update.

Full-Batch SGD. Given a fresh batch of data $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$, and current weights (\mathbf{F}, \mathbf{G}) , we have the representation gradient and corresponding SGD step:

$$\begin{aligned}\nabla_G \widehat{\mathcal{L}}(\mathbf{F}, \mathbf{G}) &= \frac{1}{n} \mathbf{F}^\top (\mathbf{F} \mathbf{G} \mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{X}) \\ \mathbf{G}_+ &= \mathbf{G} - \eta_G \nabla_G \widehat{\mathcal{L}}(\mathbf{F}, \mathbf{G}).\end{aligned}\quad (9)$$

When \mathbf{x} is isotropic $\Sigma_x = \mathbf{I}_{d_x}$, the key observation is that by multiplying both sides of (9) by \mathcal{P}_*^\perp , recalling $\mathbf{Y}^\top = \mathbf{F}_* \mathbf{G}_* \mathbf{X}^\top + \mathcal{E}^\top$, we have

$$\begin{aligned}\overline{\mathbf{G}}_+ \mathcal{P}_*^\perp &= \left(\mathbf{G} - \eta_G \mathbf{F}^\top \left((\mathbf{F} \mathbf{G} - \mathbf{F}_* \mathbf{G}_*) \widehat{\Sigma}_x - \frac{1}{n} \mathcal{E}^\top \mathbf{X} \right) \right) \mathcal{P}_*^\perp \\ &\approx (\mathbf{I}_k - \eta_G \mathbf{F}^\top \mathbf{F}) \mathbf{G} \mathcal{P}_*^\perp + \frac{\eta_G}{n} \mathbf{F}^\top \mathcal{E}^\top \mathbf{X} \mathcal{P}_*^\perp,\end{aligned}$$

where the approximate equality hinges on covariance concentration $\widehat{\Sigma}_x \approx \mathbf{I}_{d_x}$ and $\mathbf{G}_* \mathcal{P}_*^\perp = \mathbf{0}$. Therefore, in the isotropic setting, for sufficiently large $n \gtrsim d_x$, and appropriately chosen $\eta_G \approx \frac{1}{\lambda_{\max}(\mathbf{F}_*^\top \mathbf{F}_*)}$, then (omitting many details) we have the one-step contraction (Collins et al., 2021; Vaswani, 2024): defining $\Delta = \mathcal{O}\left(\sigma_\epsilon \sqrt{d_x/n}\right)$,

$$\text{dist}(\mathbf{G}_+, \mathbf{G}_*) \lesssim \left(1 - \frac{\lambda_{\min}(\mathbf{F}_*^\top \mathbf{F}_*)}{\lambda_{\max}(\mathbf{F}_*^\top \mathbf{F}_*)}\right) \text{dist}(\mathbf{G}, \mathbf{G}_*) + \Delta \quad (10)$$

Therefore, in low-noise/large-batch settings, this demonstrates SGD on the representation \mathbf{G} converges geometrically to \mathbf{G}_* (in subspace distance). However, there are clear suboptimalities to SGD. Firstly, the above analysis critically relies on $\Sigma_x = \mathbf{I}_{d_x}$ such that $\mathbf{G}_* \widehat{\Sigma}_x \mathcal{P}_*^\perp \approx \mathbf{0}$. As aforementioned, this is demonstrated to be crucial in Zhang et al. (2024b) for $\text{dist}(\mathbf{G}, \mathbf{G}_*)$ to converge using SGD. Secondly, the convergence of SGD is bottlenecked by the conditioning of \mathbf{F}_* . In fact, we show the dependence on \mathbf{F}_* in the contraction rate bound (10) cannot be improved in general, even under the most benign assumptions. Following Collins et al. (2021); Vaswani (2024), we define $\mathbf{G}_T = \text{SGD}(\mathbf{G}_0; \eta_G, T)$ as the output of alternating min-SGD, i.e. iterating (8) setting $\mathbf{P}_G = \mathbf{Q}_G = \mathbf{I}$, $\eta_F = 1$ in (4), for T steps with fixed step-size η_G starting from \mathbf{G}_0 .

Proposition 3.5. *Let $\Sigma_x = \mathbf{I}_{d_x}$, $n = \infty$. Choose any $d_x > k$, $d_y \geq k \geq 2$. Let the learner be given knowledge of \mathbf{F}_* , \mathbf{G}_* and $\text{dist}(\mathbf{G}_0, \mathbf{G}_*)$. However, assume the learner must fix $\eta_G > 0$ before observing \mathbf{G}_0 . Then, there exists $\mathbf{F}_* \in \mathbb{R}^{d_y \times k}$, $\mathbf{G}_*, \mathbf{G}_0 \in \mathbb{R}^{k \times d_x}$, such that $\mathbf{G}_T = \text{SGD}(\mathbf{G}_0; \eta_G, T)$ satisfies:*

$$\text{dist}(\mathbf{G}_T, \mathbf{G}_*) \geq \left(1 - 4 \frac{\lambda_{\min}(\mathbf{F}_*^\top \mathbf{F}_*)}{\lambda_{\max}(\mathbf{F}_*^\top \mathbf{F}_*)}\right)^T \text{dist}(\mathbf{G}_0, \mathbf{G}_*).$$

The proof can be found in Appendix B.1. Since we set $\Sigma_x = \mathbf{I}$, the lower bound also holds for the algorithm in Zhang et al. (2024b). We remark departing from a worst-case analysis to a generic performance lower bound, e.g. random initialization or varying step-sizes, is a nuanced topic even for the simple case of convex quadratics; see e.g. Bach (2024); Altschuler & Parrilo (2024). In light of Proposition 3.5 and (9), a sensible alteration might be to pre- and post-multiply $\nabla_G \widehat{\mathcal{L}}(\mathbf{F}, \mathbf{G})$ by $(\mathbf{F}^\top \mathbf{F})^{-1}$ and $\widehat{\Sigma}_x^{-1}$. These observations bring us to the proposed recipe in (5).

Stylized KFAC. By analyzing the shortcomings of the SGD update, we arrive at the proposed representation update:

$$\overline{\mathbf{G}}_+ = \mathbf{G} - \eta_G (\mathbf{F}^\top \mathbf{F})^{-1} \nabla_G \widehat{\mathcal{L}}(\mathbf{F}, \mathbf{G}) \widehat{\Sigma}_x^{-1}.$$

We can verify from (4) and (5) that $\mathbf{P}_G = \mathbf{F}^\top \mathbf{F}$ and $\mathbf{Q}_G = \widehat{\Sigma}_x$. Thus, we have recovered a stylized variant of KFAC as previewed. Our main result in this section is a local convergence guarantee.

Theorem 3.6. *Consider running (8) with $\lambda_G = 0$, $\eta_G \in [0, 1]$, and $\eta_F = 1$. Define $\bar{\sigma}^2 \triangleq \max\{1, \frac{\sigma_\epsilon^2}{\sigma_{\min}(\mathbf{F}_*)^2 \lambda_{\min}(\Sigma_x)}\}$. As long as $\text{dist}(\mathbf{G}, \mathbf{G}_*) \leq 0.01 \kappa^{-1}(\Sigma_x) \kappa^{-1}(\mathbf{F}_*)$ and $n \gtrsim \bar{\sigma}^2 (d_x + \log(1/\delta))$, we have:*

$$\text{dist}(\mathbf{G}_+, \mathbf{G}_*) \leq (1 - 0.9\eta_G) \text{dist}(\mathbf{G}, \mathbf{G}_*) + \eta_G \Delta,$$

with probability $\geq 1 - \delta$, where $\Delta \triangleq \mathcal{O}\left(\bar{\sigma} \sqrt{\frac{d_x + \log(1/\delta)}{n}}\right)$

Crucially, our contraction factor is condition-number-free, subverting the lower bound in Proposition 3.5 for sufficiently ill-conditioned \mathbf{F}_* . Therefore, setting η_G near 1 ensures a universal constant contraction rate. Curiously, our proposed stylized KFAC (8) aligns with an alternating “min-min” scheme (Jain et al., 2013; Thekumparampil et al., 2021), where \mathbf{F}, \mathbf{G} are alternately updated via solving the convex quadratic least-squares problem, by setting $\eta_F = \eta_G = 1$. However, our experiments (see Figure 5) demonstrate $\eta_G = 1$ is generally suboptimal, highlighting the flexibility of viewing KFAC as a descent method.

3.1.1. TRANSFER LEARNING

The upshot of representation learning is the ability to *transfer* (e.g. fine-tune) to a distinct, but related, task by only retraining \mathbf{F} (Du et al., 2021; Kumar et al., 2022). Assume we now have target data generated by:

$$\mathbf{x}_i^{(t)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Sigma_x^{(t)}), \quad \mathbf{y}_i^{(t)} = \mathbf{F}_*^{(t)} \mathbf{G}_* \mathbf{x}_i^{(t)} + \epsilon_i, \quad (11)$$

where $\epsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \Sigma_\epsilon)$, $\mathbf{F}_*^{(t)} \in \mathbb{R}^{d_y \times k}$. Notably, \mathbf{G}_* is shared with the “training” distribution (6). Given $\widehat{\mathbf{G}}$ (e.g. by running (8) on training task), we consider fitting the last layer \mathbf{F} given a batch of $n^{(t)}$ data from the target task (11).

275 **Lemma 3.7.** Let $\widehat{\mathbf{F}}_{\text{ls}}^{(t)} = \operatorname{argmin}_{\widehat{\mathbf{F}}} \widehat{\mathbb{E}}^{(t)}[\|\mathbf{y}^{(t)} - \widehat{\mathbf{F}}\mathbf{z}^{(t)}\|_2^2]$,
 276 $\mathbf{z}^{(t)} \triangleq \widehat{\mathbf{G}}\mathbf{x}^{(t)}$ be the optimal \mathbf{F} on the batch of $n^{(t)}$
 277 target data (11) given $\widehat{\mathbf{G}}$. Defining $\nu = \operatorname{dist}(\widehat{\mathbf{G}}, \mathbf{G}_*)$, given
 278 $n^{(t)} \gtrsim k + \log(1/\delta)$, we have with probability $\geq 1 - \delta$:

$$\begin{aligned} \mathcal{L}^{(t)}(\widehat{\mathbf{F}}_{\text{ls}}^{(t)}, \widehat{\mathbf{G}}) &\triangleq \mathbb{E} \left[\|\mathbf{y}^{(t)} - \mathbf{F}_*^{(t)} \mathbf{G}_* \mathbf{x}^{(t)}\|_2^2 \right] \\ &\lesssim \|\mathbf{F}_*^{(t)}\|_F^2 \lambda_{\max}(\Sigma_{\mathbf{x}}^{(t)}) \nu^2 + \frac{\sigma_\epsilon(d_y k + \log(1/\delta))}{n^{(t)}}. \end{aligned}$$

As hoped, the MSE of the fine-tuned predictor decomposes into a bias term scaling with the quality of $\widehat{\mathbf{G}}$, and a noise term scaling with $\dim(\mathbf{F})/n^{(t)}$. We comment the required data is $\approx k$ rather than $\approx d_x$ resulting from doing regression from scratch (Wainwright, 2019). Additionally, the noise term scales with $\dim(\mathbf{F}) = d_y k$ rather than $d_y d_x$ of the full predictor space. The transfer learning set-up (11) also reveals why data normalization (e.g. whitening, batch-norm (Ioffe & Szegedy, 2015)) can be counterproductive. To illustrate this, consider perfectly whitening the training covariates $\mathbf{v} = \Sigma_{\mathbf{x}}^{-1/2} \mathbf{x}$. By this change of variables, the ground-truth predictor changes $\mathbf{y} \approx \mathbf{F}_* \mathbf{G}_* \mathbf{x} = \mathbf{F}_* \mathbf{G}_* \Sigma_{\mathbf{x}}^{1/2} \mathbf{v}$. This is unproblematic so far—in fact, since the covariates \mathbf{v} are isotropic, SGD now may converge. However, instead of $\text{rowsp}(\mathbf{G}_*)$, the representation now converges to $\text{rowsp}(\mathbf{G}_* \Sigma_{\mathbf{x}}^{1/2})$. Deploying on the target task, since $\Sigma_{\mathbf{x}} \neq \Sigma_{\mathbf{x}}^{(t)}$, we have $\text{rowsp}(\widehat{\mathbf{G}}) \approx \text{rowsp}(\mathbf{G}_* \Sigma_{\mathbf{x}}^{1/2}) \neq \text{rowsp}(\mathbf{G}_* (\Sigma_{\mathbf{x}}^{(t)})^{1/2})$. In other words, in return for stabilizing optimization, normalizing the data destroys the shared structure of the predictor model! We illustrate this effect in Figure 3.

3.2. Single Index Learning

Assume that we observe n i.i.d. samples generated according to the following single-index model:

$$\mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Sigma_{\mathbf{x}}), \quad y_i = \sigma_*(\beta_*^\top \mathbf{x}_i) + \varepsilon_i \quad (12)$$

where $\Sigma_{\mathbf{x}} \in \mathbb{R}^{d_x \times d_x}$ is the input covariance, $\sigma_* : \mathbb{R} \rightarrow \mathbb{R}$ is the teacher activation function, $\beta_* \in \mathbb{R}^{d_x}$ is the (unknown) target direction, and ε_i is an additive noise $\varepsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_\varepsilon^2)$ independent of all other sources of randomness. We also make the following common assumption on β_* (Dicker (2016); Dobriban & Wager (2018); Tripuraneni et al. (2021a); Moniri et al. (2024); Moniri & Hassani (2024a), etc.) that ensures the covariates \mathbf{x}_i alone do not carry any information about the target direction.

Assumption 3.8. The vector β_* is drawn from $\beta_* \sim \mathcal{N}(0, d_x^{-1} \mathbf{I}_{d_x})$ independent of other sources of randomness.

In this section, we study the problem of fitting a two-layer feedforward neural network $f_{\mathbf{f}, \mathbf{G}}$ for prediction of unseen data points drawn independently from (12) at test time. When \mathbf{G} is kept at a random initialization and \mathbf{f} is trained

using ridge regression, the model coincides with a random features model (Rahimi & Recht, 2007; Montanari et al., 2019; Hu & Lu, 2023) and has repeatedly used as a toy model to study and explain various aspects of practical neural networks (see Lin & Dobriban (2021); Adlam & Pennington (2020); Tripuraneni et al. (2020); Hassani & Javanmard (2024); Bombari et al. (2023); Lee et al. (2023a); Bombari & Mondelli (2024a;b), etc.).

When the covariates are isotropic $\Sigma_{\mathbf{x}} = \mathbf{I}_{d_x}$, it is shown that a single step of full-batch SGD update on \mathbf{G} can drastically improve the performance of the model over random features as a result of *feature learning* by aligning the top right-singular-vector of the updated representation layer \mathbf{G} with the direction β_* (Damian et al., 2022; Ba et al., 2022; Moniri et al., 2024; Cui et al., 2024; Dandi et al., 2024a;b;c). In this section, we assume that the covariates are anisotropic and show that in this case, the one-step full batch SGD is suboptimal and can learn an ill-correlated direction even when the sample size n is large. We then demonstrate that the KFAC update with the preconditioners from (4) is in fact the natural fix to the full batch SGD.

Full-Batch SGD. Following the prior work, at initialization, we set $\mathbf{f} = d_h^{-1/2} \mathbf{f}_0$ with $\mathbf{f}_0 \sim \mathcal{N}(0, d_h^{-1} \mathbf{I}_{d_x})$, and $\mathbf{G} = \mathbf{G}_0$ with i.i.d. $\mathcal{N}(0, d_x^{-1})$ entries. We update \mathbf{G} with one step of full batch SGD with step size $\eta_{\mathbf{G}} = \eta \sqrt{d_h}$; i.e.,

$$\mathbf{G}_{\text{SGD}} \triangleq \mathbf{G}_0 - \eta \sqrt{d_h} \nabla_{\mathbf{G}} \widehat{\mathcal{L}}(\mathbf{f}_0, \mathbf{G}_0).$$

In the following theorem, we provide an approximation of the updated first layer \mathbf{G}_{SGD} , which is a generalization of (Ba et al., 2022, Proposition 2.1) for $\Sigma_{\mathbf{x}} \neq \mathbf{I}_{d_x}$.

Theorem 3.9. Assume that the activation function σ is $\mathcal{O}(1)$ -Lipschitz and that Assumption 3.8 holds. In the limit where n, d_x, d_h tend to infinity proportionally, the matrix \mathbf{G}_{SGD} , with probability $1 - o(1)$, satisfies

$$\|\mathbf{G}_0 + \alpha \eta \mathbf{f}_0 \beta_{\text{SGD}}^\top - \mathbf{G}_{\text{SGD}}\|_{\text{op}} \rightarrow 0,$$

in which $\alpha = \mathbb{E}_z[\sigma'(z)]$ with $z \sim \mathcal{N}(0, d_x^{-1} \operatorname{Tr}(\Sigma_{\mathbf{x}}))$, and the vector β_{SGD} is given by $\beta_{\text{SGD}} = n^{-1} \mathbf{X}^\top \mathbf{y}$.

This theorem shows that one step of full batch SGD update approximately adds a rank-one component $\alpha \eta \mathbf{f}_0 \beta_{\text{SGD}}^\top$ to the initialized weights \mathbf{G}_0 . Thus, the pre-activation features for a given input $\mathbf{x} \in \mathbb{R}^{d_x}$ after the update are given by

$$\mathbf{h}_{\text{SGD}} = \mathbf{G}_{\text{SGD}} \mathbf{x} \approx \mathbf{G}_0 \mathbf{x} + \alpha \eta (\beta_{\text{SGD}}^\top \mathbf{x}) \mathbf{f}_0 \in \mathbb{R}^{d_h}$$

where the first and second term correspond to the *random feature*, and the *learned feature* respectively. To better understand the learned feature component, note that defining $c_{*,1} = \mathbb{E}_{z \sim \mathcal{N}(0, d_x^{-1} \operatorname{Tr}(\Sigma_{\mathbf{x}}))}[\sigma'_*(z)]$, the target function $\sigma_*(\beta_*^\top \mathbf{x})$ can be decomposed as

$$\sigma_*(\beta_*^\top \mathbf{x}) = c_{*,1} \beta_{*,1}^\top \mathbf{x} + \sigma_{*,\perp}(\beta_*^\top \mathbf{x})$$

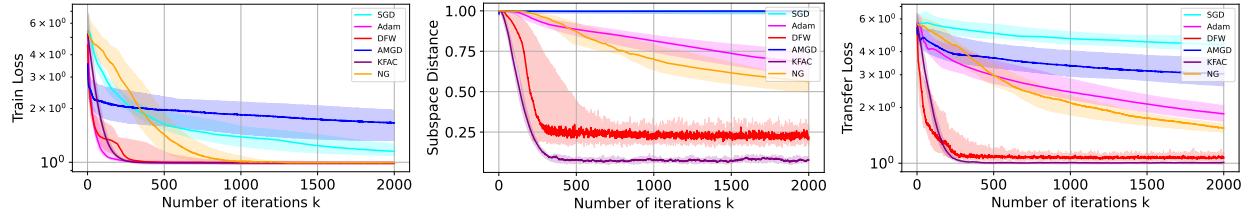


Figure 1. From left to right: the training loss, subspace distance, and transfer loss induced by various algorithms on a linear representation learning task. We note that various algorithms converge in training loss, but negligibly in subspace distance, and thus transfer loss.

satisfying $E_{\mathbf{x}} [c_{*,1}(\beta_{*}^\top \mathbf{x}) \sigma_{*,\perp}(\beta_{*}^\top \mathbf{x})] = 0$. Therefore, when $c_{*,1} \neq 0$, the target function has a *linear part*. Full batch SGD is estimating the direction of β_* using this linear part with the estimator $\beta_{\text{SGD}} = \mathbf{X}^\top \mathbf{y}/n$. However, the natural choice for this task is in fact ridge regression $\hat{\beta}_\lambda = (\hat{\Sigma}_{\mathbf{x}} + \lambda \mathbf{I}_{d_{\mathbf{x}}})^{-1} \mathbf{X}^\top \mathbf{y}/n$, and β_{SGD} is missing the prefactor $(\hat{\Sigma}_{\mathbf{x}} + \lambda \mathbf{I}_{d_{\mathbf{x}}})^{-1}$. In the isotropic case $\Sigma_{\mathbf{x}} = \mathbf{I}_{d_{\mathbf{x}}}$, we expect $\hat{\Sigma}_{\mathbf{x}} \approx \mathbf{I}_{d_{\mathbf{x}}}$ when $n \gg d_{\mathbf{x}}$. Thus, in this case the estimator β_{SGD} is roughly equivalent to the ridge estimator and can recover the direction β_* . However, in the anisotropic case, β_{SGD} is biased even when $n \gg d_{\mathbf{x}}$. To make these intuitions rigorous, we characterize in the following proposition the correlation between the learned direction β_{SGD} and the true direction β_* .

Lemma 3.10. Under the assumptions of Theorem 3.9, the correlation between β_* and β_{SGD} satisfies

$$\left| \frac{\beta_{*}^\top \beta_{\text{SGD}}}{\|\beta_{\text{SGD}}\|_2 \|\beta_*\|_2} - \frac{\frac{c_{*,1}}{d_{\mathbf{x}}} \text{Tr}(\Sigma_{\mathbf{x}})}{\sqrt{\frac{c_{*,1}^2 + \sigma_\varepsilon^2}{n} \text{Tr}(\Sigma_{\mathbf{x}}) + \frac{c_{*,1}^2}{d_{\mathbf{x}}} \text{Tr}(\Sigma_{\mathbf{x}}^2)}} \right| \rightarrow 0$$

with probability $1 - o(1)$, in which $c_{*,1} = E_z[\sigma'_*(z)]$ and $c_*^2 = E_z[\sigma_*^2(z)]$ with $z \sim N(0, d_{\mathbf{x}}^{-1} \text{Tr}(\Sigma_{\mathbf{x}}))$.

This lemma shows that the correlation is increasing in the strength of the linear component $c_{*,1}$ while keeping the signal strength c_* fixed. Also, based on this lemma, when $n \gg d_{\mathbf{x}}$, the correlation is given by $d_{\mathbf{x}}^{-1} \text{Tr}(\Sigma_{\mathbf{x}})/\sqrt{d_{\mathbf{x}}^{-1} \text{Tr}(\Sigma_{\mathbf{x}}^2)}$, which is equal to one if and only if $\Sigma_{\mathbf{x}} = \sigma^2 \mathbf{I}_{d_{\mathbf{x}}}$ for some $\sigma \in \mathbb{R}$. This means these are the only covariance matrices for which applying one step of full batch SGD update learns the correct direction of β_* .

Stylized KFAC. This time, we update \mathbf{G} using the stylized KFAC update from (5) with the regularized $\mathbf{P}_{\mathbf{G}}$. We use the same initialization as full-batch SGD. The updated representation layer in this case is given by

$$\mathbf{G}_{\text{KFAC}} \triangleq \mathbf{G}_0 - \eta \sqrt{d_{\mathbf{h}}} \nabla_{\mathbf{G}} \hat{\mathcal{L}}(\mathbf{f}_0, \mathbf{G}_0) (\mathbf{Q}_{\mathbf{G}} + \lambda_{\mathbf{G}} \mathbf{I}_{d_{\mathbf{x}}})^{-1}.$$

The preconditioning factor $(\mathbf{Q}_{\mathbf{G}} + \lambda_{\mathbf{G}} \mathbf{I}_{d_{\mathbf{x}}})^{-1}$ with $\mathbf{Q}_{\mathbf{G}} = \hat{\Sigma}_{\mathbf{x}}$ is precisely the factor required so that the direction learned by the one-step update to match the ridge regression estimator with ridge parameter $\lambda_{\mathbf{G}}$ as shown in the following immediate corollary of Theorem 3.9.

Corollary 3.11. Under the same set of assumptions as Theorem 3.9, the matrix \mathbf{G}_{KFAC} , satisfies

$$\|\mathbf{G}_0 + \alpha \eta \mathbf{f}_0 \beta_{\text{KFAC}}^\top - \mathbf{G}_{\text{KFAC}}\|_{\text{op}} \rightarrow 0$$

with probability $1 - o(1)$, where α is defined in Theorem 3.9, and $\beta_{\text{KFAC}} = (\mathbf{Q}_{\mathbf{G}} + \lambda_{\mathbf{G}} \mathbf{I}_{d_{\mathbf{x}}})^{-1} \mathbf{X}^\top \mathbf{y}/n = \hat{\beta}_{\lambda_{\mathbf{G}}}$.

Because β_{KFAC} is equivalent to ridge regression, we expect it to align well with β_* even for anisotropic $\Sigma_{\mathbf{x}}$, given a proper choice of $\lambda_{\mathbf{G}}$. The following lemma formally characterizes the correlation between β_{KFAC} and β_* for any $\lambda_{\mathbf{G}} \in \mathbb{R}$.

Lemma 3.12. Under the assumptions of Theorem 3.9, the correlation between β_* and β_{KFAC} satisfies

$$\left| \frac{\beta_{\text{KFAC}}^\top \beta_*}{\|\beta_{\text{KFAC}}\|_2 \|\beta_*\|_2} - \frac{c_{*,1} \Psi_1}{\sqrt{c_{*,1}^2 \Psi_2 + \frac{d_{\mathbf{x}}}{n} (c_{*,1}^2 + \sigma_\varepsilon^2) \Psi_3}} \right| \rightarrow 0$$

with probability $1 - o(1)$, where $c_{*,1}^2 = E_z[\sigma'_*(z)]$, $c_{*,1}^2 = E_z[\sigma_{*,\perp}^2(z)]$ with $z \sim N(0, d_{\mathbf{x}}^{-1} \text{Tr}(\Sigma_{\mathbf{x}}))$, and Ψ_1, Ψ_2, Ψ_3 are defined in (34) and depend on $\Sigma_{\mathbf{x}}$, $d_{\mathbf{x}}/n$, and $\lambda_{\mathbf{G}}$. In particular, as $\lambda_{\mathbf{G}} \rightarrow 0$ and $d_{\mathbf{x}}/n \rightarrow 0$, we have

$$\beta_{\text{KFAC}}^\top \beta_* / (\|\beta_{\text{KFAC}}\|_2 \|\beta_*\|_2) \rightarrow 1.$$

This lemma shows that when $n \gg d_{\mathbf{x}}$, and $\lambda_{\mathbf{G}} \rightarrow 0$, the one-step stylized KFAC update—unlike the one-step full-batch SGD—perfectly recovers the target direction β_* , fixing the issue with full batch SGD with anisotropic covariances.

Remark 3.13. It is well-known that, given features that align with β_* , applying least-squares on $\mathbf{Z} = \sigma(\mathbf{G}_{\text{KFAC}} \mathbf{X})$, which from Lemma 3.1 is equivalent to the KFAC \mathbf{f} -update with $\eta_{\mathbf{f}} = 1$, leverages the feature to obtain a solution with good generalization. See Appendix C.4 for more details.

4. Numerical Validation

4.1. Linear Representation Learning

We numerically study the behavior of different algorithms for a transfer learning setting (11), where the model is to be trained on data generated by $(\mathbf{F}_{*}^{\text{train}}, \mathbf{G}_{*})$, and the transfer task has data generated by $(\mathbf{F}_{*}^{\text{test}}, \mathbf{G}_{*})$, i.e. the embedding \mathbf{G}_{*} is shared, but the task heads $\mathbf{F}_{*}^{\text{train}}$ and $\mathbf{F}_{*}^{\text{test}}$ are different. The training and test covariates have anisotropic

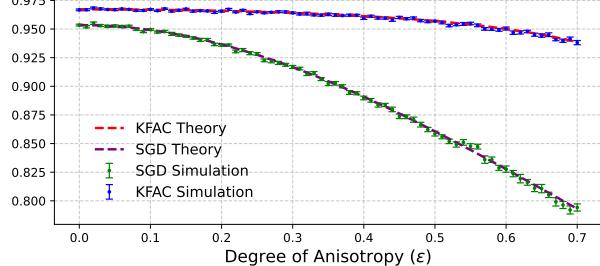
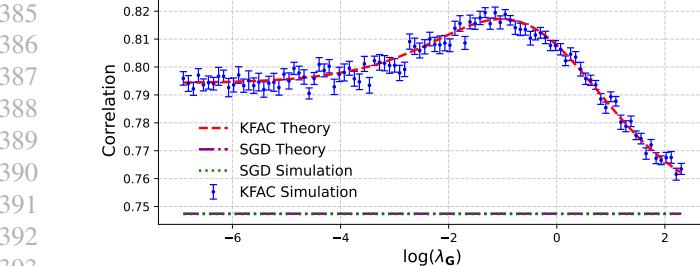


Figure 2. The correlation of the direction learned by SGD and KFAC with the true direction by numerical simulations averaged over 30 trials, and theoretical predictions. (Left) For different values of λ_G the theoretical predictions match the simulations very well. (Right) The alignment of the feature learned by SGD deteriorates as anisotropy is increased (larger ε), whereas the KFAC update remains accurate.

covariance matrices $\Sigma_{\mathbf{x}, \text{train}}$ and $\Sigma_{\mathbf{x}, \text{test}}$ respectively. Our data generation process for the training task and the transfer task are as follows:

$$\mathbf{y}_i^s = \mathbf{F}_*^s \mathbf{G}_* \mathbf{x}_i^s + \boldsymbol{\varepsilon}_i^s, \quad \mathbf{x}_i^s \stackrel{\text{i.i.d.}}{\sim} \Sigma_{\mathbf{x}, s}^{1/2} \text{Unif}(\{\pm 1\}^{d_x}), \\ \boldsymbol{\varepsilon}_i^s \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma_{\varepsilon, s}^2 \mathbf{I}_{d_y}), \quad s \in \{\text{test, train}\}, \quad (13)$$

where $\sigma_{\varepsilon, \text{train}} = 0.1$ and $\sigma_{\varepsilon, \text{test}} = 1$. We use $d_x = 100$, $d_y = 15$, $k = 8$, and batch size $n = 1024$. We present additional experiments and details in Appendix F, including discussions on the learning rates, and how \mathbf{F}_*^s , \mathbf{G}_*^s , $\Sigma_{\mathbf{x}, s}$ are precisely generated.

Head-to-head Evaluations. We track the training loss, subspace distance, and transfer loss of different algorithms during the update (Figure 1). Alongside SGD, KFAC, Adam, and NGD, we also consider Alternating Min-SGD (AMGD) (Collins et al., 2021; Vaswani, 2024), and De-bias & Feature-Whiten (DFW) (Zhang et al., 2024b) (corresponding to (5) with $\mathbf{P}_F = \mathbf{I}_{d_y}$), two algorithms studied in linear representation learning. The transfer loss is the loss incurred by fitting a least-squares $\hat{\mathbf{F}}_{\text{ls}}^{\text{test}}$ on the current \mathbf{G} iterate (see Lemma 3.7). Although various algorithms converge on $s = \text{train}$, KFAC outperforms all others in terms of subspace distance and transfer loss, as suggested by the theory.

Effect of Batch Normalization. We track the subspace distance and the training loss of AMGD (with and without batch-norm) and KFAC, see Figure 3. As theoretically predicted in Section 3.1.1, since batch-norm approximately whitens $\mathbf{x}_i^{\text{train}}$, AMGD+batch-norm converges in training loss. However, as predicted, it does not recover the correct representation, whereas KFAC does.

4.2. Single-Index Learning

Consider the single-index function learning setting of Section 3.2 with $\sigma_*(z) = z + \frac{1}{\sqrt{2}}(z^2 - 1)$, and $\sigma_\varepsilon = 1$.

Different Levels of Anisotropy. In this experiment, we set $d_x = 200$, $n = 6000$, and $d_h = 1000$ and set $\lambda_G \rightarrow 0$. For a parameter $\varepsilon \in \mathbb{R}$, we define $\Sigma_{\mathbf{x}} = \Sigma_{\mathbf{x}}^{(\varepsilon)}$ with

$$\Sigma_{\mathbf{x}}^{(\varepsilon)} = \text{diag}(\underbrace{1 + \varepsilon, \dots, 1 + \varepsilon}_{d_x/2}, \underbrace{1 - \varepsilon, \dots, 1 - \varepsilon}_{d_x/2}). \quad (14)$$

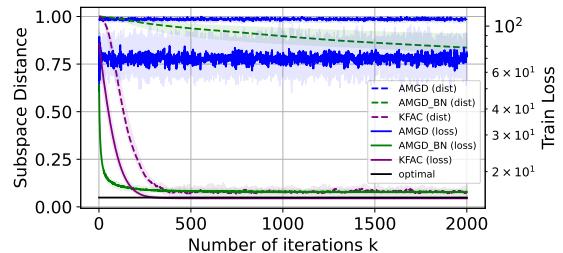


Figure 3. Subspace distance and the training loss of KFAC and AMGD (with and without batch-norm). Notably, batch-norm enables AMGD’s train loss to converge, but not its subspace distance.

For different values of ε , we simulate the KFAC and SGD updates numerically and compute their correlation with the true direction. We also theoretically predict the correlation using Lemma 3.10 and 3.12; see Figure 2 (Right). The SGD update fails to recover the true direction in highly anisotropic settings (large ε), whereas the one-step KFAC update remains accurate.

Theory vs. Simulations. We set $d_x = 900$, $n = 5000$, $d_h = 1000$, and $\Sigma_{\mathbf{x}} = \Sigma_{\mathbf{x}}^{(0.5)}$. For different λ_G , we simulate the correlation between the directions learned by KFAC and SGD with the true direction and compare it with predictions of Lemma 3.10 and 3.12; see Figure 2 (Left). We see that the theoretical results match very well with numerical simulations, even for moderately large n , d_x , and d_h . The direction learned by KFAC has a larger correlation with the true direction compared to that learned by SGD, as predicted.

5. Discussion

We study two models of feature learning in which we identify key issues of SGD-based feature learning approaches when departing from ideal settings. We then present Kronecker-Factored preconditioning—recovering variants of KFAC—to provably overcome these issues and derive improved guarantees. Our experiments on these simple models also confirm the suboptimality of full second-order methods, as well as the marginal benefit of Adam preconditioning and data normalization. We believe that analyzing properties of statistical learning problems can lead to fruitful insights into optimization and normalization schemes.

440 Impact Statement

441
442 The focus of this paper is on theoretical aspects of optimization and feature learning. We expect the results to be
443 illuminating for the optimization and learning theory com-
444 munity. We do not anticipate any negative societal impact.
445

446 References

447 Abbasi-Yadkori, Y. and Szepesvari, C. Regret bounds for
448 the adaptive control of linear quadratic systems. In *Con-
449 ference on Learning Theory*, 2011.

450 Abbe, E., Adsera, E. B., and Misiakiewicz, T. The merged-
451 staircase property: a necessary and nearly sufficient con-
452 dition for SGD learning of sparse functions on two-layer
453 neural networks. In *Conference on Learning Theory*,
454 2022.

455 Abbe, E., Adsera, E. B., and Misiakiewicz, T. SGD learning
456 on neural networks: leap complexity and saddle-to-saddle
457 dynamics. In *Conference on Learning Theory*, 2023.

458 Absil, P.-A., Mahony, R., and Sepulchre, R. *Optimization
459 algorithms on matrix manifolds*. Princeton University
460 Press, 2008.

461 Adlam, B. and Pennington, J. Understanding double de-
462 scent requires a fine-grained bias-variance decomposition.
463 In *Advances in Neural Information Processing Systems*,
464 2020.

465 Altschuler, J. and Parrilo, P. Acceleration by stepsize hedg-
466 ing: Multi-step descent and the silver stepsize schedule.
467 *Journal of the ACM*, 2024.

468 Amari, S.-i., Ba, J., Grosse, R., Li, X., Nitanda, A., Suzuki,
469 T., Wu, D., and Xu, J. When does preconditioning help
470 or hurt generalization? *arXiv preprint arXiv:2006.10732*,
471 2020.

472 Amid, E., Anil, R., and Warmuth, M. Locoprop: Enhancing
473 backprop via local loss optimization. In *International
474 Conference on Artificial Intelligence and Statistics*, 2022.

475 Anil, R., Gupta, V., Koren, T., Regan, K., and Singer, Y.
476 Scalable second order optimization for deep learning.
477 *arXiv preprint arXiv:2002.09018*, 2020.

478 Arnaboldi, L., Dandi, Y., Krzakala, F., Pesce, L., and
479 Stephan, L. Repetita iuvant: Data repetition allows SGD
480 to learn high-dimensional multi-index functions. *arXiv
481 preprint arXiv:2405.15459*, 2024.

482 Arora, S., Cohen, N., Hu, W., and Luo, Y. Implicit regu-
483 larization in deep matrix factorization. In *Advances in
484 Neural Information Processing Systems*, 2019.

485 Ba, J., Grosse, R. B., and Martens, J. Distributed second-
486 order optimization using Kronecker-factored approxima-
487 tions. In *International Conference on Learning Represen-
488 tations*, 2017.

489 Ba, J., Erdogdu, M. A., Suzuki, T., Wang, Z., Wu, D., and
490 Yang, G. High-dimensional asymptotics of feature learn-
491 ing: How one gradient step improves the representation.
492 In *Advances in Neural Information Processing Systems*,
493 2022.

494 Ba, J., Erdogdu, M. A., Suzuki, T., Wang, Z., and Wu, D.
495 Learning in the presence of low-dimensional structure: a
496 spiked random matrix perspective. In *Advances in Neural
497 Information Processing Systems*, 2024.

498 Bach, F. Scaling laws of optimization, 2024.
499 URL [https://francisbach.com/
500 scaling-laws-of-optimization/](https://francisbach.com/scaling-laws-of-optimization/).

501 Bai, Y. and Lee, J. D. Beyond linearization: On quadratic
502 and higher-order approximation of wide neural networks.
503 In *International Conference on Learning Representations*,
504 2020.

505 Barak, B., Edelman, B., Goel, S., Kakade, S., Malach, E.,
506 and Zhang, C. Hidden progress in deep learning: SGD
507 learns parities near the computational limit. In *Advances
508 in Neural Information Processing Systems*, 2022.

509 Ben Arous, G., Gheissari, R., and Jagannath, A. Online
510 stochastic gradient descent on non-convex losses from
511 high-dimensional inference. *Journal of Machine Learn-
512 ing Research*, 22(106):1–51, 2021.

513 Benzing, F. Gradient descent on neurons and its link to
514 approximate second-order optimization. In *International
515 Conference on Machine Learning*, 2022.

516 Bernstein, J. and Newhouse, L. Modular duality in deep
517 learning. *arXiv preprint arXiv:2410.21265*, 2024a.

518 Bernstein, J. and Newhouse, L. Old optimizer, new norm:
519 An anthology. *arXiv preprint arXiv:2409.20325*, 2024b.

520 Berthier, R., Montanari, A., and Zhou, K. Learning time-
521 scales in two-layers neural networks. *Foundations of
522 Computational Mathematics*, pp. 1–84, 2024.

523 Bollapragada, R., Nocedal, J., Mudigere, D., Shi, H.-J.,
524 and Tang, P. T. P. A progressive batching l-bfgs method
525 for machine learning. In *International Conference on
526 Machine Learning*, 2018.

527 Bombari, S. and Mondelli, M. How spurious features are
528 memorized: Precise analysis for random and NTK fea-
529 tures. In *International Conference on Machine Learning*,
530 2024a.

- 495 Bombari, S. and Mondelli, M. Privacy for free in the over-
496 parameterized regime. *arXiv preprint arXiv:2410.14787*,
497 2024b.
- 498 Bombari, S., Kiyani, S., and Mondelli, M. Beyond the
499 universal law of robustness: Sharper laws for random
500 features and neural tangent kernels. In *International
501 Conference on Machine Learning*, 2023.
- 502 Botev, A., Ritter, H., and Barber, D. Practical Gauss-Netwon
503 optimisation for deep learning. In *International Conference
504 on Machine Learning*, pp. 557–565, 2017.
- 505 Byrd, R. H., Hansen, S. L., Nocedal, J., and Singer, Y. A
506 stochastic quasi-Netwon method for large-scale optimization.
507 *SIAM Journal on Optimization*, 26(2):1008–1031,
508 2016.
- 509 Cai, T., Gao, R., Hou, J., Chen, S., Wang, D., He, D., Zhang,
510 Z., and Wang, L. Gram-Gauss-Netwon method: Learning
511 overparameterized neural networks for regression prob-
512 lems. *arXiv preprint arXiv:1905.11675*, 2019.
- 513 Collins, L., Hassani, H., Mokhtari, A., and Shakkottai, S.
514 Exploiting shared representations for personalized feder-
515 ated learning. In *International Conference on Machine
516 Learning*, 2021.
- 517 Collins, L., Hassani, H., Soltanolkotabi, M., Mokhtari, A.,
518 and Shakkottai, S. Provable multi-task representation
519 learning by two-layer ReLU neural networks. In *Inter-
520 national Conference on Machine Learning*, 2024.
- 521 Cui, H., Pesce, L., Dandi, Y., Krzakala, F., Lu, Y., Zde-
522 borova, L., and Loureiro, B. Asymptotics of feature
523 learning in two-layer networks after one gradient-step. In
524 *International Conference on Machine Learning*, 2024.
- 525 Dahl, G. E., Schneider, F., Nado, Z., Agarwal, N., Sastry,
526 C. S., Hennig, P., Medapati, S., Eschenhagen, R., Kasim-
527 beg, P., Suo, D., et al. Benchmarking neural network
528 training algorithms. *arXiv preprint arXiv:2306.07179*,
529 2023.
- 530 Damian, A., Lee, J., and Soltanolkotabi, M. Neural net-
531 works can learn representations with gradient descent. In
532 *Conference on Learning Theory*, 2022.
- 533 Dandi, Y., Krzakala, F., Loureiro, B., Pesce, L., and Stephan,
534 L. How two-layer neural networks learn, one (giant) step
535 at a time. *Journal of Machine Learning Research*, 25
536 (349):1–65, 2024a.
- 537 Dandi, Y., Pesce, L., Cui, H., Krzakala, F., Lu, Y. M., and
538 Loureiro, B. A random matrix theory perspective on
539 the spectrum of learned features and asymptotic gener-
540 alization capabilities. *arXiv preprint arXiv:2410.18938*,
541 2024b.
- 542 Dandi, Y., Troiani, E., Arnaboldi, L., Pesce, L., Zdeborova,
543 L., and Krzakala, F. The benefits of reusing batches for
544 gradient descent in two-layer networks: Breaking the
545 curse of information and leap exponents. In *International
546 Conference on Machine Learning*, 2024c.
- 547 Dangel, F., Harmeling, S., and Hennig, P. Modular block-
548 diagonal curvature approximations for feedforward ar-
549 chitectures. In *International Conference on Artificial
550 Intelligence and Statistics*, 2020.
- 551 Dereich, S., Graeber, R., and Jentzen, A. Non-convergence
552 of Adam and other adaptive stochastic gradient descent
553 optimization methods for non-vanishing learning rates.
554 *arXiv preprint arXiv:2407.08100*, 2024.
- 555 Dicker, L. H. Ridge regression and asymptotic minimax
556 estimation over spheres of growing dimension. *Bernoulli*,
557 pp. 1–37, 2016.
- 558 Dobriban, E. and Wager, S. High-dimensional asymptotics
559 of prediction: Ridge regression and classification. *The
560 Annals of Statistics*, 46(1):247–279, 2018.
- 561 Du, S. S., Hu, W., Kakade, S. M., Lee, J. D., and Lei,
562 Q. Few-shot learning via learning the representation,
563 provably. In *International Conference on Learning Rep-
564 resentations*, 2021.
- 565 Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient
566 methods for online learning and stochastic optimization.
567 *Journal of Machine Learning Research*, 12(7), 2011.
- 568 Frerix, T., Möllenhoff, T., Moeller, M., and Cremers, D.
569 Proximal backpropagation. In *International Conference
570 on Learning Representations*, 2018.
- 571 Fu, H., Wang, Z., Nichani, E., and Lee, J. D. Learning hier-
572 archical polynomials of multiple nonlinear features with
573 three-layer networks. *arXiv preprint arXiv:2411.17201*,
574 2024.
- 575 Ghorbani, B., Mei, S., Misiakiewicz, T., and Montanari, A.
576 Linearized two-layers neural networks in high dimension.
577 *The Annals of Statistics*, 49(2):1029–1054, 2021a.
- 578 Ghorbani, B., Mei, S., Misiakiewicz, T., and Montanari, A.
579 When do neural networks outperform kernel methods?
580 *Journal of Statistical Mechanics: Theory and Experiment*,
581 2021(12), 2021b.
- 582 Goldfarb, D. A family of variable-metric methods derived
583 by variational means. *Mathematics of computation*, 24
584 (109):23–26, 1970.
- 585 Goldfarb, D., Ren, Y., and Bahamou, A. Practical quasi-
586 Netwon methods for training deep neural networks. In *Ad-
587 vances in Neural Information Processing Systems*, 2020.

- 550 Goldt, S., Loureiro, B., Reeves, G., Krzakala, F., Mézard,
551 M., and Zdeborová, L. The Gaussian equivalence of gen-
552 erative models for learning with shallow neural networks.
553 In *Mathematical and Scientific Machine Learning*, pp.
554 426–471, 2022.
- 555 Guionnet, A., Ko, J., Krzakala, F., Mergny, P., and Zde-
556 borová, L. Spectral phase transitions in non-linear wigner
557 spiked models. *arXiv preprint arXiv:2310.14055*, 2023.
- 558 Gupta, A., Ramanath, R., Shi, J., and Keerthi, S. S. Adam
559 vs. SGD: Closing the generalization gap on image clas-
560 sification. In *OPT2021: 13th Annual Workshop on Opti-
561 mization for Machine Learning*, 2021.
- 562 Gupta, V., Koren, T., and Singer, Y. Shampoo: Precondi-
563 tioned stochastic tensor optimization. In *International
564 Conference on Machine Learning*, 2018.
- 565 Hanin, B. and Nica, M. Finite depth and width corrections
566 to the neural tangent kernel. In *International Conference
567 on Learning Representations*, 2020.
- 568 Hanson, D. L. and Wright, F. T. A bound on tail probabili-
569 ties for quadratic forms in independent random variables.
570 *The Annals of Mathematical Statistics*, 42(3):1079–1083,
571 1971.
- 572 Hassani, H. and Javanmard, A. The curse of over-
573 parametrization in adversarial training: Precise analysis
574 of robust generalization for random features regression.
575 *The Annals of Statistics*, 52(2):441–465, 2024.
- 576 Horn, R. A. and Johnson, C. R. *Matrix analysis*. Cambridge
577 university press, 2012.
- 578 Hsu, D., Kakade, S. M., and Zhang, T. Random design
579 analysis of ridge regression. In *Conference on Learning
580 Theory*, 2012.
- 581 Hu, H. and Lu, Y. M. Universality laws for high-dimensional
582 learning with random features. *IEEE Transactions on
583 Information Theory*, 69(3), 2023.
- 584 Ioffe, S. and Szegedy, C. Batch normalization: accelerating
585 deep network training by reducing internal covariate shift.
586 In *Proceedings of the 32nd International Conference on
587 International Conference on Machine Learning - Volume
588 37*, pp. 448–456. JMLR.org, 2015.
- 589 Ishikawa, S. and Karakida, R. On the parameterization of
590 second-order optimization effective towards the infinite
591 width. *arXiv preprint arXiv:2312.12226*, 2023.
- 592 Jacot, A., Gabriel, F., and Hongler, C. Neural tangent ker-
593 nel: Convergence and generalization in neural networks.
594 In *Advances in Neural Information Processing Systems*,
595 2018.
- 596 Jain, P., Netrapalli, P., and Sanghavi, S. Low-rank matrix
597 completion using alternating minimization. In *ACM Sym-
598 posium on Theory of Computing*, pp. 665–674, 2013.
- 599 Jordan, K., Jin, Y., Boza, V., Jiacheng, Y., Cecista, F., New-
600 house, L., and Bernstein, J. Muon: An optimizer for
601 hidden layers in neural networks, 2024. URL <https://kellerjordan.github.io/posts/muon/>.
- 602 Keskar, N. S. and Socher, R. Improving generalization
603 performance by switching from Adam to SGD. *arXiv
604 preprint arXiv:1712.07628*, 2017.
- 605 Kingma, D. P. and Ba, J. Adam: A method for stochastic
606 optimization. In *International Conference on Learning
607 Representations*, 2015.
- 608 Kumar, A., Raghunathan, A., Jones, R. M., Ma, T., and
609 Liang, P. Fine-tuning can distort pretrained features and
610 underperform out-of-distribution. In *International Con-
611 ference on Learning Representations*, 2022.
- 612 Large, T., Liu, Y., Huh, M., Bahng, H., Isola, P., and Bern-
613 stein, J. Scalable optimization in the modular norm. *arXiv
614 preprint arXiv:2405.14813*, 2024.
- 615 Laurent, B. and Massart, P. Adaptive estimation of a
616 quadratic functional by model selection. *The Annals
617 of Statistics*, pp. 1302–1338, 2000.
- 618 Lee, D., Moniri, B., Huang, X., Dobriban, E., and Hassani,
619 H. Demystifying disagreement-on-the-line in high dimen-
620 sions. In *International Conference on Machine Learning*,
621 2023a.
- 622 Lee, J. D., Oko, K., Suzuki, T., and Wu, D. Neural net-
623 work learns low-dimensional polynomials with SGD
624 near the information-theoretic limit. *arXiv preprint
625 arXiv:2406.01581*, 2024.
- 626 Lee, Y., Chen, A. S., Tajwar, F., Kumar, A., Yao, H., Liang,
627 P., and Finn, C. Surgical fine-tuning improves adaptation
628 to distribution shifts. In *International Conference on
629 Learning Representations*, 2023b.
- 630 Lin, L. and Dobriban, E. What causes the test error? going
631 beyond bias-variance via ANOVA. *Journal of Machine
632 Learning Research*, 22:155–1, 2021.
- 633 Lin, W., Dangel, F., Eschenhagen, R., Bae, J., Turner, R. E.,
634 and Makhzani, A. Can we remove the square-root in
635 adaptive gradient methods? a second-order perspective.
636 In *Forty-first International Conference on Machine Learn-
637 ing*, 2024.
- 638 Liu, D. C. and Nocedal, J. On the limited memory bfgs
639 method for large scale optimization. *Mathematical pro-
640 gramming*, 45(1):503–528, 1989.

- 605 Martens, J. Deep learning via Hessian-free optimization. In
606 *International Conference on Machine Learning*, 2010.
607
608 Martens, J. New insights and perspectives on the natural
609 gradient method. *Journal of Machine Learning Research*,
610 21(146):1–76, 2020.
611
612 Martens, J. and Grosse, R. Optimizing neural networks with
613 Kronecker-factored approximate curvature. In *International*
614 *Conference on Machine Learning*, 2015.
615
616 Maurer, A., Pontil, M., and Romera-Paredes, B. The benefit
617 of multitask representation learning. *Journal of Machine*
Learning Research, 17(81):1–32, 2016.
618
619 Mei, S. and Montanari, A. The generalization error of
620 random features regression: Precise asymptotics and the
621 double descent curve. *Communications on Pure and*
Applied Mathematics, 75(4):667–766, 2022.
622
623 MLCommons. Announcing the results of the inaugural
624 AlgoPerf: Training algorithms benchmark competition,
625 2024. URL [https://mlcommons.org/2024/08/
626 mlc-algoperf-benchmark-competition/](https://mlcommons.org/2024/08/mlc-algoperf-benchmark-competition/).
627
628 Moniri, B. and Hassani, H. Asymptotics of linear re-
629 gression with linearly dependent data. *arXiv preprint*
arXiv:2412.03702, 2024a.
630
631 Moniri, B. and Hassani, H. Signal-plus-noise decomposition
632 of nonlinear spiked random matrix models. *arXiv preprint*
arXiv:2405.18274, 2024b.
633
634 Moniri, B., Lee, D., Hassani, H., and Dobriban, E. A theory
635 of non-linear feature learning with one gradient step in
636 two-layer neural networks. In *International Conference*
637 *on Machine Learning*, 2024.
638
639 Montanari, A., Ruan, F., Sohn, Y., and Yan, J. The generalization
640 error of max-margin linear classifiers: High-
641 dimensional asymptotics in the overparametrized regime.
642 *arXiv preprint arXiv:1911.01544*, 2019.
643
644 Morwani, D., Shapira, I., Vyas, N., Malach, E., Kakade,
645 S., and Janson, L. A new perspective on shampoo’s
646 preconditioner. *arXiv preprint arXiv:2406.17748*, 2024.
647
648 Mousavi-Hosseini, A., Wu, D., Suzuki, T., and Erdogdu,
649 M. A. Gradient-based feature learning under structured
650 data. In *Advances in Neural Information Processing*
651 *Systems*, 2023.
652
653 Nakhleh, J., Shenouda, J., and Nowak, R. D. The effects of
654 multi-task learning on ReLU neural network functions.
655 *arXiv preprint arXiv:2410.21696*, 2024.
656
657 Nayer, S. and Vaswani, N. Fast and sample-efficient feder-
658 ated low rank matrix recovery from column-wise linear
659 and quadratic projections. *IEEE Transactions on Infor-
660 mation Theory*, 69(2):1177–1202, 2022.
661
662 Nichani, E., Damian, A., and Lee, J. D. Provable guaran-
663 tees for nonlinear feature learning in three-layer neural
664 networks. In *Advances in Neural Information Processing*
665 *Systems*, 2024a.
666
667 Nichani, E., Damian, A., and Lee, J. D. How transformers
668 learn causal structure with gradient descent. In *Inter-
669 national Conference on Machine Learning*, 2024b.
670
671 Nocedal, J. and Wright, S. J. *Numerical optimization*.
672 Springer, 1999.
673
674 Pearlmutter, B. A. Fast exact multiplication by the hessian.
675 *Neural computation*, 6(1):147–160, 1994.
676
677 Rahimi, A. and Recht, B. Random features for large-scale
678 kernel machines. In *Advances in Neural Information*
679 *Processing Systems*, 2007.
680
681 Reddi, S. J., Kale, S., and Kumar, S. On the convergence
682 of Adam and beyond. In *International Conference on*
Learning Representations, 2018.
683
684 Rudelson, M. and Vershynin, R. Hanson-wright inequality
685 and sub-gaussian concentration. *Electronic Communica-
686 tions in Probability*, 18:1–9, 2013.
687
688 Schmidt, R. M., Schneider, F., and Hennig, P. Descend-
689 ing through a crowded valley-benchmarking deep learn-
690 ing optimizers. In *International Conference on Machine*
Learning, 2021.
691
692 Schraudolph, N. N. Fast curvature matrix-vector products
693 for second-order gradient descent. *Neural computation*,
694 14(7):1723–1738, 2002.
695
696 Shi, H.-J. M., Lee, T.-H., Iwasaki, S., Gallego-Posada, J.,
697 Li, Z., Rangadurai, K., Mudigere, D., and Rabbat, M. A
698 distributed data-parallel pytorch implementation of the
699 distributed Shampoo optimizer for training neural net-
700 works at-scale. *arXiv preprint arXiv:2309.06497*, 2023.
701
702 Shi, Z., Wei, J., and Liang, Y. A theoretical analysis on
703 feature learning in neural networks: Emergence from
704 inputs and advantage over fixed features. In *International*
Conference on Learning Representations, 2022.
705
706 Silverstein, J. W. and Choi, S.-I. Analysis of the limiting
707 spectral distribution of large dimensional random mat-
708 rices. *Journal of Multivariate Analysis*, 54(2):295–309,
709 1995.
710
711 Stewart, G. W. and Sun, J.-g. *Matrix perturbation theory*.
712 Academic press, 1990.
713
714 Thekumparampil, K. K., Jain, P., Netrapalli, P., and Oh,
715 S. Sample efficient linear meta-learning by alternating
716 minimization. *arXiv preprint arXiv:2105.08306*, 2021.

- 660 Tielemans, T. and Hinton, G. Lecture 6.5-rmsprop: Divide
661 the gradient by a running average of its recent magnitude.
662 *Coursera: Neural Networks for Machine Learning*, 4(2):
663 26, 2012.
- 664 Trefethen, L. N. and Bau, D. *Numerical linear algebra*,
665 volume 181. SIAM, 2022.
- 666 Tripuraneni, N., Jordan, M., and Jin, C. On the theory
667 of transfer learning: The importance of task diversity.
668 In *Advances in Neural Information Processing Systems*,
669 2020.
- 670 Tripuraneni, N., Adlam, B., and Pennington, J. Over-
671 parameterization improves robustness to covariate shift in
672 high dimensions. In *Advances in Neural Information
673 Processing Systems*, 2021a.
- 674 Tripuraneni, N., Jin, C., and Jordan, M. Provable
675 meta-learning of linear representations. In *International
676 Conference on Machine Learning*, 2021b.
- 677 Troiani, E., Dandi, Y., Defilippis, L., Zdeborová, L.,
678 Loureiro, B., and Krzakala, F. Fundamental limits of
679 weak learnability in high-dimensional multi-index mod-
680 els. *arXiv preprint arXiv:2405.15480*, 2024.
- 681 Vaswani, N. Efficient federated low rank matrix recovery
682 via alternating GD and minimization: A simple proof.
683 *IEEE Transactions on Information Theory*, 2024.
- 684 Vershynin, R. Introduction to the non-asymptotic analysis
685 of random matrices. In Eldar, Y. C. and Kutyniok, G.
686 (eds.), *Compressed Sensing: Theory and Applications*,
687 pp. 210–268. Cambridge University Press, 2012.
- 688 Vershynin, R. *High-dimensional probability: An introduc-
689 tion with applications in data science*, volume 47. Cam-
690 bridge University Press, 2018.
- 691 Vyas, N., Morwani, D., Zhao, R., Shapira, I., Brandfon-
692 brener, D., Janson, L., and Kakade, S. M. Soap: Improv-
693 ing and stabilizing Shampoo using Adam. In *OPT 2024:
694 Optimization for Machine Learning*, 2024.
- 695 Wainwright, M. J. *High-dimensional statistics: A non-
696 asymptotic viewpoint*, volume 48. Cambridge university
697 press, 2019.
- 698 Wang, Z., Nichani, E., and Lee, J. D. Learning hierarchi-
699 cal polynomials with three-layer neural networks. In
700 *International Conference on Learning Representations*,
701 2024.
- 702 Wilson, A. C., Roelofs, R., Stern, M., Srebro, N., and Recht,
703 B. The marginal value of adaptive gradient methods in
704 machine learning. In *Advances in Neural Information
705 Processing Systems*, 2017.
- 706 Woodbury, M. A. *Inverting modified matrices*. Department
707 of Statistics, Princeton University, 1950.
- 708 Xie, Z., Wang, X., Zhang, H., Sato, I., and Sugiyama, M.
709 Adaptive inertia: Disentangling the effects of adaptive
710 learning rate and momentum. In *International Conference
711 on Machine Learning*, 2022.
- 712 Yang, G. and Hu, E. J. Feature learning in infinite-width
713 neural networks. In *International Conference on Machine
714 Learning*, 2021a.
- 715 Yang, G. and Hu, E. J. Tensor programs iv: Feature learn-
716 ing in infinite-width neural networks. In *International
717 Conference on Machine Learning*, 2021b.
- 718 Zhang, G., Martens, J., and Grosse, R. B. Fast convergence
719 of natural gradient descent for over-parameterized neural
720 networks. In *Advances in Neural Information Processing
721 Systems*, 2019.
- 722 Zhang, T. T., Lee, B. D., Ziemann, I., Pappas, G. J., and
723 Matni, N. Guarantees for nonlinear representation learn-
724 ing: Non-identical covariates, dependent data, fewer sam-
725 ples. In *International Conference on Machine Learning*,
726 2024a.
- 727 Zhang, T. T., Toso, L. F., Anderson, J., and Matni, N.
728 Sample-efficient linear representation learning from non-
729 IID non-isotropic data. In *International Conference on
730 Learning Representations*, 2024b.
- 731 Ziemann, I., Tsiamis, A., Lee, B., Jedra, Y., Matni, N., and
732 Pappas, G. J. A tutorial on the non-asymptotic theory of
733 system identification. In *IEEE Conference on Decision
734 and Control*, 2023.

715 A. Extended Background and Related Work

716 **Preconditioners for Neural Network Optimization.** A significant research effort in neural network optimization has
717 been dedicated to understanding the role of preconditioning in convergence speed and generalization. Perhaps the most
718 widespread paradigm falls under the category of *entry-wise* (“diagonal”) preconditioners, whose notable members include
719 Adam (Kingma & Ba, 2015), (diagonal) AdaGrad (Duchi et al., 2011), RMSprop (Tieleman & Hinton, 2012), and their
720 innumerable relatives and descendants (see e.g. Schmidt et al. (2021); Dahl et al. (2023) for surveys). However, diagonal
721 preconditioners inherently do not fully capture inter-parameter dependencies, which are better captured by stronger curvature
722 estimates, e.g. Gauss-Newton approximations (Botev et al., 2017; Martens, 2020), L-BFGS (Byrd et al., 2016; Bollapragada
723 et al., 2018; Goldfarb et al., 2020). Toward making non-diagonal preconditioners scalable to neural networks, many works
724 (including the above) have made use of layer-wise *Kronecker-Factored* approximations, where each layer’s curvature
725 block is factored into a Kronecker product $\mathbf{Q} \otimes \mathbf{P}$. Perhaps the two most well-known examples are Kronecker-Factored
726 Approximate Curvature (KFAC) (Martens & Grosse, 2015) and Shampoo (Gupta et al., 2018; Anil et al., 2020), where
727 approximations are made to the Fisher Information and Gauss-Newton curvature, respectively. Many works have since
728 expanded on these ideas, such as by improving practical efficiency (Ba et al., 2017; Shi et al., 2023; Jordan et al., 2024; Vyas
729 et al., 2024) and defining generalized constructions (Dangel et al., 2020; Amid et al., 2022; Benzing, 2022). An interesting
730 alternate view subsumes certain preconditioners via steepest descent with respect to layer-wise (“modular”) norms (Large
731 et al., 2024; Bernstein & Newhouse, 2024a;b). We draw a connection therein by deriving the steepest descent norm that
732 Kronecker-Factored preconditioners correspond to; see Appendix D.3.

733 **Multi-task Representation Learning (MTRL).** Toward a broader notion of generalization, the goal of MTRL is to
734 characterize the benefits of learning a *shared* representation across distinct tasks. Various works focus on the generalization
735 properties given access to an empirical risk minimizer (ERM) (Maurer et al., 2016; Du et al., 2021; Tripuraneni et al., 2020;
736 Zhang et al., 2024a), with the latter work resolving the setting where distinct tasks may have different covariate distributions.
737 Closely related formulations have been studied in the context of distribution shift (Kumar et al., 2022; Lee et al., 2023b).
738 While these works consider general non-linear representations, access to an ERM obviates the (non-convex) optimization
739 component. As such, multiple works have studied algorithms for linear representation learning (Tripuraneni et al., 2021b;
740 Collins et al., 2021; Thekumparampil et al., 2021; Nayer & Vaswani, 2022) and specific non-linear variants (Collins et al.,
741 2024; Nakhleh et al., 2024). In contrast to the ERM works, which are mostly agnostic to the covariate distribution, all the
742 listed algorithmic works assume isotropic covariates $N(\mathbf{0}, \mathbf{I})$. Zhang et al. (2024b) show that isotropy is in fact a key enabler,
743 and propose an adjustment to handle general covariances. In this paper, we show that many prior linear representation
744 learning algorithms belong to the same family of (preconditioned) optimizers. We then propose an algorithm coinciding
745 with KFAC that achieves the first condition-number-free convergence rate.

746 **Nonlinear Feature Learning.** In the early phase of training, neural networks are shown to be essentially equivalent to
747 the kernel methods, and can be described by the neural tangent kernel (NTK). See Jacot et al. (2018); Mei & Montanari
748 (2022); Hu & Lu (2023). However, kernel methods are inherently limited and have a sample complexity superlinear in the
749 input dimension d for learning nonlinear functions (Ghorbani et al., 2021a;b). The main reason for this limitation is that
750 kernel methods use a set of fixed features that are not task specific. There has been a lot of interest in studying the benefits
751 of feature learning from a theoretical perspective (Bai & Lee (2020); Hanin & Nica (2020); Yang & Hu (2021a); Shi et al.
752 (2022); Abbe et al. (2022), etc.). In a setting with isotropic covariates $N(\mathbf{0}, \mathbf{I})$, it is shown that even a one-step of SGD update
753 on the first layer of a two-layer neural networks can learn good enough features to provide a significant sample complexity
754 improvement over kernel methods assuming that the target function has some low-dimensional structure (Damian et al.,
755 2022; Ba et al., 2022; Moniri et al., 2024; Cui et al., 2024; Dandi et al., 2024a;b;c; Arnaboldi et al., 2024; Lee et al., 2024)
756 and this has became a very popular model for studying feature learning. These results were later extended to three-layer
757 neural networks in which the first layer is kept at random initialization and the second layer is updated using one step of
758 SGD (Wang et al., 2024; Nichani et al., 2024a; Fu et al., 2024). Recently, Ba et al. (2024); Mousavi-Hosseini et al. (2023)
759 considered an anisotropic case where the covariance contains a planted signal about the target function and showed that a
760 single step of SGD can leverage this to better learn the target function. However, the general case of anisotropic covariate
761 distributions remains largely unexplored. In this paper, we study feature learning with two-layer neural networks with
762 general anisotropic covariates in single-index models and that one-step of SGD update has inherent limitations in this setting,
763 and the natural fix will coincide with applying the KFAC layer-wise preconditioner.

770 B. Proofs and Additional Details for Section 3.1

771 B.1. Convergence Rate Lower Bound of SGD

773 Our goal is to establish the following lower bound construction.

774 **Proposition 3.5.** *Let $\Sigma_x = I_{d_x}$, $n = \infty$. Choose any $d_x > k$, $d_y \geq k \geq 2$. Let the learner be given knowledge of \mathbf{F}_* , \mathbf{G}_* and $\text{dist}(\mathbf{G}_0, \mathbf{G}_*)$. However, assume the learner must fix $\eta_G > 0$ before observing \mathbf{G}_0 . Then, there exists $\mathbf{F}_* \in \mathbb{R}^{d_y \times k}$, $\mathbf{G}_*, \mathbf{G}_0 \in \mathbb{R}^{k \times d_x}$, such that $\mathbf{G}_T = SGD(\mathbf{G}_0; \eta_G, T)$ satisfies:*

$$778 \quad 779 \quad 780 \quad \text{dist}(\mathbf{G}_T, \mathbf{G}_*) \geq \left(1 - 4 \frac{\lambda_{\min}(\mathbf{F}_*^\top \mathbf{F}_*)}{\lambda_{\max}(\mathbf{F}_*^\top \mathbf{F}_*)}\right)^T \text{dist}(\mathbf{G}_0, \mathbf{G}_*).$$

781 *Proof of Proposition 3.5.* We prove the lower bound by construction. First, we write out the one-step SGD update given
782 step size η_G .

$$\begin{aligned} 784 \quad \mathbf{G}_+ &= \mathbf{G} - \eta_G \nabla_{\mathbf{G}} \hat{\mathcal{L}}(\mathbf{F}, \mathbf{G}) \\ 785 \quad &= \mathbf{G} - \frac{1}{n} \mathbf{F}^\top (\mathbf{F} \mathbf{G} \mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{X}) \\ 786 \quad &= \mathbf{G} - \mathbf{F}^\top (\mathbf{F} \mathbf{G} - \mathbf{F}_* \mathbf{G}_*). \quad (\mathbf{Y}^\top = \mathbf{F}_* \mathbf{G}_* \mathbf{X}^\top + \mathcal{E}^\top, n = \infty, \Sigma_x = \mathbf{I}) \\ 787 \quad &\mathbf{G}_+ = \text{Ortho}(\mathbf{G}_+). \\ 788 \end{aligned}$$

789 We recall \mathbf{F} is given by the \mathbf{F} -update in (5) with $\eta_F = 1$, which by Lemma 3.1 is equivalent to setting \mathbf{F} to the least-squares
790 solution conditional on \mathbf{G} :

$$\begin{aligned} 793 \quad \mathbf{F} &= \mathbf{Y}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} = \mathbf{F}_* \mathbf{G}_* \Sigma_x \mathbf{G}^\top (\mathbf{G} \Sigma_x \mathbf{G}^\top)^{-1} \quad (\mathbf{z} = \mathbf{Gx}, n = \infty) \\ 794 \quad &= \mathbf{F}_* \mathbf{G}_* \mathbf{G}^\top (\mathbf{G} \mathbf{G}^\top)^{-1} = \mathbf{F}_* \mathbf{G}_* \mathbf{G}^\top. \quad (\mathbf{G} \text{ row-orthonormal}) \\ 795 \end{aligned}$$

796 Therefore, plugging in \mathbf{F} into the SGD update yields:

$$\begin{aligned} 797 \quad \mathbf{G}_+ &= \mathbf{G} - \eta_G \mathbf{F}^\top (\mathbf{F} \mathbf{G} - \mathbf{F}_* \mathbf{G}_*) \\ 798 \quad &= \mathbf{G} - \eta_G \mathbf{G} \mathbf{G}_*^\top \mathbf{F}_*^\top (\mathbf{F}_* \mathbf{G}_* \mathbf{G}^\top \mathbf{G} - \mathbf{F}_* \mathbf{G}_*). \\ 799 \end{aligned}$$

800 Before proceeding, let us present the construction of \mathbf{F}_* , \mathbf{G}_* . We focus on the case $d_x = 3$, $k = 2$, as it will be clear the
801 construction is trivially embedded to arbitrary $d_x > k \geq 2$. We observe that \mathbf{F}_* only appears in the SGD update in the form
802 $\mathbf{F}_*^\top \mathbf{F}_* \in \mathbb{R}^{k \times k}$, thus $d_y \geq k$ can be set arbitrarily as long as $\mathbf{F}_*^\top \mathbf{F}_*$ satisfies our specifications. Set \mathbf{F}_* , \mathbf{G}_* such that

$$803 \quad \mathbf{F}_*^\top \mathbf{F}_* = \begin{bmatrix} 1 - \lambda & 0 \\ 0 & \lambda \end{bmatrix}, \lambda \in (0, 1/2], \quad \mathbf{G}_* = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}.$$

804 Accordingly, the initial representation \mathbf{G}_0 (which the learner is not initially given) will have form

$$805 \quad \mathbf{G}_0 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \sqrt{1 - \varepsilon_0^2} & \varepsilon_0 \end{bmatrix}, \text{ or } \begin{bmatrix} \sqrt{1 - \varepsilon_0^2} & 0 & \varepsilon_0 \\ 0 & 1 & 0 \end{bmatrix}.$$

806 We prove all results with the first form of \mathbf{G}_0 , as all results will hold for the second with the only change swapping $\lambda, 1 - \lambda$.
807 It is clear that we may extend to arbitrary $d_x > k \geq 2$ by setting:

$$808 \quad \mathbf{F}_*^\top \mathbf{F}_* = \begin{bmatrix} (1 - \lambda) \mathbf{I}_{k-1} & \mathbf{0} \\ \mathbf{0} & \lambda \end{bmatrix}, \quad \mathbf{G}_* = [\mathbf{I}_k \quad \mathbf{0}_{d_x-k}], \quad \mathbf{G}_0 = \begin{bmatrix} \mathbf{I}_{k-1} & \mathbf{0} & \cdots & \mathbf{0} \\ 0 & \sqrt{1 - \varepsilon_0^2} & \varepsilon_0 & \mathbf{0} \end{bmatrix}.$$

809 Returning to the $d_x = 3$, $k = 2$ case, we first prove the following invariance result.

810 **Lemma B.1.** *Given $\mathbf{G}_0 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \sqrt{1 - \varepsilon_0^2} & \varepsilon_0 \end{bmatrix}$, then for any $t \geq 0$, $\mathbf{G}_t = \begin{bmatrix} 1 & 0 & 0 \\ 0 & c_1 & c_2 \end{bmatrix}$ for some $c_1^2 + c_2^2 = 1$. Furthermore,
811 we have $\text{dist}(\mathbf{G}_t, \mathbf{G}_*) = |c_2|$.*

825 *Proof of Lemma B.1.* This follows by induction. The base case follows by definition of \mathbf{G}_0 . Now given $\mathbf{G}_t = \begin{bmatrix} 1 & 0 & 0 \\ 0 & c_1 & c_2 \end{bmatrix}$
 826 for some c_1, c_2 , we observe that
 827

$$828 \quad 829 \quad 830 \quad \mathbf{G}_t \mathbf{G}_\star^\top = \begin{bmatrix} 1 & 0 \\ 0 & c_1 \end{bmatrix}, \quad \mathbf{F}_\star^\top \mathbf{F}_\star = \begin{bmatrix} 1 - \lambda & 0 \\ 0 & \lambda \end{bmatrix}.$$

831 Notably, we may write

$$\begin{aligned} 832 \quad \mathbf{G}_{t+1} &= \mathbf{G}_t - \eta_{\mathbf{G}} \mathbf{G}_t \mathbf{G}_\star^\top \mathbf{F}_\star^\top (\mathbf{F}_\star \mathbf{G}_\star \mathbf{G}_t^\top \mathbf{G}_t - \mathbf{F}_\star \mathbf{G}_\star) \\ 833 &= (\mathbf{I}_k - \eta_{\mathbf{G}} \mathbf{G}_t \mathbf{G}_\star^\top \mathbf{F}_\star^\top \mathbf{F}_\star \mathbf{G}_\star \mathbf{G}_t^\top) \mathbf{G}_t + \eta_{\mathbf{G}} \mathbf{G}_t \mathbf{G}_\star^\top \mathbf{F}_\star^\top \mathbf{F}_\star \mathbf{G}_\star \\ 834 &= \left(\mathbf{I}_k - \eta_{\mathbf{G}} \begin{bmatrix} 1 - \lambda & 0 \\ 0 & c_1^2 \lambda \end{bmatrix} \right) \mathbf{G}_t + \eta_{\mathbf{G}} \begin{bmatrix} 1 - \lambda & 0 \\ 0 & c_1 \lambda \end{bmatrix} \mathbf{G}_\star \\ 835 &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & c_1(1 + \eta_{\mathbf{G}} c_2^2 \lambda) & c_2(1 - \eta_{\mathbf{G}} c_1^2 \lambda) \end{bmatrix} \\ 836 &\quad \mathbf{G}_{t+1} = \text{Ortho}(\mathbf{G}_{t+1}). \end{aligned} \tag{15}$$

837 Therefore, \mathbf{G}_{t+1} shares the same support as \mathbf{G}_t , and by the orthonormalization step, the squared entries of the second row
 838 of \mathbf{G}_{t+1} equal 1, completing the induction step.
 839

840 To prove the second claim, we see that $\mathcal{P}_\star^\perp = \begin{bmatrix} \mathbf{0}_2 & \\ & 1 \end{bmatrix}$, and since \mathbf{G}_t is by assumption row-orthonormal, we have
 841

$$842 \quad \text{dist}(\mathbf{G}_t, \mathbf{G}_\star) = \|\mathbf{G}_t \mathcal{P}_\star^\perp\| = \left\| \begin{bmatrix} 0 & 0 & c_2 \end{bmatrix}^\top \right\|_{\text{op}} = |c_2|,$$

843 completing the proof. \square

844 With these facts in hand, we prove the following stability limit of the step-size, and the consequences for the contraction rate.
 845

846 **Lemma B.2.** If $\eta_{\mathbf{G}} \geq \frac{4}{1-\lambda}$, then for any given $\text{dist}(\mathbf{G}_0, \mathbf{G}_\star)$ we may find \mathbf{G}_0 such that $\limsup_t \text{dist}(\mathbf{G}_t, \mathbf{G}_\star) \geq \frac{1}{2}$.

847 *Proof of Lemma B.2.* By assumption $\lambda \leq 1/2$ and thus $\lambda \leq 1 - \lambda$. Evaluating Lemma B.1 instead on $\mathbf{G}_0 =$
 848 $\begin{bmatrix} \sqrt{1 - \varepsilon_0^2} & 0 & \varepsilon_0 \\ 0 & 1 & 0 \end{bmatrix}$, writing out (15) yields symmetrically:

$$\begin{aligned} 849 \quad \text{dist}(\mathbf{G}_t, \mathbf{G}_\star) &= |c_2| \\ 850 \quad \mathbf{G}_{t+1} &= \begin{bmatrix} c_1(1 + \eta_{\mathbf{G}} c_2^2(1 - \lambda)) & 0 & c_2(1 - \eta_{\mathbf{G}} c_1^2(1 - \lambda)) \\ 0 & 1 & 0 \end{bmatrix} \\ 851 \quad \mathbf{G}_{t+1} &= \text{Ortho}(\mathbf{G}_{t+1}). \end{aligned}$$

852 We first observe that regardless of $\eta_{\mathbf{G}}$, the norm of the first row \mathbf{G}_{t+1} is always greater than 1 pre-orthonormalization. Let
 853 us define $\omega = \eta_{\mathbf{G}}(1 - \lambda)$. Then, the squared-norm of the first row satisfies:

$$854 \quad (c_1(1 + \omega c_2^2))^2 + (c_2(1 - \omega c_1^2))^2 = 1 + \omega^2 c_1^2 c_2^2.$$

855 Therefore, the norm is strictly bounded away from 1 when $\omega > 0$ and either $c_1, c_2 \neq 0$ by the constraint $c_1^2 + c_2^2 = 1$.
 856 Importantly, this implies that regardless of the step-size taken, the resulting first-row norm of \mathbf{G}_+ must exceed 1 prior to
 857 orthonormalization. Given this property, we observe that for $\omega \geq 1/c_1^2$, we have:
 858

$$859 \quad \frac{|1 - \omega c_1^2|}{|1 + \omega c_2^2|} = \frac{\omega c_1^2 - 1}{1 + \omega c_2^2}.$$

860 When this ratio is greater than 1, we are guaranteed that the first-row coefficients c'_1, c'_2 of \mathbf{G}_{t+1} post-orthonormalization
 861 satisfy $c'_2/c'_1 > c_2/c_1$, and recall from Lemma B.1 $\text{dist}(\mathbf{G}_{t+1}, \mathbf{G}_\star) = c'_2$, and thus $\text{dist}(\mathbf{G}_{t+1}, \mathbf{G}_\star) > \text{dist}(\mathbf{G}_t, \mathbf{G}_\star)$.
 862 Rearranging the above ratio, this is equivalent to the condition $\omega = \eta_{\mathbf{G}}(1 - \lambda) \geq \frac{2}{c_1^2 - c_2^2}$, $\omega = \eta_{\mathbf{G}}(1 - \lambda) \geq 1/c_1^2$.
 863

864 Setting $c_1^2 = 3/4, c_2^2 = 1/4$, this implies for $\eta_{\mathbf{G}} \geq \frac{4}{1-\lambda}$, the moment $\text{dist}(\mathbf{G}_t, \mathbf{G}_\star) \leq c_2 = 1/2$, then we are guaranteed
 865 $\text{dist}(\mathbf{G}_{t+1}, \mathbf{G}_\star) > \text{dist}(\mathbf{G}_t, \mathbf{G}_\star)$, and thus $\limsup_{t \rightarrow \infty} \text{dist}(\mathbf{G}_t, \mathbf{G}_\star) \geq \frac{1}{2}$, irregardless of $\text{dist}(\mathbf{G}_0, \mathbf{G}_\star)$. \square

Now, to finish the construction of the lower bound, Lemma B.2 establishes that $\eta_G \leq \frac{4}{1-\lambda}$ is necessary for convergence (though not sufficient!). This implies that when we plug back in $G_0 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \sqrt{1-\varepsilon_0^2} & \varepsilon_0 \end{bmatrix}$, we have G_{t+1} :

$$G_{t+1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & c_1(1 + \eta_G c_2^2 \lambda) & c_2(1 - \eta_G c_1^2 \lambda) \end{bmatrix}$$

$$G_{t+1} = \text{Ortho}(G_{t+1}) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & c'_1 & c'_2 \end{bmatrix}.$$

We have trivially $1 - \eta_G c_1^2 \lambda \geq 1 - \eta_G \lambda$. Therefore, for $\lambda \leq 1/5$ such that $\frac{\lambda}{1-\lambda} \leq 1/4$, we have $1 - \eta_G \lambda \geq 1 - \frac{4\lambda}{1-\lambda} \geq 0$. As shown in the proof of Lemma B.2, the norm of the second row pre-orthonormalization is strictly greater than 1, and thus:

$$\text{dist}(G_{t+1}, G_*) = c'_2 \geq c_2(1 - \eta_G \lambda c_1^2) \geq (1 - \eta_G \lambda) \text{dist}(G_t, G_*) \geq \left(1 - 4 \frac{\lambda}{1-\lambda}\right) \text{dist}(G_t, G_*).$$

Applying this recursively to G_0 yields the desired lower bound. \square

B.2. Proof of Theorem 3.6

Recall that running an iteration of stylized KFAC (5) with $\lambda_F, \lambda_G = 0, \eta_F = 1$ yields:

$$\begin{aligned} \bar{G}_+ &= G - \eta_G P_G^{-1} \nabla_G \hat{\mathcal{L}}(F, G) (Q_G + \lambda_G I_{d_X})^{-1} \\ &= G - \eta_G (F^\top F)^{-1} F^\top (F G \hat{\Sigma}_x - F_* G_* \hat{\Sigma}_x - \frac{1}{n} \mathcal{E}^\top X) \hat{\Sigma}_x^{-1} \\ &= G - \eta_G (F^\top F)^{-1} F^\top (F G - F_* G_*) + (F^\top F)^{-1} F^\top \mathcal{E}^\top X (X^\top X)^{-1}, \end{aligned} \quad (16)$$

where the matrix F is given by

$$\begin{aligned} F &= F_{\text{prev}} - \eta_F P_F^{-1} \nabla_F \hat{\mathcal{L}}(F_{\text{prev}}, G) (Q_F + \lambda_F I_{d_h})^{-1} \\ &= Y^\top Z (Z^\top Z)^{-1} \\ &= F_* G_* X^\top Z (Z^\top Z)^{-1} + \mathcal{E}^\top Z (Z^\top Z)^{-1}, \end{aligned}$$

recalling that $z \triangleq Gx$. Focusing on the representation update, we have

$$\bar{G}_+ \mathcal{P}_*^\perp = (1 - \eta_G) G \mathcal{P}_*^\perp + \eta_G (F^\top F)^{-1} F^\top \mathcal{E}^\top X (X^\top X)^{-1}.$$

Therefore, to prove a one-step contraction of $\text{rowsp}(G_+)$ toward $\text{rowsp}(G_*)$, we require two main components:

- Bounding the noise term $\eta_G (F^\top F)^{-1} F^\top \mathcal{E}^\top X (X^\top X)^{-1}$.
- Bounding the orthonormalization factor; the subspace distance measures distance between two orthonormalized bases (a.k.a. elements of the Stiefel manifold (Absil et al., 2008)), while a step of SGD or KFAC does not inherently conform to the Stiefel manifold, and thus the “off-manifold” shift must be considered when computing $\text{dist}(G_+, G_*)$. This amounts to bounding the “R”-factor of the QR-decomposition (Trefethen & Bau, 2022) of \bar{G}_+ .

Thanks to the left-preconditioning by $(F^\top F)^{-1}$, the contraction factor is essentially determined by $(1 - \eta_G)$; however, the second point about the “off-manifold” shift is what prevents us from setting $\eta_G = 1$.

Bounding the noise term

We start by observing $\mathcal{E}^\top X (X^\top X)^{-1} = \sum_{i=1}^n \varepsilon_i x_i^\top (\sum_{i=1}^n x_i x_i^\top)^{-1}$ (and thus $(F^\top F)^{-1} F^\top \mathcal{E}^\top X (X^\top X)^{-1}$) is a least-squares error-like term, and thus can be bounded by standard self-normalized martingale arguments. In particular, defining $\bar{F} \triangleq (F^\top F)^{-1} F^\top$ we may decompose

$$\|\bar{F} \mathcal{E}^\top X (X^\top X)^{-1}\|_{\text{op}} \leq \left\| \bar{F} \mathcal{E}^\top X (X^\top X)^{-1/2} \right\|_{\text{op}} \lambda_{\min}(X^\top X)^{-1/2},$$

where the first factor is the aforementioned self-normalized martingale (see e.g. Abbasi-Yadkori & Szepesvári (2011); Ziemann et al. (2023)), and the second can be bounded by standard covariance lower-tail bounds. Toward bounding the first factor, we invoke a high-probability self-normalized bound:

935 **Lemma B.3** (cf. Ziemann et al. (2023, Theorem 4.1)). Let $\{\mathbf{v}_i, \mathbf{w}_i\}_{i \geq 1}$ be a $\mathbb{R}^{d_v} \times \mathbb{R}^{d_w}$ -valued process and $\{\mathcal{F}_i\}_{i \geq 1}$ be
 936 a filtration such that $\{\mathbf{v}_i\}_{i \geq 1}$ is adapted to $\{\mathcal{F}_i\}_{i \geq 1}$, $\{\mathbf{w}_i\}_{i \geq 1}$ is adapted to $\{\mathcal{F}_i\}_{i \geq 2}$, and $\{\mathbf{w}_i\}_{i \geq 1}$ is a σ^2 -subgaussian
 937 martingale difference sequence⁶. Fix (non-random) positive-definite matrix \mathbf{Q} . For $k \geq 1$, define $\widehat{\Sigma}_k \triangleq \sum_{i=1}^k \mathbf{v}_i \mathbf{v}_i^\top$. Then,
 938 given any fixed $n \in \mathbb{N}_+$, with probability at least $1 - \delta$:

$$940 \quad 941 \quad 942 \quad 943 \quad 944 \quad \left\| \sum_{i=1}^n \mathbf{w}_i \mathbf{v}_i^\top (\mathbf{Q} + \widehat{\Sigma}_n)^{-1/2} \right\|_{\text{op}}^2 \leq 4\sigma^2 \log \left(\frac{\det(\mathbf{Q} + \widehat{\Sigma}_n)}{\det(\mathbf{Q})} \right) + 13d_w\sigma^2 + 8\sigma^2 \log(1/\delta). \quad (17)$$

945 Instantiating this for Gaussian $\mathbf{w}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{w}})$, $\mathbf{v}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{v}})$, we may set $\mathbf{Q} \approx \Sigma_{\mathbf{v}}$ to yield:

946 **Lemma B.4.** Consider the quantities defined in Lemma B.3 and assume $\mathbf{w}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{w}})$, $\mathbf{v}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{v}})$, defining
 947 $\sigma_{\mathbf{w}}^2 \triangleq \lambda_{\max}(\Sigma_{\mathbf{w}})$, $\sigma_{\mathbf{v}}^2 \triangleq \lambda_{\max}(\Sigma_{\mathbf{v}})$. Then, as long as $n \gtrsim \frac{18.27}{c^2} (d_v + \log(1/\delta))$, with probability at least $1 - \delta$:

$$950 \quad 951 \quad 952 \quad 953 \quad 954 \quad 955 \quad \left\| \sum_{i=1}^n \mathbf{w}_i \mathbf{v}_i^\top (\widehat{\Sigma}_n)^{-1/2} \right\|_{\text{op}}^2 \leq 8d_v \log \left(\frac{1+c}{1-c} \right) \sigma_{\mathbf{w}}^2 + 26d_w\sigma_{\mathbf{w}}^2 + 16\sigma_{\mathbf{w}}^2 \log(1/\delta)$$

$$\lambda_{\min}(\widehat{\Sigma}_n) \geq (1-c) \lambda_{\min}(\Sigma_{\mathbf{v}}).$$

956 *Proof of Lemma B.4.* We observe that if $\widehat{\Sigma}_n \succeq \mathbf{Q}$, then

$$957 \quad 958 \quad 959 \quad 2\widehat{\Sigma}_n \succeq \mathbf{Q} + \widehat{\Sigma}_n \implies (\widehat{\Sigma}_n)^{-1} \preceq 2(\mathbf{Q} + \widehat{\Sigma}_n)^{-1}.$$

960 This implies

$$962 \quad 963 \quad 964 \quad 965 \quad \mathbf{1}\{\widehat{\Sigma}_n \succeq \mathbf{Q}\} \left\| \sum_{i=1}^n \mathbf{w}_i \mathbf{v}_i^\top (\widehat{\Sigma}_n)^{-1/2} \right\|^2 \leq 2\mathbf{1}\{\widehat{\Sigma}_n \succeq \mathbf{Q}\} \left\| \sum_{i=1}^n \mathbf{w}_i \mathbf{v}_i^\top (\mathbf{Q} + \widehat{\Sigma}_n)^{-1/2} \right\|^2. \quad (18)$$

966 Let us consider the event:

$$967 \quad 968 \quad (1-c)\Sigma_{\mathbf{v}} \preceq \widehat{\Sigma}_n \preceq (1+c)\Sigma_{\mathbf{v}},$$

969 which by Lemma E.2 occurs with probability at least $1 - \delta$ as long as $n \gtrsim \frac{18.27}{c^2} (d_v + \log(1/\delta))$. This immediately
 970 establishes the latter desired inequality. Setting $\mathbf{Q} = (1-c)n\Sigma_{\mathbf{v}}$ and conditioning on the above event, we observe that by
 971 definition $\widehat{\Sigma}_n \succeq \mathbf{Q}$, and

$$973 \quad 974 \quad 975 \quad 976 \quad 977 \quad 978 \quad 979 \quad 980 \quad \begin{aligned} \log \left(\frac{\det(\mathbf{Q} + \widehat{\Sigma}_n)}{\det(\mathbf{Q})} \right) &= \log \det(I_{d_v} + \widehat{\Sigma}_n (\mathbf{Q})^{-1}) \\ &\leq \log \det \left(\left(1 + \frac{1+c}{1-c}\right) I_{d_v} \right) \\ &\leq d_v \log \left(\frac{1+c}{1-c} \right). \end{aligned}$$

981 Plugging this into Lemma B.3, applied to the RHS of (18), we get our desired result. \square

984 Therefore, instantiating $\mathbf{w} \rightarrow \bar{\mathbf{F}}\boldsymbol{\varepsilon}$, $\mathbf{v} \rightarrow \mathbf{x}$, $(\mathbf{x}_i, \boldsymbol{\varepsilon}_i)_{i \geq 1}$ is a $\mathbb{R}^{d_x} \times \mathbb{R}^k$ -valued process. Furthermore, since we assumed out
 985 of convenience that \mathbf{F}, \mathbf{G}_+ are computed on independent batches of data, we have that $\bar{\mathbf{F}}\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \bar{\mathbf{F}}\Sigma_{\boldsymbol{\varepsilon}}\bar{\mathbf{F}}^\top)$. In order to
 986 complete the noise term bound, it suffices to provide a uniform bound on $\|\bar{\mathbf{F}}\| = 1/\sigma_{\min}(\mathbf{F})$ in terms of \mathbf{F}_* .

988 ⁶See Appendix E.3 for discussion of formalism. It suffices for our purposes to consider $\mathbf{w} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{w}})$.

990 **Lemma B.5.** Assume the following conditions hold:

$$991 \quad n \gtrsim \max \left\{ k + \log(1/\delta), \sigma_\epsilon^2 \frac{d_y + k + \log(1/\delta)}{\sigma_{\min}(\mathbf{F}_*)^2 \lambda_{\min}(\boldsymbol{\Sigma}_x)} \right\}$$

$$992 \quad \text{dist}(\mathbf{G}, \mathbf{G}_*) \leq \frac{2}{5} \kappa(\mathbf{F}_*)^{-1} \kappa(\boldsymbol{\Sigma}_x)^{-1},$$

993
994
995
996 then with probability at least $1 - \delta$, we have $\|\bar{\mathbf{F}}\| = 1 / \sigma_{\min}(\mathbf{F}) \leq 2 \sigma_{\min}(\mathbf{F}_*)^{-1}$.

997
998 *Proof of Lemma B.5.* Recall we may write \mathbf{F} as

$$999 \quad \mathbf{F} = \mathbf{F}_* \mathbf{G}_* \mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} + \mathcal{E}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1}$$

$$1000 \quad = \mathbf{F}_* \mathbf{G}_* \mathbf{G}^\top + \mathbf{F}_* \mathbf{G}_* (\mathbf{I}_{d_x} - \mathbf{G}^\top \mathbf{G}) \mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} + \mathcal{E}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \quad (19)$$

1001 By Weyl's inequality for singular values (Horn & Johnson, 2012), we have

$$1002 \quad \sigma_{\min}(\mathbf{F}) \geq \sigma_{\min}(\mathbf{F}_* \mathbf{G}_* \mathbf{G}^\top) - \sigma_{\max}(\mathbf{F}_* \mathbf{G}_* (\mathbf{I}_{d_x} - \mathbf{G}^\top \mathbf{G}) \mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} + \mathcal{E}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1})$$

1003 Since $\mathbf{G}_* \mathbf{G}^\top$ is an orthogonal matrix, the first term is equal to $\sigma_{\min}(\mathbf{F}_*)$. On the other hand, applying triangle inequality on
1004 the second term, for $n \gtrsim k + \log(1/\delta)$ we have:

$$1005 \quad \|\mathbf{F}_* \mathbf{G}_* (\mathbf{I}_{d_x} - \mathbf{G}^\top \mathbf{G}) \mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1}\|_{\text{op}} \leq \|\mathbf{F}_*\|_{\text{op}} \|\mathbf{G}_* (\mathbf{I}_{d_x} - \mathbf{G}^\top \mathbf{G})\|_{\text{op}} \|\mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1}\|_{\text{op}}$$

$$1006 \quad \leq \|\mathbf{F}_*\|_{\text{op}} \text{dist}(\mathbf{G}, \mathbf{G}_*) \left(\frac{5}{4} \|\boldsymbol{\Sigma}_x\|_{\text{op}} \lambda_{\min}(\mathbf{G} \boldsymbol{\Sigma}_x \mathbf{G}^\top) \right)$$

$$1007 \quad \leq \frac{5}{4} \|\mathbf{F}_*\|_{\text{op}} \text{dist}(\mathbf{G}, \mathbf{G}_*) \kappa(\boldsymbol{\Sigma}_x),$$

1008 where we used covariance concentration for the second inequality Lemma E.2 and the trivial bound $\lambda_{\min}(\mathbf{A} \boldsymbol{\Sigma} \mathbf{A}^\top) \geq \lambda_{\min}(\boldsymbol{\Sigma})$
1009 for the last inequality. In turn, we may bound:

$$1010 \quad \|\mathcal{E}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1}\|_{\text{op}} \leq \left\| \mathcal{E}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1/2} \right\|_{\text{op}} \lambda_{\min}(\mathbf{Z}^\top \mathbf{Z})^{-1/2}$$

$$1011 \quad \lesssim \left\| \mathcal{E}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1/2} \right\|_{\text{op}} n^{-1/2} \lambda_{\min}(\boldsymbol{\Sigma}_x)^{-1/2} \quad (\text{Lemma E.2})$$

$$1012 \quad \lesssim \sigma_\epsilon \sqrt{\frac{d_y + k + \log(1/\delta)}{\lambda_{\min}(\boldsymbol{\Sigma}_x) n}}.$$

1013 Therefore, setting $\text{dist}(\mathbf{G}, \mathbf{G}_*) \leq \frac{2}{5} \kappa(\mathbf{F}_*)^{-1} \kappa(\boldsymbol{\Sigma}_x)^{-1}$, and $n \gtrsim \sigma_\epsilon^2 \frac{d_y + k + \log(1/\delta)}{\sigma_{\min}(\mathbf{F}_*)^2 \lambda_{\min}(\boldsymbol{\Sigma}_x)}$, we have $\sigma_{\min}(\mathbf{F}) \geq \frac{1}{2} \sigma_{\min}(\mathbf{F}_*)$,
1014 which leads to our desired bound on $\|\bar{\mathbf{F}}\|$. \square

1015 With a bound on $\|\bar{\mathbf{F}}\|_{\text{op}}$, bounding the noise term is a straightforward application of Lemma B.4.

1016 **Proposition B.6** (KFAC noise term bound). Let the conditions in Lemma B.5 hold. In addition, assume $n \gtrsim d_x + \log(1/\delta)$.
1017 Then, with probability at least $1 - \delta$:

$$1018 \quad \|\bar{\mathbf{F}} \mathcal{E}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}\|_{\text{op}} \lesssim \sigma_\epsilon \sqrt{\frac{d_x + k + \log(1/\delta)}{\sigma_{\min}(\mathbf{F}_*)^2 \lambda_{\min}(\boldsymbol{\Sigma}_x) n}}.$$

1019 *Proof of Proposition B.6.* Condition on the event of Lemma B.5. Then, assuming $n \gtrsim d_x + \log(1/\delta)$, we may apply
1020 covariance concentration (Lemma E.2) on $\hat{\boldsymbol{\Sigma}}_x$ and Lemma B.4 to bound the noise term by:

$$1021 \quad \|\bar{\mathbf{F}} \mathcal{E}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}\|_{\text{op}} \leq \left\| \bar{\mathbf{F}} \mathcal{E}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1/2} \right\|_{\text{op}} \lambda_{\min}(\mathbf{X}^\top \mathbf{X})^{-1/2}$$

$$1022 \quad \leq \|\bar{\mathbf{F}}\|_{\text{op}} \left\| \mathcal{E}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1/2} \right\|_{\text{op}} n^{-1/2} \lambda_{\min}(\boldsymbol{\Sigma}_x)^{-1/2} \quad (\text{Lemma E.2})$$

$$1023 \quad \lesssim \sigma_\epsilon \sqrt{\frac{d_x + k + \log(1/\delta)}{\sigma_{\min}(\mathbf{F}_*)^2 \lambda_{\min}(\boldsymbol{\Sigma}_x) n}}, \quad (\text{Lemma B.4})$$

1024 which completes the proof. \square

1045 This completes the bound on the noise term. We proceed to the orthonormalization factor.
 1046

1047 Bounding the orthonormalization factor

1048 Toward bounding the orthonormalization factor from (8). Defining $\bar{\mathbf{G}}_+$ as the updated representation pre-orthonormalization,
 1049 we write $\bar{\mathbf{G}}_+ = \mathbf{R}\mathbf{G}_+$, where \mathbf{G}_+ is the orthonormalized representation and $\mathbf{R} \in \mathbb{R}^{k \times k}$ is the corresponding orthonormaliza-
 1050 tion factor. Therefore, defining the shorthand $\text{Sym}(\mathbf{A}) = \mathbf{A} + \mathbf{A}^\top$, we have
 1051

$$\begin{aligned} \mathbf{R}\mathbf{R}^\top &= \mathbf{R}\mathbf{G}_+(\mathbf{R}\mathbf{G}_+)^\top \\ &= (\mathbf{G} - \eta_{\mathbf{G}}(\mathbf{F}^\top \mathbf{F})^{-1}\mathbf{F}^\top(\mathbf{F}\mathbf{G} - \mathbf{F}_*\mathbf{G}_*) + (\mathbf{F}^\top \mathbf{F})^{-1}\mathbf{F}^\top \mathcal{E}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1})(\dots)^\top && \text{(from (16))} \\ &\succ \mathbf{I}_k - \eta_{\mathbf{G}} \text{Sym} \left(\underbrace{\bar{\mathbf{F}}(\mathbf{F}\mathbf{G} - \mathbf{F}_*\mathbf{G}_*)\mathbf{G}^\top}_{\triangleq \Gamma_1} \right) + \eta_{\mathbf{G}} \text{Sym} \left(\underbrace{\bar{\mathbf{F}}\mathcal{E}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{G}^\top}_{\triangleq \Gamma_2} \right) \\ &\quad - \eta_{\mathbf{G}}^2 \text{Sym} \left(\bar{\mathbf{F}}(\mathbf{F}\mathbf{G} - \mathbf{F}_*\mathbf{G}_*) (\bar{\mathbf{F}}\mathcal{E}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1})^\top \right), \end{aligned}$$

1060 where the strictly inequality comes from discarding the positive-definite “diagonal” terms of the expansion. Therefore, by
 1061 Weyl’s inequality for symmetric matrices (Horn & Johnson, 2012), we have:

$$1062 \lambda_{\min}(\mathbf{R}\mathbf{R}^\top) \geq 1 - 2\eta_{\mathbf{G}} \left(\|\Gamma_1\|_{\text{op}} + \|\Gamma_2\|_{\text{op}} + \eta_{\mathbf{G}} \|\Gamma_1\|_{\text{op}} \|\Gamma_2\|_{\text{op}} \right).$$

1065 Toward bounding $\|\Gamma_1\|_{\text{op}}$, let the conditions of Lemma B.5 hold. Then,

$$\begin{aligned} 1066 \|\Gamma_1\|_{\text{op}} &= \|\bar{\mathbf{F}}(\mathbf{F}\mathbf{G} - \mathbf{F}_*\mathbf{G}_*)\mathbf{G}^\top\|_{\text{op}} \\ 1067 &= \|(\mathbf{F}^\top \mathbf{F})^{-1}\mathbf{F}^\top \mathbf{F}\mathbf{G}\mathbf{G}^\top - (\mathbf{F}^\top \mathbf{F})^{-1}\mathbf{F}^\top \mathbf{F}_*\mathbf{G}_*\mathbf{G}^\top\|_{\text{op}} \\ 1068 &= \|\bar{\mathbf{F}}(\mathbf{F}_*\mathbf{G}_*(\mathbf{I}_{d_x} - \mathbf{G}^\top \mathbf{G})\mathbf{X}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1} + \mathcal{E}^\top \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1})\|_{\text{op}} && \text{(from (19))} \\ 1069 &\leq \frac{5}{4} \sigma_{\min}(\mathbf{F}) \|\mathbf{F}_*\|_{\text{op}} \text{dist}(\mathbf{G}, \mathbf{G}_*) \kappa(\Sigma_x) + \sigma_{\min}(\mathbf{F}) \sigma_\epsilon^2 \sqrt{\frac{k + \log(1/\delta)}{\lambda_{\min}(\Sigma_x)n}} \\ 1070 &\leq \underbrace{\frac{5}{2} \kappa(\mathbf{F}_*) \text{dist}(\mathbf{G}, \mathbf{G}_*) \kappa(\Sigma_x)}_{\triangleq \gamma_1} + \sigma_\epsilon^2 \sqrt{\frac{k + \log(1/\delta)}{\sigma_{\min}(\mathbf{F}_*)^2 \lambda_{\min}(\Sigma_x)n}}. && \text{(Lemma B.5)} \\ 1071 & \\ 1072 & \\ 1073 & \\ 1074 & \\ 1075 & \\ 1076 & \\ 1077 & \end{aligned}$$

1078 Similarly, letting the conditions of Proposition B.6 hold, we have

$$\begin{aligned} 1079 \|\Gamma_2\|_{\text{op}} &= \|\bar{\mathbf{F}}\mathcal{E}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{G}^\top\|_{\text{op}} \\ 1080 &\leq \sigma_\epsilon \sqrt{\frac{d_x + k + \log(1/\delta)}{\sigma_{\min}(\mathbf{F}_*)^2 \lambda_{\min}(\Sigma_x)n}} \\ 1081 &\triangleq \gamma_2. && \text{(Proposition B.6)} \\ 1082 & \\ 1083 & \\ 1084 & \\ 1085 & \end{aligned}$$

1086 We observe that γ_2 will always dominate the second term of the bound on $\|\Gamma_1\|_{\text{op}}$, and therefore:

$$1087 \lambda_{\min}(\mathbf{R}\mathbf{R}^\top) \geq 1 - 2\eta_{\mathbf{G}}(\gamma_1 + 2\gamma_2 + \eta_{\mathbf{G}}(\gamma_1 + \gamma_2)\gamma_2).$$

1088 Therefore, we have the following bound on the orthonormalization factor.

1089 **Proposition B.7.** *Let the following conditions hold:*

$$\begin{aligned} 1090 n &\gtrsim \max \left\{ d_x + \log(1/\delta), \frac{\sigma_\epsilon^2}{\gamma_2^2} \frac{d_x + \log(1/\delta)}{\sigma_{\min}(\mathbf{F}_*)^2 \lambda_{\min}(\Sigma_x)} \right\} \\ 1091 \text{dist}(\mathbf{G}, \mathbf{G}_*) &\leq \frac{2}{5\gamma_1} \kappa(\mathbf{F}_*)^{-1} \kappa(\Sigma_x)^{-1}. \\ 1092 & \\ 1093 & \\ 1094 & \\ 1095 & \\ 1096 & \end{aligned}$$

1097 Then, with probability at least $1 - \delta$, we have the following bound on the orthonormalization factor:

$$1098 \sigma_{\min}(\mathbf{R}) \geq \sqrt{1 - 2\eta_{\mathbf{G}}(\gamma_1 + 2\gamma_2 + \eta_{\mathbf{G}}(\gamma_1 + \gamma_2)\gamma_2)}. \\ 1099$$

1100 The constants γ_1, γ_2 will be instantiated to control the deflation of the contraction factor $1 - \eta_G \implies 1 - c\eta_G$ due to the
 1101 orthonormalization factor.

1102

1103 Completing the bound

1104 We are almost ready to complete the proof. By instantiating the noise bound Proposition B.6 and the orthonormalization
 1105 factor bound Proposition B.7, we have:

$$\begin{aligned} 1107 \quad \|\mathbf{G}_+ \mathcal{P}_\star^\perp\|_{\text{op}} &= \|\mathbf{R}^{-1} ((1 - \eta_G) \mathbf{G} \mathcal{P}_\star^\perp + \eta_G (\mathbf{F}^\top \mathbf{F})^{-1} \mathbf{F}^\top \mathcal{E}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1})\|_{\text{op}} \\ 1108 \quad &\leq \frac{1 - \eta_G}{\sigma_{\min}(\mathbf{R})} \|\mathbf{G} \mathcal{P}_\star^\perp\|_{\text{op}} + \eta_G \|(\mathbf{F}^\top \mathbf{F})^{-1} \mathbf{F}^\top \mathcal{E}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}\|_{\text{op}} \\ 1110 \quad &\leq \frac{1 - \eta_G}{\sqrt{1 - 2\eta_G(\gamma_1 + 2\gamma_2 + \eta_G(\gamma_1 + \gamma_2)\gamma_2)}} \|\mathbf{G} \mathcal{P}_\star^\perp\|_{\text{op}} + \eta_G \sigma_\varepsilon \sqrt{\frac{d_x + k + \log(1/\delta)}{\sigma_{\min}(\mathbf{F}_\star)^2 \lambda_{\min}(\Sigma_x) n}}. \end{aligned}$$

1111 To understand the effective deflation of the convergence rate, we prove the following numerical helper lemma.
 1112

1113 **Lemma B.8.** Given $c, d \in (0, 1)$ and $\varepsilon \in (0, 1/2)$, if $\varepsilon \geq c$, then the following holds:

1114

$$\frac{1 - d}{\sqrt{1 - cd}} < 1 - (1 - \varepsilon)d.$$

1115

1116 Additionally, as long as $\varepsilon \leq 1 - \frac{1 - \sqrt{1-d}}{d}$, then $1 - (1 - \varepsilon)d \leq \sqrt{1 - d}$.

1117

1118 *Proof of Lemma B.8:* squaring both sides of the desired inequality and re-arranging some terms, we arrive at

1119

$$\begin{aligned} 1120 \quad c &\leq \frac{1}{d} \left(1 - \frac{(1 - d)^2}{(1 - (1 - \varepsilon)d)^2} \right) \\ 1121 \quad &= \frac{1}{d} \left(1 - \underbrace{\frac{1 - d}{1 - (1 - \varepsilon)d}}_{<1} \right) \underbrace{\left(1 + \frac{1 - d}{1 - (1 - \varepsilon)d} \right)}_{>1}. \end{aligned}$$

1122 To certify the above inequality, it suffices to lower-bound the RHS. Since $c, d \in (0, 1)$, the last factor is at least 1, such that
 1123 we have

1124

$$\begin{aligned} 1125 \quad \frac{1}{d} \left(1 - \frac{1 - d}{1 - (1 - \varepsilon)d} \right) \left(1 + \frac{1 - d}{1 - (1 - \varepsilon)d} \right) &> \frac{1}{d} \left(1 - \frac{1 - d}{1 - (1 - \varepsilon)d} \right) \\ 1126 \quad &= \frac{1}{d} \frac{(1 - \varepsilon)d}{1 - (1 - \varepsilon)d} \\ 1127 \quad &> \varepsilon. \end{aligned}$$

1128 Therefore, $c \leq \varepsilon$ is sufficient for certifying the desired inequality. The latter claim follows by squaring and rearranging
 1129 terms to yield the quadratic inequality:

1130

$$(1 - \varepsilon)^2 d - 2(1 - \varepsilon) + 1 \leq 0,$$

1131

1132 Setting $\lambda := 1 - \varepsilon$, the solution interval is $\lambda \in \left(\frac{1 - \sqrt{1-d}}{d}, \frac{1 + \sqrt{1-d}}{d}\right)$. The upper limit is redundant as it exceeds 1 and
 1133 $\varepsilon \in (0, 1)$, leaving the lower limit as the condition on ε proposed in the lemma.

1134

1135 Plugging in $\eta_G = d \in (0, 1]$ and $2(\gamma_1 + 2\gamma_2 + \eta_G(\gamma_1 + \gamma_2)\gamma_2) \leq (\gamma_1 + 2\gamma_2 + (\gamma_1 + \gamma_2)\gamma_2 = c$, we try candidate values
 1136 $\gamma_1 = 1/40, \gamma_2 = 1/100$ and set $\varepsilon = c$ to get:

1137

$$\frac{1 - \eta_G}{\sqrt{1 - 2\eta_G(\gamma_1 + 2\gamma_2 + \eta_G(\gamma_1 + \gamma_2)\gamma_2)}} < (1 - 0.9\eta_G).$$

1138

1139 Plugging in our candidate values of γ_1, γ_2 into the burn-in conditions of Proposition B.7 finishes the proof of Theorem 3.6.

1140

1141

1155 **Theorem 3.6.** Consider running (8) with $\lambda_G = 0$, $\eta_G \in [0, 1]$, and $\eta_F = 1$. Define $\bar{\sigma}^2 \triangleq \max\{1, \frac{\sigma_\epsilon^2}{\sigma_{\min}(\mathbf{F}_*)^2 \lambda_{\min}(\Sigma_x)}\}$. As
1156 long as $\text{dist}(\mathbf{G}, \mathbf{G}_*) \leq 0.01\kappa^{-1}(\Sigma_x)\kappa^{-1}(\mathbf{F}_*)$ and $n \gtrsim \bar{\sigma}^2(d_x + \log(1/\delta))$, we have:
1157

$$1158 \quad \text{dist}(\mathbf{G}_+, \mathbf{G}_*) \leq (1 - 0.9\eta_G)\text{dist}(\mathbf{G}, \mathbf{G}_*) + \eta_G\Delta,$$

1159
1160 with probability $\geq 1 - \delta$, where $\Delta \triangleq \mathcal{O}\left(\bar{\sigma}\sqrt{\frac{d_x + \log(1/\delta)}{n}}\right)$
1161

1162 B.3. Multi-Task and Transfer Learning

1163 We first discuss how the ideas in our “single-task” setting directly translate to multi-task learning. For example, taking our
1164 proposed algorithm template in (8), an immediate idea is, given the current task heads and shared representation $(\{\mathbf{F}^{(t)}\}, \mathbf{G})$,
1165 to form task-specific preconditioners formed locally on each task’s batch data:
1166

$$1168 \quad \mathbf{P}_F^{(t)} = \widehat{\mathbf{E}}^{(t)}[\mathbf{z}\mathbf{z}^\top], \quad \mathbf{P}_G^{(t)} = \mathbf{F}^{(t)^\top} \mathbf{F}^{(t)}, \quad \mathbf{Q}_G^{(t)} = \widehat{\mathbf{E}}^{(t)}[\mathbf{x}\mathbf{x}^\top],$$

1169 and perform a local update on $\mathbf{F}^{(t)}$, \mathbf{G} before a central agent averages the resulting updated \mathbf{G} :
1170

$$\begin{aligned} 1172 \quad \mathbf{F}_+^{(t)} &= \mathbf{F}^{(t)} - \eta_F \nabla_{\mathbf{F}} \widehat{\mathcal{L}}(\mathbf{F}^{(t)}, \mathbf{G}^{(t)}) \mathbf{Q}_F^{(t)^{-1}} \\ 1173 \quad \mathbf{G}_+^{(t)} &= \mathbf{G} - \eta_G \mathbf{P}_G^{(t)^{-1}} \nabla_{\mathbf{G}} \widehat{\mathcal{L}}^{(t)}(\mathbf{F}_+^{(t)}, \mathbf{G}) \mathbf{Q}_G^{(t)^{-1}}, \quad t \in [T] \\ 1175 \quad \mathbf{G}_+ &= \frac{1}{T} \sum_{t=1}^T \mathbf{G}_+^{(t)}. \end{aligned}$$

1178 However, this presumes $\mathbf{F}^{(t)}$ are invertible, i.e. the task-specific dimension $d_y > k$. As opposed to the single-task setting,
1179 where as stated Remark 3.4 we are really viewing \mathbf{F} as the concatenation of $\mathbf{F}^{(t)}$ to make recovering the representation
1180 a well-posed problem, in multi-task settings d_y may often be small, e.g. $d_y = 1$ (Tripuraneni et al., 2021b; Du et al.,
1181 2021; Collins et al., 2021; Thekumparampil et al., 2021). Therefore, (pseudo)-inverting away $(\mathbf{F}^{(t)^\top} \mathbf{F}^{(t)})^\dagger$ may be highly
1182 suboptimal. However, we observe that writing out the representation gradient (9), as long as we invert away $\mathbf{Q}_G^{(t)^{-1}}$ first,
1183 then we have:
1184

$$\begin{aligned} 1186 \quad \mathbf{G}_+^{(t)} &= \mathbf{G} - \eta_G \nabla_{\mathbf{G}} \widehat{\mathcal{L}}^{(t)}(\mathbf{F}_+^{(t)}, \mathbf{G}) \mathbf{Q}_G^{(t)^{-1}} \\ 1188 \quad \mathbf{G}_+ \mathcal{P}_*^\perp &= \frac{1}{T} \sum_{t=1}^T \mathbf{G}_+^{(t)} \mathcal{P}_*^\perp \\ 1190 \quad &= \left(\mathbf{I}_k - \eta_G \frac{1}{T} \sum_{t=1}^T \mathbf{F}^{(t)^\top} \mathbf{F}^{(t)} \right) \mathbf{G} \mathcal{P}_*^\perp + (\text{task-averaged}) \text{ noise term}. \end{aligned}$$

1194 Since by assumption $\frac{1}{T} \sum_{t=1}^T \mathbf{F}^{(t)^\top} \mathbf{F}^{(t)}$ is full-rank (otherwise recovering the rank k representation is impossible), then
1195 suggestively, we may instead invert away the *task-averaged* preconditioner $\mathbf{P}_G = \frac{1}{T} \sum_{t=1}^T \mathbf{F}^{(t)^\top} \mathbf{F}^{(t)}$ on the task-averaged
1196 $\mathbf{G}^{(t)}$ descent direction before taking a representation step \mathbf{G}_+ . To summarize, we propose the following two-stage
1197 preconditioning:
1198

$$1199 \quad \mathbf{F}_+^{(t)} = \mathbf{F}^{(t)} - \eta_F \nabla_{\mathbf{F}} \widehat{\mathcal{L}}(\mathbf{F}^{(t)}, \mathbf{G}^{(t)}) \mathbf{Q}_F^{(t)^{-1}} \tag{20}$$

$$1201 \quad \mathbf{D}^{(t)} = \nabla_{\mathbf{G}} \widehat{\mathcal{L}}^{(t)}(\mathbf{F}_+^{(t)}, \mathbf{G}) \mathbf{Q}_G^{(t)^{-1}}, \quad t \in [T] \tag{21}$$

$$1203 \quad \mathbf{G}_+ = \mathbf{G} - \eta_G \mathbf{P}_G^{-1} \left(\frac{1}{T} \sum_{t=1}^T \mathbf{D}^{(t)} \right) \tag{22}$$

$$1206 \quad \text{such that } \mathbf{G}_+ \mathcal{P}_*^\perp = (1 - \eta_G) \mathbf{G} \mathcal{P}_*^\perp + (\text{task-averaged}) \text{ noise term}. \tag{23}$$

1207 The exact same tools used in the proof of Theorem 3.6 apply here, with the requirement of a few additional standard tools to
1208 study the “task-averaged” noise term(s). As an example, we refer to Zhang et al. (2024b) for some candidates. However,
1209

we note the qualitative behavior is unchanged. As such, since we are using n data points per each of T tasks to update the gradient, the scaling of the noise term goes from $\mathcal{O}(\sigma_\epsilon \sqrt{d_x/N})$ in our bounds to $\mathcal{O}(\sigma_\epsilon \sqrt{d_x/NT})$.

We remark that in the multi-task setting, where each task may have differing covariances and task-heads $\mathbf{F}_\star^{(t)}$, the equivalence of our stylized KFAC variant and the alternating min-min algorithm proposed in Jain et al. (2013); Thekumparampil et al. (2021) breaks down. In particular, the alternating min-min algorithm no longer in general admits \mathbf{G} iterates that can be expressed as a product of matrices as in (5) or (20), and rather can only be stated in vectorized space $\text{vec}(\mathbf{G})$. This means that whereas (20) can be solved as T parallel small matrix multiplication problems, the alternating min-min algorithm nominally requires operating in the vectorized-space $d_x k$.

Transfer Learning

We first prove the proposed fine-tuning generalization bound.

Lemma 3.7. *Let $\widehat{\mathbf{F}}_{ls}^{(t)} = \underset{\widehat{\mathbf{F}}}{\operatorname{argmin}} \widehat{\mathbb{E}}^{(t)} [\|\mathbf{y}^{(t)} - \widehat{\mathbf{F}} \mathbf{z}^{(t)}\|_2^2]$, $\mathbf{z}^{(t)} \triangleq \widehat{\mathbf{G}} \mathbf{x}^{(t)}$ be the optimal \mathbf{F} on the batch of $n^{(t)}$ target data (11) given $\widehat{\mathbf{G}}$. Defining $\nu = \text{dist}(\widehat{\mathbf{G}}, \mathbf{G}_\star)$, given $n^{(t)} \gtrsim k + \log(1/\delta)$, we have with probability $\geq 1 - \delta$:*

$$\begin{aligned} \mathcal{L}^{(t)}(\widehat{\mathbf{F}}_{ls}^{(t)}, \widehat{\mathbf{G}}) &\triangleq \mathbb{E} \left[\|\mathbf{y}^{(t)} - \mathbf{F}_\star^{(t)} \mathbf{G}_\star \mathbf{x}^{(t)}\|_2^2 \right] \\ &\lesssim \|\mathbf{F}_\star^{(t)}\|_F^2 \lambda_{\max}(\Sigma_x^{(t)}) \nu^2 + \frac{\sigma_\epsilon (d_y k + \log(1/\delta))}{n^{(t)}}. \end{aligned}$$

Proof of Lemma 3.7. We observe that we may write:

$$\begin{aligned} \mathcal{L}^{(t)}(\widehat{\mathbf{F}}_{ls}^{(t)}, \widehat{\mathbf{G}}) &= \mathbb{E} \left[\|\mathbf{y}^{(t)} - \widehat{\mathbf{F}}_{ls}^{(t)} \widehat{\mathbf{G}} \mathbf{x}^{(t)}\|_2^2 \right] \\ &= \mathbb{E} \left[\|(\widehat{\mathbf{F}}_{ls}^{(t)} \widehat{\mathbf{G}} - \mathbf{F}_\star^{(t)} \mathbf{G}_\star) \mathbf{x}^{(t)}\|_2^2 \right] \\ &= \|(\widehat{\mathbf{F}}_{ls}^{(t)} \widehat{\mathbf{G}} - \mathbf{F}_\star^{(t)} \mathbf{G}_\star) \Sigma_x^{(t)}{}^{1/2}\|_F^2. \end{aligned}$$

Now writing out the definition of $\widehat{\mathbf{F}}_{ls}^{(t)}$, defining $\mathbf{z}^{(t)} = \widehat{\mathbf{G}} \mathbf{x}^{(t)}$, we have

$$\begin{aligned} \widehat{\mathbf{F}}_{ls}^{(t)} &= \underset{\widehat{\mathbf{F}}}{\operatorname{argmin}} \widehat{\mathbb{E}}^{(t)} [\|\mathbf{y}^{(t)} - \widehat{\mathbf{F}} \widehat{\mathbf{G}} \mathbf{x}^{(t)}\|_2^2] \\ &= \mathbf{Y}^{(t)\top} \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1} \\ &= \mathbf{F}_\star^{(t)} \mathbf{G}_\star \mathbf{X}^{(t)\top} \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1} + \mathcal{E}^\top \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1} \\ &= \mathbf{F}_\star^{(t)} \mathbf{G}_\star \widehat{\mathbf{G}}^\top + \mathbf{F}_\star^{(t)} \mathbf{G}_\star \mathcal{P}_{\widehat{\mathbf{G}}}^\perp \mathbf{X}^{(t)\top} \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1} + \mathcal{E}^\top \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1}, \end{aligned}$$

where $\mathcal{P}_{\widehat{\mathbf{G}}}^\perp = \mathbf{I}_{d_x} - \widehat{\mathbf{G}}^\top \widehat{\mathbf{G}}$ is the projection matrix onto the rowspace of $\widehat{\mathbf{G}}$, using the fact that $\widehat{\mathbf{G}}$ is row-orthonormal (8). Therefore, plugging in the last line into error expression, we have

$$\begin{aligned} \mathcal{L}^{(t)}(\widehat{\mathbf{F}}_{ls}^{(t)}, \widehat{\mathbf{G}}) &= \|(\widehat{\mathbf{F}}_{ls}^{(t)} \widehat{\mathbf{G}} - \mathbf{F}_\star^{(t)} \mathbf{G}_\star) \Sigma_x^{(t)}{}^{1/2}\|_F^2 \\ &\leq 2 \left\| \left(\mathbf{F}_\star^{(t)} \mathbf{G}_\star \widehat{\mathbf{G}}^\top + \mathbf{F}_\star^{(t)} \mathbf{G}_\star \mathcal{P}_{\widehat{\mathbf{G}}}^\perp \mathbf{X}^{(t)\top} \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1} \widehat{\mathbf{G}} - \mathbf{F}_\star^{(t)} \mathbf{G}_\star \right) \Sigma_x^{(t)}{}^{1/2} \right\|_F^2 \\ &\quad + 2 \left\| \mathcal{E}^\top \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1} \widehat{\mathbf{G}} \Sigma_x^{(t)}{}^{1/2} \right\|_F^2. \end{aligned} \quad ((a+b)^2 \leq 2a^2 + 2b^2)$$

Focusing on the first term, we have:

$$\begin{aligned} &\mathbf{F}_\star^{(t)} \mathbf{G}_\star \widehat{\mathbf{G}}^\top \widehat{\mathbf{G}} + \mathbf{F}_\star^{(t)} \mathbf{G}_\star \mathcal{P}_{\widehat{\mathbf{G}}}^\perp \mathbf{X}^{(t)\top} \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1} \widehat{\mathbf{G}} - \mathbf{F}_\star^{(t)} \mathbf{G}_\star \\ &= \mathbf{F}_\star^{(t)} \mathbf{G}_\star \mathcal{P}_{\widehat{\mathbf{G}}}^\perp \left(\mathbf{X}^{(t)\top} \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1} \widehat{\mathbf{G}} - \mathbf{I}_{d_x} \right). \end{aligned}$$

By a covariance concentration argument Lemma E.2, since $\mathbf{X}^{(t)\top} \mathbf{Z}^{(t)}$ and $\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)}$ are rank- k matrices, as long as $n^{(t)} \gtrsim k + \log(1/\delta)$, we have with probability at least $1 - \delta$:

$$\mathbf{X}^{(t)\top} \mathbf{Z}^{(t)} \approx n^{(t)} \Sigma_x^{(t)} \widehat{\mathbf{G}}^\top, \quad (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1} \approx n^{(t)} \widehat{\mathbf{G}} \Sigma_x^{(t)} \widehat{\mathbf{G}}^\top,$$

1265 and thus

$$\begin{aligned}
& \|\mathbf{F}_*^{(t)} \mathbf{G}_* \mathcal{P}_{\widehat{\mathbf{G}}}^\perp \left(\mathbf{X}^{(t)\top} \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1} \widehat{\mathbf{G}} - \mathbf{I}_{d_x} \right) \Sigma_{\mathbf{x}}^{(t)1/2} \|_F \\
& \approx \|\mathbf{F}_*^{(t)} \mathbf{G}_* \mathcal{P}_{\widehat{\mathbf{G}}}^\perp \Sigma_{\mathbf{x}}^{(t)1/2} \left(\Sigma_{\mathbf{x}}^{(t)1/2} \widehat{\mathbf{G}}^\top (\widehat{\mathbf{G}} \Sigma_{\mathbf{x}}^{(t)} \widehat{\mathbf{G}}^\top)^{-1} \widehat{\mathbf{G}} \Sigma_{\mathbf{x}}^{(t)1/2} - \mathbf{I}_{d_x} \right) \|_F \\
& \lesssim \|\mathbf{F}_*^{(t)}\|_F \left\| \mathbf{G}_* \mathcal{P}_{\widehat{\mathbf{G}}}^\perp \right\|_{\text{op}} \left\| \Sigma_{\mathbf{x}}^{(t)1/2} \right\|_{\text{op}} \left\| \Sigma_{\mathbf{x}}^{(t)1/2} \widehat{\mathbf{G}}^\top (\widehat{\mathbf{G}} \Sigma_{\mathbf{x}}^{(t)} \widehat{\mathbf{G}}^\top)^{-1} \widehat{\mathbf{G}} \Sigma_{\mathbf{x}}^{(t)1/2} - \mathbf{I}_{d_x} \right\|_{\text{op}} \\
& \leq \|\mathbf{F}_*^{(t)}\|_F \text{dist}(\widehat{\mathbf{G}}, \mathbf{G}_*) \lambda_{\max}(\Sigma_{\mathbf{x}}^{(t)})^{1/2},
\end{aligned}$$

1275 where in the last line we applied the definition $\text{dist}(\widehat{\mathbf{G}}, \mathbf{G}_*) = \left\| \mathbf{G}_* \mathcal{P}_{\widehat{\mathbf{G}}}^\perp \right\|_{\text{op}} = \left\| \widehat{\mathbf{G}} \mathcal{P}_{\mathbf{G}_*}^\perp \right\|_{\text{op}}$, and the fact that the matrix
1276
1277 $\mathcal{P} \triangleq \Sigma_{\mathbf{x}}^{(t)1/2} \widehat{\mathbf{G}}^\top (\widehat{\mathbf{G}} \Sigma_{\mathbf{x}}^{(t)} \widehat{\mathbf{G}}^\top)^{-1} \widehat{\mathbf{G}} \Sigma_{\mathbf{x}}^{(t)1/2}$ can be verified to be a projection matrix $\mathcal{P}^2 = \mathcal{P}$, $\mathcal{P}^\top = \mathcal{P}$, such that
1278 $\mathcal{P} - \mathbf{I} = \mathcal{P}^\perp$ is also an orthogonal projection and $\|\mathcal{P}^\perp\|_{\text{op}} = 1$. Now, we analyze the noise term:
1279

$$\mathcal{E}^\top \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1} \widehat{\mathbf{G}} \Sigma_{\mathbf{x}}^{(t)1/2} \approx \mathcal{E}^\top \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1/2} (n^{(t)})^{-1/2} \left(\widehat{\mathbf{G}} \Sigma_{\mathbf{x}}^{(t)} \widehat{\mathbf{G}}^\top \right)^{-1/2} \widehat{\mathbf{G}} \Sigma_{\mathbf{x}}^{(t)1/2},$$

1283 where we observed $\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)} = n \mathbf{G} \widehat{\Sigma}_{\mathbf{x}} \mathbf{G}^\top$ and applied covariance concentration. Now, defining the (compact) SVD of
1284 $\widehat{\mathbf{G}} \Sigma_{\mathbf{x}}^{(t)1/2} = \mathbf{U}_{\mathbf{z}} \mathbf{D}_{\mathbf{z}} \mathbf{V}_{\mathbf{z}}^\top$, we find
1285

$$\begin{aligned}
\|\mathcal{E}^\top \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1} \widehat{\mathbf{G}} \Sigma_{\mathbf{x}}^{(t)1/2}\|_F & \lesssim \frac{1}{\sqrt{n^{(t)}}} \|\mathcal{E}^\top \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1/2} \mathbf{U}_{\mathbf{z}} \mathbf{V}_{\mathbf{z}}^\top\|_F \\
& \lesssim \frac{1}{\sqrt{n^{(t)}}} \|\mathcal{E}^\top \mathbf{Z}^{(t)} (\mathbf{Z}^{(t)\top} \mathbf{Z}^{(t)})^{-1/2}\|_F \\
& \lesssim \frac{1}{\sqrt{n^{(t)}}} \sigma_{\mathcal{E}} \sqrt{d_y k + \log(1/\delta)},
\end{aligned}$$

1294 for $n^{(t)} \gtrsim k + \log(1/\delta)$. The last line comes from the Frobenius norm variants of Lemma B.3 and Lemma B.4 (see Ziermann
1295 et al. (2023, Theorem 4.1) or Zhang et al. (2024b, Lemma A.3) for details). Putting the two bounds together yields the
1296 desired result. \square

C. Proofs and Additional Details for Section 3.2

C.1. Proof of Theorem 3.9

1300 **Theorem 3.9.** Assume that the activation function σ is $\mathcal{O}(1)$ -Lipschitz and that Assumption 3.8 holds. In the limit where
1301 n, d_x, d_h tend to infinity proportionally, the matrix \mathbf{G}_{SGD} , with probability $1 - o(1)$, satisfies

$$\|\mathbf{G}_0 + \alpha \eta \mathbf{f}_0 \boldsymbol{\beta}_{\text{SGD}}^\top - \mathbf{G}_{\text{SGD}}\|_{\text{op}} \rightarrow 0,$$

1307 in which $\alpha = \mathbb{E}_z[\sigma'(z)]$ with $z \sim N(0, d_x^{-1} \text{Tr}(\Sigma_{\mathbf{x}}))$, and the vector $\boldsymbol{\beta}_{\text{SGD}}$ is given by $\boldsymbol{\beta}_{\text{SGD}} = n^{-1} \mathbf{X}^\top \mathbf{y}$.

1310 *Proof.* To prove this theorem, we first note that

$$\nabla_{\mathbf{G}} \widehat{\mathcal{L}}(\mathbf{f}_0, \mathbf{G}_0) = -\frac{1}{n} \sum_{i=1}^n \left(y_i - \frac{1}{\sqrt{d_h}} \mathbf{f}_0^\top \sigma(\mathbf{G}_0 \mathbf{x}_i) \right) \left(\frac{1}{\sqrt{d_h}} \mathbf{f}_0 \odot \sigma'(\mathbf{G}_0 \mathbf{x}_i) \right) \mathbf{x}_i^\top.$$

1315 Adopting the matrix notation $\mathbf{X} = [\mathbf{x}_1 | \dots | \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d_x}$ and $\mathbf{y} = [y_1, \dots, y_n]^\top \in \mathbb{R}^n$, we can write
1316

$$\nabla_{\mathbf{G}} \widehat{\mathcal{L}}(\mathbf{f}_0, \mathbf{G}_0) = -\frac{1}{n} \left[\left(d_h^{-1/2} \mathbf{f}_0 \mathbf{y}^\top - d_h^{-1} \mathbf{f}_0 \mathbf{f}_0^\top \sigma(\mathbf{G}_0 \mathbf{X}^\top) \right) \odot \sigma'(\mathbf{G}_0 \mathbf{X}^\top) \right] \mathbf{X}. \quad (24)$$

Let $z \sim N(0, d_x^{-1} \text{Tr}(\Sigma_x))$ and define $\alpha = E_z[\sigma'(z)]$, and $\sigma_\perp : \mathbb{R} \rightarrow \mathbb{R}$ as $\sigma_\perp(z) = \sigma(z) - \alpha z$. This function satisfies $E_z[\sigma'_\perp(z)] = 0$. With this, we decompose the gradient into three components as $\nabla_{\mathbf{G}} \hat{\mathcal{L}}(\mathbf{f}_0, \mathbf{G}_0) = \mathbf{T}_1 + \mathbf{T}_2 + \mathbf{T}_3$ with

$$\begin{aligned}\mathbf{T}_1 &= -\alpha d_h^{-1/2} \mathbf{f}_0 \left(\frac{\mathbf{X}^\top \mathbf{y}}{n} \right)^\top, \quad \mathbf{T}_2 = -n^{-1} d_h^{-1/2} [\mathbf{f}_0 \mathbf{y}^\top \odot \sigma'_\perp(\mathbf{G}_0 \mathbf{X}^\top)] \mathbf{X}, \\ \mathbf{T}_3 &= n^{-1} d_h^{-1} [(\mathbf{f}_0 \mathbf{f}_0^\top \sigma(\mathbf{G}_0 \mathbf{X}^\top)) \odot \sigma'(\mathbf{G}_0 \mathbf{X}^\top)] \mathbf{X}.\end{aligned}$$

We will analyze each of these components separately.

- **Term 1:** For this term, using the facts that $\|\mathbf{f}_0\|_2 = \mathcal{O}(1)$ and $\|\mathbf{X}^\top \mathbf{y}/n\|_2 = \mathcal{O}(1)$, we have

$$\|\mathbf{T}_1\|_{\text{op}} = \mathcal{O}(d_h^{-1/2}).$$

- **Term 2:** To analyze this term, note that

$$\mathbf{f}_0 \mathbf{y}^\top \odot \sigma'_\perp(\mathbf{G}_0 \mathbf{X}^\top) = \text{diag}(\mathbf{f}_0) \sigma'_\perp(\mathbf{G}_0 \mathbf{X}^\top) \text{diag}(\mathbf{y}),$$

which gives

$$\|\mathbf{f}_0 \mathbf{y}^\top \odot \sigma'_\perp(\mathbf{G}_0 \mathbf{X}^\top)\|_{\text{op}} \leq \|\mathbf{f}_0\|_\infty \|\mathbf{y}\|_\infty \|\sigma'_\perp(\mathbf{G}_0 \mathbf{X}^\top)\|_{\text{op}}.$$

Using basic concentration arguments, we have $\|\mathbf{f}_0\|_\infty = \tilde{\mathcal{O}}(d_h^{-1/2})$, and $\|\mathbf{y}\|_\infty = \tilde{\mathcal{O}}(1)$, with probability $1 - o(1)$. By construction of $\sigma_\perp(\cdot)$, the matrix $\sigma'_\perp(\mathbf{G}_0 \mathbf{X}^\top)$ has mean zero entries, thus using (Vershynin, 2012, Theorem 5.44), we have $\|\sigma'_\perp(\mathbf{G}_0 \mathbf{X}^\top)\|_{\text{op}} = \tilde{\mathcal{O}}(d_h^{1/2} + n^{1/2})$ with probability $1 - o(1)$. Thus, the norm of \mathbf{T}_2 can be upper bounded as

$$\|\mathbf{T}_2\|_{\text{op}} = \tilde{\mathcal{O}} \left(\frac{1}{d_h} \left(1 + \sqrt{\frac{d_h}{n}} \right) \left(1 + \sqrt{\frac{d_x}{n}} \right) \right) = \tilde{\mathcal{O}}(d_h^{-1}).$$

- **Term 3:** Similar to the second term, note that

$$(\mathbf{f}_0 \mathbf{f}_0^\top \sigma(\mathbf{G}_0 \mathbf{X}^\top)) \odot \sigma'(\mathbf{G}_0 \mathbf{X}^\top) = \text{diag}(\mathbf{f}_0) \sigma'(\mathbf{G}_0 \mathbf{X}^\top) \text{diag}(\mathbf{f}_0^\top \sigma(\mathbf{G}_0 \mathbf{X}^\top)).$$

Thus, the norm of the third term can be upper bounded as

$$\|\mathbf{T}_3\|_{\text{op}} = \left\| \frac{1}{n d_h} [(\mathbf{f}_0 \mathbf{f}_0^\top \sigma(\mathbf{G}_0 \mathbf{X}^\top)) \odot \sigma'(\mathbf{G}_0 \mathbf{X}^\top)] \mathbf{X} \right\|_{\text{op}} \leq n^{-1} d_h^{-1} \|\mathbf{X}\|_{\text{op}} \|\mathbf{f}_0\|_\infty \|\mathbf{f}_0^\top \sigma(\mathbf{G}_0 \mathbf{X}^\top)\|_\infty \|\sigma'(\mathbf{G}_0 \mathbf{X}^\top)\|_{\text{op}}.$$

To analyze the right hand side, note that assuming that σ is $\mathcal{O}(1)$ -Lipschitz, the entries of $\sigma'(\mathbf{G}_0 \mathbf{X}^\top)$ are bounded by the Lipschitz constant, and we have $\|\sigma'(\mathbf{G}_0 \mathbf{X}^\top)\|_{\text{op}} = \mathcal{O}(\sqrt{n d_h})$. Also, using a simple orderwise analysis we have $\|\mathbf{f}_0^\top \sigma(\mathbf{G}_0 \mathbf{X}^\top)\|_\infty = \tilde{\mathcal{O}}(1)$, which gives

$$\|\mathbf{T}_3\|_{\text{op}} = \tilde{\mathcal{O}} \left(\frac{1}{d_h} \left(1 + \sqrt{\frac{d_x}{n}} \right) \right) = \tilde{\mathcal{O}}(d_h^{-1}).$$

To wrap up, note that $d_h^{1/2} \|\mathbf{T}_1\|_{\text{op}} = \mathcal{O}(1)$, whereas $d_h^{1/2} \|\mathbf{T}_2\|_{\text{op}}$ and $d_h^{1/2} \|\mathbf{T}_3\|_{\text{op}} = o(1)$. Thus, with probability $1 - o(1)$ we have

$$\mathbf{G}_{\text{SGD}} = \mathbf{G}_0 + \eta d_h^{1/2} \nabla_{\mathbf{G}} \hat{\mathcal{L}}(\mathbf{f}_0, \mathbf{G}_0) = \mathbf{G}_0 + \alpha \eta \mathbf{f}_0 (n^{-1} \mathbf{X}^\top \mathbf{y})^\top + \Delta$$

with $\|\Delta\|_{\text{op}} = o(1)$, finishing the proof. \square

1375 **C.2. Proof of Lemma 3.10**

1376 **Lemma 3.10.** Under the assumptions of Theorem 3.9, the correlation between β_* and β_{SGD} satisfies

$$1379 \quad \left| \frac{\beta_*^\top \beta_{\text{SGD}}}{\|\beta_{\text{SGD}}\|_2 \|\beta_*\|_2} - \frac{\frac{c_{*,1}}{d_x} \text{Tr}(\Sigma_x)}{\sqrt{\frac{c_*^2 + \sigma_\varepsilon^2}{n} \text{Tr}(\Sigma_x) + \frac{c_{*,1}^2}{d_x} \text{Tr}(\Sigma_x^2)}} \right| \rightarrow 0$$

1382 with probability $1 - o(1)$, in which $c_{*,1} = E_z[\sigma'_*(z)]$ and $c_*^2 = E_z[\sigma_*^2(z)]$ with $z \sim N(0, d_x^{-1} \text{Tr}(\Sigma_x))$.

1384 *Proof.* Recall that $\beta_{\text{SGD}} = \frac{1}{n} \mathbf{X}^\top \mathbf{y}$ and $\mathbf{y} = \sigma_*(\mathbf{X}\beta_*) + \varepsilon$ where $\varepsilon = [\varepsilon_1, \dots, \varepsilon_n]^\top$. Therefore, with probability $1 - o(1)$
1386 we have

$$1388 \quad \beta_{\text{SGD}}^\top \beta_* = \frac{1}{n} (\mathbf{X}\beta_*)^\top \mathbf{y} = \frac{1}{n} (\mathbf{X}\beta_*)^\top \sigma_*(\mathbf{X}\beta_*) + o(1)$$

1390 where we have used the fact that ε is mean zero. Thus, using the weak law of large numbers,

$$1392 \quad \beta_{\text{SGD}}^\top \beta_* \rightarrow E_z[z\sigma_*(z)] = d_x^{-1} \text{Tr}(\Sigma_x) E_z[\sigma'_*(z)] = c_{*,1} d_x^{-1} \text{Tr}(\Sigma_x) \quad (25)$$

1394 in probability, where $z \sim N(0, d_x^{-1} \text{Tr}(\Sigma_x))$. Similarly, $\|\beta_{\text{SGD}}\|_2^2$ can be written as

$$1396 \quad \|\beta_{\text{SGD}}\|_2^2 = n^{-2} \mathbf{y}^\top \mathbf{X} \mathbf{X}^\top \mathbf{y} = n^{-2} \varepsilon^\top \mathbf{X} \mathbf{X}^\top \varepsilon + n^{-2} \sigma_*(\mathbf{X}\beta_*)^\top \mathbf{X} \mathbf{X}^\top \sigma_*(\mathbf{X}\beta_*) + o(1).$$

1397 We will analyze each of the two remaining term separately. For the first term, recall that ε is independent of \mathbf{X} . Using the
1398 Hanson-Wright inequality (Theorem E.6) we have

$$1400 \quad n^{-2} \varepsilon^\top \mathbf{X} \mathbf{X}^\top \varepsilon = \sigma_\varepsilon^2 n^{-1} \text{Tr}(\mathbf{X} \mathbf{X}^\top / n) + o(1) = \sigma_\varepsilon^2 n^{-1} \text{Tr}(\Sigma_x) + o(1).$$

1402 For the second term, note that $\mathbf{X}\beta_*$ is a vector with i.i.d. elements $x_i^\top \beta_*$, each of them distributed according to
1403 $N(0, \beta_*^\top \Sigma_x \beta_*)$. Let z be a random variable distributed as $z \sim N(0, \beta_*^\top \Sigma_x \beta_*)$. We decompose the function σ_* into
1404 a linear and a nonlinear part as

$$1406 \quad \sigma_*(z) = c_{*,1} z + \sigma_{*,\perp}(z). \quad (26)$$

1408 This decomposition satisfies

$$1410 \quad E_z[\sigma_{*,\perp}(z)] = E_z[\sigma_*(z)] = 0$$

$$1411 \quad E_z[z\sigma_{*,\perp}(z)] = E_z[z\sigma_*(z)] - c_{*,1} E_z[z^2] = E_z[z\sigma_*(z)] - c_{*,1} \beta_*^\top \Sigma_x \beta_* = 0,$$

1413 where the last equality is due to Stein's lemma (Lemma E.7). This shows that the random variables z and $\sigma_{*,\perp}(z)$ are
1414 uncorrelated. With this, we have

$$1416 \quad n^{-2} \sigma_*(\mathbf{X}\beta_*)^\top \mathbf{X} \mathbf{X}^\top \sigma_*(\mathbf{X}\beta_*) = n^{-2} (c_{*,1} \mathbf{X}\beta_* + \sigma_{*,\perp}(\mathbf{X}\beta_*))^\top \mathbf{X} \mathbf{X}^\top (c_{*,1} \mathbf{X}\beta_* + \sigma_{*,\perp}(\mathbf{X}\beta_*))$$

$$1418 \quad = c_{*,1}^2 n^{-2} \beta_*^\top (\mathbf{X}^\top \mathbf{X})^2 \beta_* + 2c_{*,1} n^{-2} (\mathbf{X}\beta_*)^\top \mathbf{X} \mathbf{X}^\top \sigma_{*,\perp}(\mathbf{X}\beta_*) + n^{-2} \sigma_{*,\perp}(\mathbf{X}\beta_*)^\top \mathbf{X} \mathbf{X}^\top \sigma_{*,\perp}(\mathbf{X}\beta_*). \quad (27)$$

1420 For the first term in this sum, by assumption 3.8 and the Hanson-Wright inequality (Theorem E.6), we can write

$$1422 \quad c_{*,1}^2 n^{-2} \beta_*^\top (\mathbf{X}^\top \mathbf{X})^2 \beta_* = c_{*,1}^2 d_x^{-1} \text{Tr}(n^{-2} (\mathbf{X}^\top \mathbf{X})^2) + o(1) = c_{*,1}^2 d_x^{-1} \text{Tr}(\Sigma_x^2) + c_{*,1}^2 n^{-1} d_x^{-1} \text{Tr}^2(\Sigma_x),$$

1424 where in the last we plugged in the second Wishart moment. For the second term in (27), although by construction $\mathbf{X}\beta_*$
1425 and $\sigma_{*,\perp}(\mathbf{X}\beta_*)$ are uncorrelated, the vector $\mathbf{X}\beta_*$ and the matrix $\mathbf{X} \mathbf{X}^\top$ are dependent, which complicates the analysis. To
1426 resolve this issue, we define $\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{X}\beta_* \beta_*^\top$ which satisfies $\mathbf{X}\beta_* \perp \tilde{\mathbf{X}}$ and write

$$1428 \quad \mathbf{X} \mathbf{X}^\top = \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top + \tilde{\mathbf{X}} \beta_* \beta_*^\top \mathbf{X}^\top + \mathbf{X} \beta_* \beta_*^\top \tilde{\mathbf{X}}^\top + \mathbf{X} \beta_* \beta_*^\top \mathbf{X}^\top.$$

1430 Thus, the second term in (27) can be written as

$$1432 n^{-2}(\mathbf{X}\beta_*)^\top \mathbf{X}\mathbf{X}^\top \sigma_{*,\perp}(\mathbf{X}\beta_*) = n^{-2}(\mathbf{X}\beta_*)^\top [\tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top + \tilde{\mathbf{X}}\beta_*\beta_*^\top \mathbf{X}^\top + \mathbf{X}\beta_*\beta_*^\top \tilde{\mathbf{X}}^\top + \mathbf{X}\beta_*\beta_*^\top \mathbf{X}^\top] \sigma_{*,\perp}(\mathbf{X}\beta_*) \\ 1433 = \underbrace{n^{-2}(\mathbf{X}\beta_*)^\top \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top \sigma_{*,\perp}(\mathbf{X}\beta_*)}_{T_1} + \underbrace{n^{-2}(\mathbf{X}\beta_*)^\top [\tilde{\mathbf{X}}\beta_*\beta_*^\top \mathbf{X}^\top + \mathbf{X}\beta_*\beta_*^\top \tilde{\mathbf{X}}^\top + \mathbf{X}\beta_*\beta_*^\top \mathbf{X}^\top] \sigma_{*,\perp}(\mathbf{X}\beta_*)}_{T_2}.$$

1437 The term T_1 can be shown to be $o(1)$ by using the Hanson-Wright inequality (Theorem E.6) and noting that $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top$ is
1438 independent of $\mathbf{X}\beta_*$, and also the fact that $\mathbf{X}\beta_*$ and $\sigma_{*,\perp}(\mathbf{X}\beta_*)$ are orthogonal, by construction. Similarly, T_2 can also be
1439 shown to be $o(1)$ using a similar argument. For this, we also use that fact that by construction we have $\tilde{\mathbf{X}}\beta_* \perp\!\!\!\perp \mathbf{X}\beta_*$ which
1440 alongside $E[\tilde{\mathbf{X}}\beta_*] = \mathbf{0}_n$ proves that $n^{-1}(\mathbf{X}\beta_*)^\top \tilde{\mathbf{X}}\beta_* = o(1)$ and $n^{-1}(\tilde{\mathbf{X}}\beta_*)^\top \sigma_{*,\perp}(\tilde{\mathbf{X}}\beta_*) = o(1)$. Hence,

$$1442 2c_{*,1}n^{-2}(\mathbf{X}\beta_*)^\top \mathbf{X}\mathbf{X}^\top \sigma_{*,\perp}(\mathbf{X}\beta_*) \rightarrow 0.$$

1444 For the third term in (27), we can use a similar argument and replace \mathbf{X} with $\tilde{\mathbf{X}}$ to show that

$$1446 n^{-2}\sigma_{*,\perp}(\mathbf{X}\beta_*)^\top \mathbf{X}\mathbf{X}^\top \sigma_{*,\perp}(\mathbf{X}\beta_*) \rightarrow E_z[\sigma_{*,\perp}(z)]^2 n^{-1} \text{Tr}(\Sigma_x).$$

1448 Putting everything together, we have

$$1450 \|\beta_{\text{SGD}}\|_2^2 = c_{*,1}^2 d_x^{-1} \text{Tr}(\Sigma_x^2) + \sigma_\varepsilon^2 n^{-1} \text{Tr}(\Sigma_x) + E_z[\sigma_{*,\perp}(z)]^2 n^{-1} \text{Tr}(\Sigma_x) + c_{*,1}^2 n^{-1} d_x^{-1} \text{Tr}^2(\Sigma_x) + o(1) \\ 1451 = c_{*,1}^2 d_x^{-1} \text{Tr}(\Sigma_x^2) + n^{-1} \text{Tr}(\Sigma_x) (\sigma_\varepsilon^2 + E_z[\sigma_{*,\perp}(z)]^2 + c_{*,1}^2 d_x^{-1} \text{Tr}(\Sigma_x)) + o(1).$$

1453 Note that given the decomposition (26), we have

$$1455 E_z[\sigma_{*,\perp}^2(z)] = E_z[\sigma_{*,\perp}^2(z)] + c_{*,1}^2 d_x^{-1} \text{Tr}(\Sigma_x)$$

1457 given the orthogonality of the linear and nonlinear terms. Hence,

$$1459 \|\beta_{\text{SGD}}\|_2^2 = c_{*,1}^2 d_x^{-1} \text{Tr}(\Sigma_x^2) + n^{-1} \text{Tr}(\Sigma_x) (\sigma_\varepsilon^2 + c_{*,1}^2) + o(1),$$

1461 which alongside (25) proves the lemma. \square

C.3. Proof of Lemma 3.12

1465 **Lemma 3.12.** *Under the assumptions of Theorem 3.9, the correlation between β_* and β_{KFAC} satisfies*

$$1467 \left| \frac{\beta_{\text{KFAC}}^\top \beta_*}{\|\beta_{\text{KFAC}}\|_2 \|\beta_*\|_2} - \frac{c_{*,1} \Psi_1}{\sqrt{c_{*,1}^2 \Psi_2 + \frac{d_x}{n} (c_{*,1}^2 + \sigma_\varepsilon^2) \Psi_3}} \right| \rightarrow 0$$

1471 with probability $1 - o(1)$, where $c_{*,1}^2 = E_z[\sigma_{*,\perp}^2(z)]$, $c_{*,1}^2 = E_z[\sigma_{*,\perp}^2(z)]$ with $z \sim N(0, d_x^{-1} \text{Tr}(\Sigma_x))$, and Ψ_1, Ψ_2, Ψ_3 are
1472 defined in (34) and depend on Σ_x , d_x/n , and λ_G . In particular, as $\lambda_G \rightarrow 0$ and $d_x/n \rightarrow 0$, we have

$$1474 \beta_{\text{KFAC}}^\top \beta_* / (\|\beta_{\text{KFAC}}\|_2 \|\beta_*\|_2) \rightarrow 1.$$

1476 *Proof.* Recall that $\mathbf{Q}_G = n^{-1}\mathbf{X}^\top \mathbf{X}$ and let $\mathbf{R} = (\mathbf{Q}_G + \lambda_G \mathbf{I}_{d_x})^{-1}$. The inner product of β_* and β_{KFAC} is given by

$$1478 \beta_*^\top \beta_{\text{KFAC}} = n^{-1} \beta_*^\top \mathbf{R} \mathbf{X}^\top \sigma_*(\mathbf{X}\beta_*) + o(1),$$

1480 where we have used the fact that ε is mean zero and is independent of all other randomness in the problem. Defining
1481 $\bar{\mathbf{R}} = (\mathbf{X}\mathbf{X}^\top/n + \lambda_G \mathbf{I}_n)^{-1}$, we can use the push-through identity to rewrite the inner product as

$$1483 \beta_*^\top \beta_{\text{KFAC}} = n^{-1} (\mathbf{X}\beta_*)^\top \bar{\mathbf{R}} \sigma_*(\mathbf{X}\beta_*) + o(1).$$

1485 Note that $\mathbf{X}\beta_*$ is a vector with i.i.d. elements $\mathbf{x}_i^\top \beta_*$, each of them distributed according to $N(0, \beta_*^\top \Sigma_{\mathbf{x}} \beta_*)$. Using the same
 1486 decomposition for σ_* as the one used in the proof of Lemma 3.10 in (26), we have
 1487

$$\begin{aligned} 1488 \quad \beta_*^\top \beta_{\text{KFAC}} &= \frac{1}{n} (\mathbf{X}\beta_*)^\top \bar{\mathbf{R}} (c_{*,1} \mathbf{X}\beta_* + \sigma_{*,\perp}(\mathbf{X}\beta_*)) + o(1) \\ 1489 \\ 1490 &= c_{*,1} d_{\mathbf{x}}^{-1} \operatorname{Tr} (\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda_G n \mathbf{I}_n)^{-1} \mathbf{X}) + n^{-1} (\mathbf{X}\beta_*)^\top \bar{\mathbf{R}} \sigma_{*,\perp}(\mathbf{X}\beta_*) + o(1), \end{aligned} \quad (28)$$

1492 where for the first term we have used Assumption 3.8 and the Hanson-Wright inequality (Theorem E.6). To analyze
 1493 the second term, note that although by construction $\mathbf{X}\beta_*$ and $\sigma_{*,\perp}(\mathbf{X}\beta_*)$ are uncorrelated, the vectors $\mathbf{X}\beta_*$ and $\bar{\mathbf{R}}$ are
 1494 dependent, which complicates the analysis. To resolve this issue, we use the same trick used in the proof of Lemma 3.10 and
 1495 set $\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{X}\beta_* \beta_*^\top$ which satisfies $\mathbf{X}\beta_* \perp \tilde{\mathbf{X}}$, and use it to write
 1496

$$\begin{aligned} 1497 \quad \bar{\mathbf{R}}^{-1} &= n^{-1} \mathbf{X}\mathbf{X}^\top + \lambda_G \mathbf{I}_n = n^{-1} (\tilde{\mathbf{X}} + \mathbf{X}\beta_* \beta_*^\top) (\tilde{\mathbf{X}} + \mathbf{X}\beta_* \beta_*^\top)^\top + \lambda_G \mathbf{I}_n \\ 1498 \\ 1499 &= (n^{-1} \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top + \lambda_G \mathbf{I}_n) + n^{-1} (\mathbf{X}\beta_*)(\mathbf{X}\beta_*)^\top + n^{-1} (\tilde{\mathbf{X}} \beta_*)(\mathbf{X}\beta_*)^\top + n^{-1} (\mathbf{X}\beta_*)(\tilde{\mathbf{X}} \beta_*)^\top. \end{aligned}$$

1501 Defining $\tilde{\mathbf{R}} = (\tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top / n + \lambda_G \mathbf{I}_n)^{-1} \in \mathbb{R}^{n \times n}$, $\mathbf{V} = [n^{-1/2} \mathbf{X}\beta_* | n^{-1/2} \tilde{\mathbf{X}} \beta_*] \in \mathbb{R}^{n \times 2}$, and
 1502

$$1504 \quad \mathbf{D} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix},$$

1507 we have $\bar{\mathbf{R}} = \tilde{\mathbf{R}} + \mathbf{V}\mathbf{D}\mathbf{V}^\top$. Using Woodbury matrix identity (Theorem E.8), $\bar{\mathbf{R}}$ is given by
 1508

$$1509 \quad \bar{\mathbf{R}} = \tilde{\mathbf{R}} - \tilde{\mathbf{R}} \mathbf{V} (\mathbf{D}^{-1} + \mathbf{V}^\top \tilde{\mathbf{R}} \mathbf{V})^{-1} \mathbf{V}^\top \tilde{\mathbf{R}} \quad (29)$$

1511 and plugging this expression into the second term in (28) gives
 1512

$$\begin{aligned} 1513 \quad n^{-1} (\mathbf{X}\beta_*)^\top \bar{\mathbf{R}} \sigma_{*,\perp}(\mathbf{X}\beta_*) &= n^{-1} (\mathbf{X}\beta_*)^\top (\tilde{\mathbf{R}} - \tilde{\mathbf{R}} \mathbf{V} (\mathbf{D}^{-1} + \mathbf{V}^\top \tilde{\mathbf{R}} \mathbf{V})^{-1} \mathbf{V}^\top \tilde{\mathbf{R}}) \sigma_{*,\perp}(\mathbf{X}\beta_*) + o(1) \\ 1514 \\ 1515 &= n^{-1} (\mathbf{X}\beta_*)^\top \tilde{\mathbf{R}} \sigma_{*,\perp}(\mathbf{X}\beta_*) - n^{-1} (\mathbf{X}\beta_*)^\top \tilde{\mathbf{R}} (\mathbf{V} (\mathbf{D}^{-1} + \mathbf{V}^\top \tilde{\mathbf{R}} \mathbf{V})^{-1} \mathbf{V}^\top) \tilde{\mathbf{R}} \sigma_{*,\perp}(\mathbf{X}\beta_*) + o(1). \end{aligned}$$

1517 The first term can be shown to be $o(1)$ in probability by using the fact that $\tilde{\mathbf{R}}$ is independent of $\mathbf{X}\beta_*$ and the orthogonality of
 1518 $\mathbf{X}\beta_*$ and $\sigma_{*,\perp}(\mathbf{X}\beta_*)$. To analyze the second term, first note that the elements of the matrix $\Omega = (\mathbf{D}^{-1} + \mathbf{V}^\top \tilde{\mathbf{R}} \mathbf{V})^{-1} \in \mathbb{R}^{2 \times 2}$
 1519 can all be shown to be $\mathcal{O}(1)$ by a simple norm argument. Moreover,
 1520

$$\begin{aligned} 1521 \quad n^{-1} (\mathbf{X}\beta_*)^\top \tilde{\mathbf{R}} \mathbf{V} (\mathbf{D}^{-1} + \mathbf{V}^\top \tilde{\mathbf{R}} \mathbf{V})^{-1} \mathbf{V}^\top \tilde{\mathbf{R}} \sigma_{*,\perp}(\mathbf{X}\beta_*) \\ 1522 \\ 1523 &= n^{-2} (\mathbf{X}\beta_*)^\top \tilde{\mathbf{R}} \left(\Omega_{11} (\mathbf{X}\beta_*) (\mathbf{X}\beta_*)^\top + \Omega_{12} (\mathbf{X}\beta_*) (\tilde{\mathbf{X}} \beta_*)^\top + \Omega_{21} (\tilde{\mathbf{X}} \beta_*) (\mathbf{X}\beta_*)^\top + \Omega_{22} (\tilde{\mathbf{X}} \beta_*) (\tilde{\mathbf{X}} \beta_*)^\top \right) \tilde{\mathbf{R}} \sigma_{*,\perp}(\mathbf{X}\beta_*) \end{aligned}$$

1525 where Ω_{ij} are the elements of the matrix Ω . We analyze each term in this sum separately and show that all of them are $o(1)$.
 1526

- **First Term.** Using a simple norm argument, $n^{-1} (\mathbf{X}\beta_*)^\top \tilde{\mathbf{R}} (\mathbf{X}\beta_*) = \mathcal{O}(1)$. Also, by construction of $\sigma_{*,\perp}$, we have

$$1529 \quad n^{-1} (\mathbf{X}\beta_*)^\top \tilde{\mathbf{R}} \sigma_{*,\perp}(\mathbf{X}\beta_*) \rightarrow 0.$$

1531 Thus, the whole term is $o(1)$.

- **Second Term.** Similar to the first term, we have $n^{-1} (\mathbf{X}\beta_*)^\top \tilde{\mathbf{R}} (\mathbf{X}\beta_*) = \mathcal{O}(1)$. Also, $n^{-1} (\tilde{\mathbf{X}} \beta_*)^\top \tilde{\mathbf{R}} \sigma_{*,\perp}(\mathbf{X}\beta_*) \rightarrow 0$ in probability, using the weak law of large numbers by noting that $\sigma_{*,\perp}(\mathbf{X}\beta_*)$ is independent of $n^{-1} (\tilde{\mathbf{X}} \beta_*)^\top \tilde{\mathbf{R}}$ by construction and that it has mean zero. Hence, the whole term is $o(1)$.

- **Third Term.** By construction, the vector $\mathbf{X}\beta_*$ is independent of $\tilde{\mathbf{R}} (\tilde{\mathbf{X}} \beta_*)$ and has mean zero, which gives $n^{-1} (\mathbf{X}\beta_*)^\top \tilde{\mathbf{R}} (\tilde{\mathbf{X}} \beta_*) \rightarrow 0$. Also, using a simple norm argument, we have $n^{-1} (\mathbf{X}\beta_*)^\top \tilde{\mathbf{R}} \sigma_{*,\perp}(\mathbf{X}\beta_*) = \mathcal{O}(1)$ which proves that the third term is also $o(1)$.

- **Fourth Term.** This term can be shown to be $o(1)$ with an argument very similar to the argument for the third term.

Putting these all together and using (28), we have

$$\beta_{\star}^{\top} \beta_{\text{KFAC}} = c_{\star,1} d_x^{-1} \operatorname{Tr}((\mathbf{X}^{\top} \mathbf{X}/n) \mathbf{R}) + o(1). \quad (30)$$

Next we move to the analysis of the squared ℓ_2 -norm of the vector β_{KFAC} . By decomposing the function σ_{\star} into a linear and an orthogonal nonlinear component similar to the one used for the analysis of the inner product term above, we write

$$\begin{aligned} \|\beta_{\text{KFAC}}\|_2^2 &= n^{-2} \mathbf{y}^{\top} \mathbf{X} \mathbf{R}^2 \mathbf{X}^{\top} \mathbf{y} = n^{-2} \boldsymbol{\varepsilon}^{\top} \mathbf{X} \mathbf{R}^2 \mathbf{X}^{\top} \boldsymbol{\varepsilon} + c_{\star,1}^2 n^{-2} \beta_{\star}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{R}^2 \mathbf{X}^{\top} \mathbf{X} \beta_{\star} \\ &\quad + n^{-2} \sigma_{\star,\perp}(\mathbf{X} \beta_{\star})^{\top} \mathbf{X} \mathbf{R}^2 \mathbf{X}^{\top} \sigma_{\star,\perp}(\mathbf{X} \beta_{\star}) + 2 c_{\star,1} n^{-2} \beta_{\star}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{R}^2 \mathbf{X}^{\top} \sigma_{\star,\perp}(\mathbf{X} \beta_{\star}). \end{aligned}$$

We will analyze each of these terms separately.

- **First Term.** Recalling that $\boldsymbol{\varepsilon} \sim N(0, \sigma_{\varepsilon}^2 \mathbf{I}_n)$ independent of all randomness in the problem, using the Hanson-Wright inequality (Theorem E.6) we have

$$n^{-2} \boldsymbol{\varepsilon}^{\top} \mathbf{X} \mathbf{R}^2 \mathbf{X}^{\top} \boldsymbol{\varepsilon} = \sigma_{\varepsilon}^2 n^{-1} \operatorname{Tr}((\mathbf{X}^{\top} \mathbf{X}/n) \mathbf{R}^2) + o(1).$$

- **Second Term.** Using Assumption 3.8, and by the Hanson-Wright inequality we have

$$n^{-2} c_{\star,1}^2 \beta_{\star}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{R}^2 \mathbf{X}^{\top} \mathbf{X} \beta_{\star} = c_{\star,1}^2 d_x^{-1} \operatorname{Tr}[(\mathbf{X}^{\top} \mathbf{X}/n)^2 \mathbf{R}^2] + o(1)$$

- **Third Term.** Note that $\bar{\mathbf{R}}$ and $\mathbf{X} \beta_{\star}$ are dependent. Note that $\mathbf{X} \mathbf{R}^2 \mathbf{X}^{\top} = \mathbf{X} \mathbf{X}^{\top} \bar{\mathbf{R}} = n \bar{\mathbf{R}} - \lambda_G n \bar{\mathbf{R}}^2$. Thus, an almost identical argument to the argument used above for the analysis of $\beta_{\star}^{\top} \beta_{\text{KFAC}}$ using $\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{X} \beta_{\star} \beta_{\star}^{\top}$ gives

$$n^{-2} \sigma_{\star,\perp}(\mathbf{X} \beta_{\star})^{\top} \mathbf{X} \mathbf{R}^2 \mathbf{X}^{\top} \sigma_{\star,\perp}(\mathbf{X} \beta_{\star}) = \mathbb{E}_z[\sigma_{\star,\perp}^2(z)] \cdot n^{-1} \operatorname{Tr}((\mathbf{X}^{\top} \mathbf{X}/n) \mathbf{R}^2) + o(1)$$

- **Fourth Term.** This term can readily be shown to be $o(1)$ in the analysis of $\beta_{\star}^{\top} \beta_{\text{KFAC}}$; i.e.,

$$2 c_{\star,1} n^{-2} \beta_{\star}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{R}^2 \mathbf{X}^{\top} \sigma_{\star,\perp}(\mathbf{X} \beta_{\star}) = o(1).$$

Putting everything together, we find

$$\|\beta_{\text{KFAC}}\|_2^2 = c_{\star,1}^2 d_x^{-1} \operatorname{Tr}[(\mathbf{X}^{\top} \mathbf{X}/n)^2 \mathbf{R}^2] + (\sigma_{\varepsilon}^2 + \mathbb{E}_z[\sigma_{\star,\perp}^2(z)]) n^{-1} \operatorname{Tr}((\mathbf{X}^{\top} \mathbf{X}/n) \mathbf{R}^2) \quad (31)$$

Now, given (30) and (31), defining $c_{\star,>1} = \mathbb{E}_z[\sigma_{\star,\perp}^2(z)]$, we have

$$\frac{\beta_{\star}^{\top} \beta_{\text{KFAC}}}{\|\beta_{\text{KFAC}}\| \|\beta_{\star}\|} = \frac{c_{\star,1} d_x^{-1} \operatorname{Tr}((\mathbf{X}^{\top} \mathbf{X}/n) \mathbf{R})}{\sqrt{c_{\star,1}^2 d_x^{-1} \operatorname{Tr}[(\mathbf{X}^{\top} \mathbf{X}/n)^2 \mathbf{R}^2] + (\sigma_{\varepsilon}^2 + c_{\star,>1}) n^{-1} \operatorname{Tr}((\mathbf{X}^{\top} \mathbf{X}/n) \mathbf{R}^2)}}. \quad (32)$$

Thus, noting that if $d_x/n \rightarrow 0$ and $\lambda_G \rightarrow 0$, we have $\mathbf{R} \rightarrow \Sigma_{\mathbf{x}}^{-1}$, we find

$$\lim_{\lambda \rightarrow 0} \lim_{d_x/n \rightarrow \infty} \frac{\beta_{\star}^{\top} \beta_{\text{KFAC}}}{\|\beta_{\text{KFAC}}\| \|\beta_{\star}\|} = 1,$$

proving the second part of the lemma. For the first part, we define $m(z) : \mathbb{R} \rightarrow \mathbb{R}$ as the limiting Stieltjes transform of the empirical eigenvalue distribution of $n^{-1} \mathbf{X}^{\top} \mathbf{X}$; i.e.,

$$m(z) = \lim_{d_x/n \rightarrow \infty} d_x^{-1} \operatorname{Tr} \left[(\mathbf{X}^{\top} \mathbf{X}/n - z \mathbf{I}_{d_x})^{-1} \right] \quad (33)$$

where the limit is taken under the assumption that $d_x/n \rightarrow \phi > 0$. For a general covariance matrix $\Sigma_{\mathbf{x}}$, $m(z)$ does not have a closed form except for very special cases; however, it can be efficiently computed. See Section E.8 for more details. The derivative of the function m is given by

$$m'(z) = - \lim_{d_x/n \rightarrow \infty} d_x^{-1} \operatorname{Tr} \left[(\mathbf{X}^{\top} \mathbf{X}/n - z \mathbf{I}_{d_x})^{-2} \right].$$

1595 We can write all the traces appearing in (32) in terms of the function m and its derivative:

$$1596 \quad d_x^{-1} \operatorname{Tr} ((\mathbf{X}^\top \mathbf{X}/n) \mathbf{R}) = d_x^{-1} \operatorname{Tr} ((\mathbf{X}^\top \mathbf{X}/n + \lambda_G \mathbf{I}_{d_x} - \lambda_G \mathbf{I}_{d_x}) \mathbf{R}) = d_x^{-1} \operatorname{Tr} (\mathbf{I}_{d_x} - \lambda_G \mathbf{R}) = 1 - \lambda_G m(-\lambda_G),$$

$$1599 \quad d_x^{-1} \operatorname{Tr} ((\mathbf{X}^\top \mathbf{X}/n)^2 \mathbf{R}^2) = d_x^{-1} \operatorname{Tr} ((\mathbf{X}^\top \mathbf{X}/n + \lambda_G \mathbf{I}_{d_x} - \lambda_G \mathbf{I}_{d_x})^2 \mathbf{R}^2) = 1 - \lambda_G^2 m'(-\lambda_G) - 2\lambda_G m(-\lambda_G),$$

$$1600 \quad n^{-1} \operatorname{Tr} ((\mathbf{X}^\top \mathbf{X}/n) \mathbf{R}^2) = n^{-1} \operatorname{Tr} (\mathbf{R} - \lambda_G \mathbf{R}^2) = \phi m(-\lambda_G) + \phi \lambda_G m'(-\lambda_G).$$

1603 With these, the correlation is given by

$$1605 \quad \frac{\beta_{\text{KFAC}}^\top \beta_\star}{\|\beta_{\text{KFAC}}\|_2 \|\beta_\star\|_2} = \frac{c_{\star,1}[1 - \lambda_G m(-\lambda_G)]}{\sqrt{c_{\star,1}^2[1 - \lambda_G^2 m'(-\lambda_G) - 2\lambda_G m(-\lambda_G)] + \phi(c_{\star,>1}^2 + \sigma_\varepsilon^2)[m(-\lambda_G) + \lambda_G m'(-\lambda_G)]}},$$

1609 which defining

$$\begin{aligned} 1610 \quad \Psi_1 &= 1 - \lambda_G m(-\lambda_G) \\ 1611 \quad \Psi_2 &= 1 - \lambda_G^2 m'(-\lambda_G) - 2\lambda_G m(-\lambda_G) \\ 1612 \quad \Psi_3 &= m(-\lambda_G) + \lambda_G m'(-\lambda_G) \end{aligned} \tag{34}$$

1615 concludes the proof. \square

1617 C.4. From Feature Learning to Generalization

1619 In Section 3.2, we showed that after one step of SGD and KFAC, the first layer weights will become approximately equal to

$$1620 \quad \hat{\mathbf{G}}_a \approx \hat{\mathbf{G}}_0 + \alpha \eta \mathbf{f}_0 \hat{\beta}_a^\top, \quad a \in \{\text{SGD, KFAC}\}. \tag{35}$$

1622 Given Lemma 3.10 and Lemma 3.12, we argued that compared to SGD, the weights obtained by the KFAC algorithm are
1623 more aligned to the true direction β_\star . Given a nontrivial alignment between the weights and the target direction, the second
1624 layer \mathbf{f} can be trained using least squares (or based on Lemma 3.1, equivalently using one step of the KFAC update on \mathbf{f} with
1625 $\eta_f = 1$) with $\Theta(d)$ samples to achieve good generalization performance (See e.g., Ba et al. (2022, Theorem 11) and Dandi
1626 et al. (2024c, Section 3.4)). The existence of nontrivial alignment of the learned weights and the true direction in a single
1627 index model is often called *weak recovery* and has been subject to extensive investigation (see e.g., Ben Arous et al. (2021);
1628 Dandi et al. (2024c); Troiani et al. (2024); Arnaboldi et al. (2024), etc.).

1629 To see this, consider the feature matrix $\mathbf{Z}_a \in \mathbb{R}^{n \times d_h}$ as $\mathbf{Z}_a = \sigma(\mathbf{X} \hat{\mathbf{G}}_a^\top)$, where the activation function is applied element-
1630 wise. Based on equation (35), this matrix can be written as

$$1632 \quad \mathbf{Z}_a \approx \sigma \left(\mathbf{X} \mathbf{G}_0^\top + \alpha \eta (\mathbf{X} \hat{\beta}_a) \mathbf{f}_0^\top \right).$$

1634 This is an example of a random matrix in which a nonlinear function is applied element-wise to a random component plus a
1635 rank-one signal component which has been studied in the literature (Guionnet et al., 2023; Moniri et al., 2024; Moniri &
1636 Hassani, 2024b). In particular, by Taylor expanding the activation function, the feature matrix \mathbf{Z}_a can be written as

$$1639 \quad \mathbf{Z}_a \approx \sigma(\mathbf{X} \mathbf{G}_0^\top) + \sum_{k=1}^{\ell} \frac{\alpha^k \eta^k}{k!} \left(\sigma^{(k)}(\mathbf{X} \mathbf{G}_0^\top) \right) \odot \left((\mathbf{X} \hat{\beta}_a)^{\circ k} \mathbf{f}_0^{\circ k \top} \right) + \mathcal{E}_\ell,$$

1641 where \circ denotes element-wise power and \mathcal{E}_ℓ is the reminder term. Let $\eta = n^\alpha$ for some $\alpha \in [0, 0.5]$. Given α , the
1642 integer ℓ is chosen to be large enough so that the operator norm of \mathcal{E}_ℓ is $o(d_x^{1/2})$ and the reminder term is negligible
1643 compared to $\sigma(\mathbf{X} \mathbf{G}_0^\top)$. By a simple concentration argument, the matrix $\sigma^{(k)}(\mathbf{X} \mathbf{G}_0^\top)$ can be replaced with its mean
1644 $E(\sigma^{(k)}(\mathbf{X} \mathbf{G}_0^\top)) = \mu \mathbf{1} \mathbf{1}^\top$ to get

$$1647 \quad \mathbf{Z}_a \approx \sigma(\mathbf{X} \mathbf{G}_0^\top) + \sum_{k=1}^{\ell} \frac{\alpha^k \eta^k \mu}{k!} (\mathbf{X} \hat{\beta}_a)^{\circ k} \mathbf{f}_0^{\circ k \top}.$$

1650 The first term $\sigma(\mathbf{X}\mathbf{G}_0^\top)$ is the feature matrix of a random feature model and based on the Gaussian Equivalence Theorem
 1651 (GET) (see e.g., Goldt et al. (2022); Hu & Lu (2023); Dandi et al. (2024a); Moniri et al. (2024)), we can linearize it; i.e.,
 1652 we can replace it with $\alpha\mathbf{X}\mathbf{G}_0^\top + \mathbf{N}$ where \mathbf{N} is a properly scaled independent Gaussian noise. The vectors $(\mathbf{X}\hat{\beta}_a)^{\circ k}$ are
 1653 nonlinear functions of the covariates with different degrees. The least squares estimator $\hat{\mathbf{f}}_a$ is then fit on the features \mathbf{Z}_a in a
 1654 way that $\hat{\mathcal{L}}(\hat{\mathbf{f}}_a, \hat{\mathbf{G}}_a)$ is minimized; i.e.

$$y \approx \sigma(\mathbf{X}\mathbf{G}_0^\top)\hat{\mathbf{f}}_a + \sum_{k=1}^{\ell} \frac{\alpha^k \eta^k \mu(\mathbf{f}_0^{\circ k \top} \hat{\mathbf{f}}_a)}{k!} (\mathbf{X}\hat{\beta}_a)^{\circ k}. \quad (36)$$

1661 Based on the GET, the random feature component can only learn linear functions with sample complexity of learning
 1662 $n = \Theta(d_h) = \Theta(d_x)$. When η is large enough and β_a is aligned to β_* , with the finite dimensional correction to the random
 1663 features model, the model can also represent nonlinear functions $(\mathbf{x}^\top \hat{\beta}_a)^k$ of degree $k \leq \ell$ by matching the coefficients
 1664 $\alpha^k \eta^k \mu(\mathbf{f}_0^{\circ k \top} \hat{\mathbf{f}}_a)/k!$ with the Taylor coefficients of the teacher function $\sigma_*(\mathbf{x}^\top \hat{\beta}_*)$.

1666 Although we have provided a complete proof sketch for providing generalization guarantees given weight alignment, a
 1667 complete analysis require tedious computations and is beyond the scope of this work as we mainly focus on feature learning
 1668 properties of different optimization algorithms.

1670 D. Additional Information Regarding Kronecker-Factored Preconditioners

1672 Here, we provide some additional background information regarding key Kronecker-Factored preconditioning methods,
 1673 including their derivation and relations to various methods in the literature. We recall the running example of a fully-
 1674 connected net omitting biases, introducing layer-wise dimensionality and a final non-linear layer (e.g. softmax) for
 1675 completeness:

$$f_\theta(\mathbf{x}) = \phi(\mathbf{W}_L \sigma(\mathbf{W}_{L-1} \cdots \sigma(\mathbf{W}_1 \mathbf{x}) \cdots)), \quad \mathbf{W}_\ell \in \mathbb{R}^{d_\ell \times d_{\ell-1}}, d_0 = d_x. \quad (37)$$

1680 As before, we define θ as the concatenation of $\theta_\ell = \text{vec}(\mathbf{W}_\ell)$, $\ell \in [L]$. We define an expected loss induced by the neural
 1681 network $\mathcal{L}(\theta) = E_{(\mathbf{x}, \mathbf{y})}[\ell(f_\theta(\mathbf{x}), \mathbf{y})]$, and its batch counterpart $\hat{\mathcal{L}}(\theta)$. Here, we define the family of Kronecker-Factored
 1682 preconditioned optimizers as those that update weights in the following fashion:

$$\mathbf{W}_{\ell+} = \mathbf{W}_\ell - \eta \mathbf{P}_\ell^{-1} \nabla_{\mathbf{W}_\ell} \hat{\mathcal{L}}(\theta) \mathbf{Q}_\ell^{-1}, \quad \ell \in [L],$$

1688 where $\mathbf{P}_\ell \in \mathbb{R}^{d_\ell \times d_\ell}$, $\mathbf{Q}_\ell \in \mathbb{R}^{d_{\ell-1} \times d_{\ell-1}}$, $\ell \in [L]$ are square matrices. For simplicity, we ignore moving parts such as
 1689 momentum, damping exponents, adaptive learning rate schedules/regularization etc. We now demonstrate the basic principles
 1690 and derivation of certain notable members of these preconditioning methods on the feedforward network (37).

1692 D.1. Kronecker-Factored Approximate Curvature KFAC

1694 As described in the main paper, KFAC (Martens & Grosse, 2015) is at its core an approximation to natural gradient
 1695 descent. Given that we are approximating NGD, a crucial presumption on $f_\theta(\mathbf{x})$ and $\mathcal{L}(\theta)$ is that the network output $f_\theta(\mathbf{x})$
 1696 parameterizes a conditional distribution $p(\mathbf{y}|\mathbf{x}; \theta)$, and $\mathcal{L}(\theta) \propto E_{(\mathbf{x}, \mathbf{y})}[-\log p(\mathbf{y}|\mathbf{x}; \theta)]$ is the corresponding negative log-
 1697 likelihood. As such, KFAC is technically only applicable to settings where such an interpretation exists. However, this notably
 1698 subsumes cases $\mathcal{L}(\theta) = E_{(\mathbf{x}, \mathbf{y})}[\ell(f_\theta(\mathbf{x}), \mathbf{y})]$, where $\ell(\cdot)$ is a strictly convex function in $f_\theta(\mathbf{x})$, as this admits an interpretation
 1699 as $f_\theta(\mathbf{x})$ parameterizing an exponential family distribution. In particular, the square-loss regression case $\ell(\hat{\mathbf{y}}, \mathbf{y}) = \|\hat{\mathbf{y}} - \mathbf{y}\|^2$
 1700 corresponds to a conditionally-Gaussian predictive distribution with fixed variance $\hat{\mathbf{y}}(\mathbf{x}) \sim N(f_\theta(\mathbf{x}), \sigma^2 \mathbf{I})$, and if $\phi(\cdot)$
 1701 is a softmax layer and $\ell(\hat{\mathbf{y}}, \mathbf{y}) = \text{CrossEnt}(\hat{\mathbf{y}}, \mathbf{y})$, the multi-class classification case corresponds to a conditionally-
 1702 multinomial predictive distribution.

1703 Defining $\mathbf{h}_\ell = \mathbf{W}_\ell \mathbf{z}_{\ell-1}$, $\mathbf{z}_\ell = \sigma(\mathbf{h})$, $\mathbf{z}_0 = \mathbf{x}$, the Fisher Information of the predictive distribution $p(\mathbf{y}|\mathbf{x}; \theta)$ at θ can be

expressed in block form:

$$\begin{aligned} \mathbf{FI}(\boldsymbol{\theta}) &\triangleq \mathbb{E}_{\mathbf{x}} \left[\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \left(\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^{\top} \right] && \text{(recall } \boldsymbol{\theta} \text{ is vec-ed parameters)} \\ &= \begin{bmatrix} \mathbb{E}_{\mathbf{x}} \left[\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}_1} \left(\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}_1} \right)^{\top} \right] & \dots & \mathbb{E}_{\mathbf{x}} \left[\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}_1} \left(\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}_L} \right)^{\top} \right] \\ \vdots & \ddots & \vdots \\ \mathbb{E}_{\mathbf{x}} \left[\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}_L} \left(\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}_1} \right)^{\top} \right] & \dots & \mathbb{E}_{\mathbf{x}} \left[\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}_L} \left(\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}_L} \right)^{\top} \right] \end{bmatrix} \end{aligned}$$

Looking at the (i, j) th block, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{x}} \left[\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}_i} \left(\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}_j} \right)^{\top} \right] &= \mathbb{E}_{\mathbf{x}} \left[\text{vec} \left(\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \mathbf{W}_i} \right) \text{vec} \left(\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \mathbf{W}_j} \right)^{\top} \right] \\ &= \mathbb{E}_{\mathbf{x}} [(\mathbf{z}_{i-1} \otimes \mathbf{g}_i)(\mathbf{z}_{j-1} \otimes \mathbf{g}_j)^{\top}] && (\mathbf{g}_\ell \triangleq -\frac{\partial p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \mathbf{h}_\ell}) \\ &= \mathbb{E}_{\mathbf{x}} [(\mathbf{z}_{i-1} \mathbf{z}_{j-1}^{\top}) \otimes (\mathbf{g}_i \mathbf{g}_j^{\top})], && \text{(Lemma E.1, item 2)} \end{aligned}$$

where the second line comes from writing out the backpropagation formula. KFAC makes two key approximations:

1. The matrix $\mathbf{FI}(\boldsymbol{\theta})^{-1}$ is approximated by a block-diagonal, and hence so is $\mathbf{FI}(\boldsymbol{\theta})$. We note the original formulation of KFAC in Martens & Grosse (2015) also supports a tridiagonal inverse approximation.
2. The vectors $\mathbf{z}_{\ell-1}$ and \mathbf{g}_ℓ are independent for all $\ell \in [L]$, such that

$$\mathbb{E}_{\mathbf{x}} [(\mathbf{z}_{\ell-1} \mathbf{z}_{\ell-1}^{\top}) \otimes (\mathbf{g}_\ell \mathbf{g}_\ell^{\top})] = \mathbb{E}[\mathbf{z}_{\ell-1} \mathbf{z}_{\ell-1}^{\top}] \otimes \mathbb{E}[\mathbf{g}_\ell \mathbf{g}_\ell^{\top}].$$

Now replacing the true expectation with the empirical estimate, and defining $\mathbf{P}_\ell = \widehat{\mathbb{E}}[\mathbf{g}_\ell \mathbf{g}_\ell^{\top}]$, $\mathbf{Q}_\ell = \widehat{\mathbb{E}}[\mathbf{z}_{\ell-1} \mathbf{z}_{\ell-1}^{\top}]$ completes the Kronecker-Factored approximation to the Fisher Information. It is clear to see from the derivation that, as we previewed in the introduction and expressed emphatically in Martens & Grosse (2015), this approximation is *never* expected to be tight.

Some related preconditioners

Having introduced KFAC, we introduce some related preconditioners. Notably, it has been noted that computing \mathbf{g}_ℓ requires a backwards gradient computation, whereas \mathbf{z}_ℓ only requires a forward pass. In particular, various works have recovered the *right* preconditioner \mathbf{Q}_ℓ of KFAC via various notions of “local” (layer-wise) losses. Notably, these alternative views allow KFAC-like preconditioning to extend beyond the negative-log-likelihood interpretation.

- LocoProp, square-loss case (Amid et al., 2022):

$$\begin{aligned} \text{Update rule : } \mathbf{W}_{\ell+} &= \underset{\mathbf{W}}{\operatorname{argmin}} \frac{1}{2} \widehat{\mathbb{E}} [\|\mathbf{W}\mathbf{z}_{\ell-1} - \mathbf{h}_\ell\|^2] + \frac{1}{2\eta} \|\mathbf{W} - \mathbf{W}_\ell\|_F^2 \\ &= \mathbf{W}_\ell - \eta \nabla_{\mathbf{W}_\ell} \widehat{\mathcal{L}}(\boldsymbol{\theta}) \left(\mathbf{I}_{d_{\ell-1}} + \eta \widehat{\mathbb{E}}[\mathbf{z}_{\ell-1} \mathbf{z}_{\ell-1}^{\top}] \right)^{-1}. \end{aligned}$$

As noted in Amid et al. (2022), this update is also closely related to ProxProp (Frerix et al., 2018).

- FOOF (Benzing, 2022):

$$\begin{aligned} \text{Update rule : } \Delta \mathbf{W}_\ell &= \underset{\Delta \mathbf{W}}{\operatorname{argmin}} \widehat{\mathbb{E}} [\|\Delta \mathbf{W} \mathbf{z}_{\ell-1} - \eta \mathbf{g}_\ell\|^2] + \frac{\lambda}{2} \|\Delta \mathbf{W}\|_F^2 && \left(\mathbf{g}_\ell = \frac{\partial \ell(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y})}{\partial \mathbf{h}_\ell} \right) \\ &= \eta \nabla_{\mathbf{W}_\ell} \widehat{\mathcal{L}}(\boldsymbol{\theta}) \left(\widehat{\mathbb{E}}[\mathbf{z}_{\ell-1} \mathbf{z}_{\ell-1}^{\top}] + \lambda \mathbf{I}_{d_{\ell-1}} \right)^{-1}, \end{aligned}$$

$$\mathbf{W}_{\ell+} = \mathbf{W}_\ell - \Delta \mathbf{W}_\ell.$$

Interestingly, we note that these right-preconditioner-only variants subsume the DFW algorithm for two-layer linear representation learning proposed in Zhang et al. (2024b); thus we may see the guarantee therein as support of the above algorithms from a feature learning perspective, albeit weaker than Theorem 3.6.

D.2. Shampoo

Shampoo is designed to be a Kronecker-Factored approximation of the full AdaGrad preconditioner, which we recall is the running sum of the outer-product of loss gradients. Turning off the AdaGrad accumulator and instead considering the empirical batch estimate $\hat{E}[\nabla_{\theta} \ell(f_{\theta}(\mathbf{x}), \mathbf{y}) \nabla_{\theta} \ell(f_{\theta}(\mathbf{x}), \mathbf{y})^\top]$, the curvature matrix being estimated can also be viewed as the Gauss-Newton matrix $E_{(\mathbf{x}, \mathbf{y})}[\nabla_{\theta} \ell(f_{\theta}(\mathbf{x}), \mathbf{y}) \nabla_{\theta} \ell(f_{\theta}(\mathbf{x}), \mathbf{y})^\top]$. As documented in various works (see e.g. Martens (2020)), the (generalized) Gauss-Newton matrix in many cases is related or equal to the Fisher Information, establishing a link between the target curvatures of KFAC and Shampoo.

However, the Shampoo preconditioners differ from KFAC's. Let us define the $\mathbf{z}_\ell, \mathbf{h}_\ell$ as before, and $\mathbf{g}_\ell = \frac{\partial \ell(f_{\theta}(\mathbf{x}), \mathbf{y})}{\partial \mathbf{h}_\ell}$. Then, the Shampoo preconditioners are given by

$$\mathbf{P}_\ell = \hat{E}[\mathbf{g}_\ell \mathbf{z}_{\ell-1}^\top (\mathbf{g}_\ell \mathbf{z}_{\ell-1}^\top)^\top]^{1/4}, \quad \mathbf{Q}_\ell = \hat{E}[\mathbf{z}_{\ell-1} \mathbf{g}_\ell^\top (\mathbf{z}_{\ell-1} \mathbf{g}_\ell^\top)^\top]^{1/4}.$$

Notably, Shampoo takes the fourth root in the preconditioners, as its target is the AdaGrad preconditioner which is (modulo scaling) the square-root of the empirical Gauss-Newton matrix—analogous to the square-root of the second moment in Adam. Whether the target curvature should be the square-root or not of the Gauss-Newton matrix is the topic of recent discussion (Morwani et al., 2024; Lin et al., 2024).

D.3. Kronecker-Factored Preconditioners and the Modular Norm

The “modular norm” (Large et al., 2024; Bernstein & Newhouse, 2024a;b) is a recently introduced notion that provides a general recipe for producing different optimization algorithms that act layer-wise. By specifying different norms customized for different kinds of layers (e.g. feed-forward, residual, convolutional etc.), one in principle has the flexibility to customize an optimizer to handle the different kinds of curvature induced by different parameter spaces. Given a choice of norm on the weight tensor \mathbf{W}_ℓ , the descent direction is returned by *steepest descent* with respect to that norm. To introduce steepest descent, we require a few definitions (cf. Bernstein & Newhouse (2024a)):

Definition D.1 (Dual norms, steepest direction). Given a norm $\|\cdot\|$ defined over a finite-dimensional real vector space \mathcal{V} . The dual norm $\|\cdot\|_\dagger$ is defined by

$$\|\mathbf{v}\|_\dagger = \max_{\|\mathbf{u}\|=1} \langle \mathbf{u}, \mathbf{v} \rangle.$$

With $\mathbf{g} \in \mathcal{V}$ and a “sharpness” parameter $\eta > 0$, the steepest direction(s) are given by the following variational representation:

$$\operatorname{argmin}_{\mathbf{d}} \left[\langle \mathbf{g}, \mathbf{d} \rangle + \frac{1}{2\eta} \|\mathbf{d}\|^2 \right] = -\eta \|\mathbf{g}\|_\dagger \cdot \operatorname{argmax}_{\|\mathbf{u}\|=1} \langle \mathbf{g}, \mathbf{u} \rangle.$$

Here we focus on finite-dimensional normed spaces, but note that these concepts extend *mutatis mutandis* to general Banach spaces. The aforementioned works derive various standard optimizers by choosing different norms, including *induced matrix norms* $\|\mathbf{W}_\ell\|_{\alpha \rightarrow \beta} = \max_{\mathbf{x}} \frac{\|\mathbf{W}_\ell \mathbf{x}\|_\beta}{\|\mathbf{x}\|_\alpha}$, applied to a given layer's weight space, for example (Bernstein & Newhouse, 2024b):

- SGD: induced by Frobenius (Euclidean) norm $\|\cdot\| = \|\cdot\|_F$. Note the Frobenius norm is *not* an induced matrix norm.
- Sign-descent (“ideal” Adam with EMA on moments turned off): induced by $\|\cdot\| = \|\cdot\|_{\ell_1 \rightarrow \ell_\infty}$.
- Shampoo (“ideal” variant with moment accumulator turned off): induced by $\|\cdot\| = \|\cdot\|_{\ell_2 \rightarrow \ell_2} = \|\cdot\|_{\text{op}}$.

Therefore, in light of this characterization, a natural question to ask is *what norm induces a given Kronecker-Factored preconditioner* (which includes Shampoo). We provide a simple derivation that determines the norm.

1815 **Proposition D.2** (Kronecker-Factored matrix norm). *Recall the fully-connected network (37). Given preconditioners
1816 $\{(\mathbf{P}_\ell, \mathbf{Q}_\ell)\}_{\ell=1}^L$, where $\mathbf{P}_\ell \in \mathbb{R}^{d_\ell \times d_\ell}$, $\mathbf{Q}_\ell \in \mathbb{R}^{d_{\ell-1} \times d_{\ell-1}}$, $\ell \in [L]$ are invertible square matrices. Then, the layer-wise
1817 Kronecker-Factored update:*

$$1819 \quad \mathbf{W}_{\ell+} = \mathbf{W}_\ell - \eta \mathbf{P}_\ell^{-1} \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}) \mathbf{Q}_\ell^{-1}, \quad \ell \in [L]$$

1821 is equivalent to layer-wise steepest descent with norm $\|\mathbf{M}_\ell\| \triangleq \|\mathbf{P}_\ell^\top \mathbf{M}_\ell \mathbf{Q}_\ell^\top\|_F$:

$$1823 \quad \operatorname{argmin}_{\mathbf{M}} \left[\langle \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}), \mathbf{M} \rangle + \frac{1}{2\eta} \|\mathbf{M}\|^2 \right] = -\eta \mathbf{P}_\ell^{-1} \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}) \mathbf{Q}_\ell^{-1}.$$

1827 *Proof of Proposition D.2.* It is straightforward to verify $\|\mathbf{M}\| \triangleq \|\mathbf{P}^\top \mathbf{M} \mathbf{Q}^\top\|_F$ for invertible \mathbf{P}, \mathbf{Q} satisfies the axioms of a
1828 norm. It remains to verify the steepest descent direction:

$$1830 \quad \operatorname{argmin}_{\mathbf{M}} \left[\langle \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}), \mathbf{M} \rangle + \frac{1}{2\eta} \|\mathbf{M}\|^2 \right] = -\eta \|\nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta})\|_\dagger \cdot \operatorname{argmax}_{\|\mathbf{M}\|=1} \langle \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}), \mathbf{M} \rangle = \mathbf{P}_\ell^{-1} \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}) \mathbf{Q}_\ell^{-1}.$$

1833 We start by writing:

$$\begin{aligned} 1835 \quad \|\nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta})\|_\dagger &\triangleq \max_{\|\mathbf{M}\|=1} \langle \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}), \mathbf{M} \rangle \\ 1836 \\ 1837 \quad &= \max_{\|\mathbf{P}_\ell^\top \mathbf{M} \mathbf{Q}_\ell^\top\|_F=1} \operatorname{Tr}(\mathbf{M}^\top \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta})) \\ 1838 \\ 1839 \quad &= \max_{\|\mathbf{D}\|_F=1} \operatorname{Tr}(\mathbf{Q}_\ell^{-1} \mathbf{D}^\top \mathbf{P}_\ell^{-1} \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta})) && (\mathbf{P}_\ell^\top \mathbf{M} \mathbf{Q}_\ell^\top \rightarrow \mathbf{D}) \\ 1840 \\ 1841 \quad &= \|\mathbf{P}_\ell^{-1} \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}) \mathbf{Q}_\ell^{-1}\|_F. && \text{(trace cyclic property, } \|\cdot\|_F \text{ is self-dual)} \\ 1842 \\ 1843 \end{aligned}$$

1844 Similarly, it is straightforward to verify that the maximizing matrix is:

$$1846 \quad \operatorname{argmax}_{\|\mathbf{M}\|=1} \langle \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}), \mathbf{M} \rangle = \frac{\mathbf{P}_\ell^{-1} \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}) \mathbf{Q}_\ell^{-1}}{\|\mathbf{P}_\ell^{-1} \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}) \mathbf{Q}_\ell^{-1}\|_F},$$

1849 such that plugging it into the steepest descent expression yields:

$$\begin{aligned} 1851 \quad \operatorname{argmin}_{\mathbf{M}} \left[\langle \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}), \mathbf{M} \rangle + \frac{1}{2\eta} \|\mathbf{M}\|^2 \right] &= -\eta \|\mathbf{P}_\ell^{-1} \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}) \mathbf{Q}_\ell^{-1}\|_F \cdot \frac{\mathbf{P}_\ell^{-1} \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}) \mathbf{Q}_\ell^{-1}}{\|\mathbf{P}_\ell^{-1} \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}) \mathbf{Q}_\ell^{-1}\|_F} \\ 1852 \\ 1853 \quad &= -\eta \mathbf{P}_\ell^{-1} \nabla_{\mathbf{W}_\ell} \mathcal{L}(\boldsymbol{\theta}) \mathbf{Q}_\ell^{-1}, \\ 1854 \\ 1855 \end{aligned}$$

1856 as required. □

1857
1858 We remark that for complex-valued matrices, the above holds without modification for the Hermitian transpose \mathbf{A}^H . Notably,
1859 the layer-wise norm corresponding to Kronecker-Factored preconditioning is not an induced matrix norm, though modified
1860 optimizers can certainly be derived via induced-norm variants, such as a “Mahalonobis-to-Mahalonobis” induced norm:
1861

$$\begin{aligned} 1863 \quad \|\mathbf{M}\|_{\mathbf{Q}^{-1} \rightarrow \mathbf{P}} &\triangleq \max_{\mathbf{x}} \frac{\sqrt{(\mathbf{M}\mathbf{x})^\top \mathbf{P}(\mathbf{M}\mathbf{x})}}{\sqrt{\mathbf{x}^\top \mathbf{Q}^{-1} \mathbf{x}}} && (\mathbf{P}, \mathbf{Q} \succ \mathbf{0}) \\ 1864 \\ 1865 \quad &= \max_{\|\mathbf{x}\|=1} \|\mathbf{P}^{1/2} \mathbf{M} \mathbf{Q}^{1/2}\| \\ 1866 \\ 1867 \quad &= \left\| \mathbf{P}^{1/2} \mathbf{M} \mathbf{Q}^{1/2} \right\|_{\text{op}}. \\ 1868 \\ 1869 \end{aligned}$$

1870 E. Auxiliary Results

1871 E.1. Properties of Kronecker Product

1873 Recall the definition of the Kronecker Product: given $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{p \times q}$

$$1875 \quad \mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} A_{11}\mathbf{B} & \cdots & A_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ A_{m1}\mathbf{B} & \cdots & A_{mn}\mathbf{B} \end{bmatrix} \in \mathbb{R}^{mp \times nq}.$$

1879 Complementarily, the vectorization operator $\text{vec}(\mathbf{A})$ is defined by stacking the columns of \mathbf{A} on top of each other (i.e.
1880 column-major order)

$$1882 \quad \text{vec}(\mathbf{A}) = [A_{11} \ \cdots \ A_{m1} \ \cdots \ A_{1n} \ \cdots \ A_{mn}]^\top \in \mathbb{R}^{mn}.$$

1883 We now introduce some fundamental facts about the Kronecker Product.

1884 **Lemma E.1** (Kronecker-Product Properties). *The following properties hold:*

- 1887 1. $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$. Holds for Moore-Penrose pseudoinverse † as well.
- 1888 2. For size-compliant $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$, we have $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{B} \otimes \mathbf{D})$.
- 1889 3. $\text{vec}(\mathbf{AXB}) = (\mathbf{B}^\top \otimes \mathbf{A})\text{vec}(\mathbf{X})$.

1892 E.2. Covariance Concentration

1893 We often use the following Gaussian covariance concentration result.

1895 **Lemma E.2** (Gaussian covariance concentration). *Let $\mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \Sigma_x)$ for $i = 1, \dots, n$, where $\mathbf{x}_i \in \mathbb{R}^d$. Defining the
1896 empirical covariance matrix $\widehat{\Sigma}_x \triangleq \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, as long as $n \geq \frac{18.27}{c^2}(d + \log(1/\delta))$, we have with probability at least
1897 $1 - \delta$,*

$$1899 \quad (1 - c)\Sigma_x \preceq \widehat{\Sigma}_x \preceq (1 + c)\Sigma_x.$$

1901 *Proof of Lemma E.2.* The result follows essentially from combining a by-now standard concentration inequality for Gaussian
1902 quadratic forms and a covering number argument. To be precise, we observe that

$$1904 \quad \left\| \widehat{\Sigma}_x - \Sigma_x \right\|_{\text{op}} \leq c\|\Sigma_x\| \implies (1 - c)\Sigma_x \preceq \widehat{\Sigma}_x \preceq (1 + c)\Sigma_x.$$

1906 Therefore, it suffices to establish a concentration bound on $\|\widehat{\Sigma}_x - \Sigma_x\|$ and invert for $c\|\Sigma_x\|$. To do so, we recall a standard
1907 covering argument (see e.g. [Vershynin \(2018, Chapter 4\)](#)) yields: given an ε -covering of \mathbb{S}^{d-1} , $\mathcal{N} \triangleq \mathcal{N}(\mathbb{S}^{d-1}, \|\cdot\|_2, \varepsilon)$, the
1908 operator norm of a symmetric matrix Σ is bounded by

$$1910 \quad \|\Sigma\| \leq \frac{1}{1 - 2\varepsilon} \max_{\mathbf{u} \in \mathcal{N}} \mathbf{u}^\top \Sigma \mathbf{u},$$

1912 where the corresponding covering number is bounded by:

$$1914 \quad |\mathcal{N}(\mathbb{S}^{d-1}, \|\cdot\|_2, \varepsilon)| \leq \left(1 + \frac{2}{\varepsilon}\right)^d.$$

1917 As such it suffices to provide a concentration bound on $\mathbf{u}^\top \Sigma \mathbf{u}$ for each $\mathbf{u} \in \mathcal{N}$ and then union-bound. Toward establishing
1918 this, we first state the Gaussian quadratic form concentration bound due to [Hsu et al. \(2012\)](#), which is in turn an instantiation
1919 of a chi-squared concentration bound from [Laurent & Massart \(2000\)](#).

1920 **Proposition E.3** (Prop. 1 in [\(Hsu et al., 2012\)](#)). *Let $\mathbf{A} \in \mathbb{R}^{m \times d}$ be a fixed matrix. Let $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ be a mean-zero,
1921 isotropic Gaussian random vector. For any $\delta \in (0, 1)$, we have*

$$1923 \quad \mathbb{P}[\|\mathbf{Ag}\|^2 > \text{Tr}(\mathbf{A}^\top \mathbf{A}) + 2\sqrt{\text{Tr}((\mathbf{A}^\top \mathbf{A})^2) \log(1/\delta)} + 2\|\mathbf{A}^\top \mathbf{A}\|_{\text{op}} \log(1/\delta)] \leq \delta.$$

Now, given $\mathbf{u} \in \mathbb{S}^{d-1}$, setting $\mathbf{A} = \mathbf{u}^\top \Sigma^{1/2}$ such that $\mathbf{u}^\top \Sigma^{1/2} \mathbf{g} \stackrel{d}{=} \mathbf{u}^\top \mathbf{x}$, instantiating Proposition E.3 yields:

$$\mathsf{P} \left[\mathbf{u}^\top \widehat{\Sigma}_{\mathbf{x}} \mathbf{u} > \mathbf{u}^\top \Sigma_{\mathbf{x}} \mathbf{u} + 2\mathbf{u}^\top \Sigma_{\mathbf{x}} \mathbf{u} \sqrt{\frac{\log(1/\delta)}{n}} + 2\mathbf{u}^\top \Sigma_{\mathbf{x}} \mathbf{u} \frac{\log(1/\delta)}{n} \right] \leq \delta.$$

Put another way, this says with probability at least $1 - \delta$:

$$\mathbf{u}^\top (\widehat{\Sigma}_{\mathbf{x}} - \Sigma) \mathbf{u} \leq 2\mathbf{u}^\top \Sigma_{\mathbf{x}} \mathbf{u} \left(\sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n} \right).$$

Taking a union bound over $\mathbf{u} \in \mathcal{N}$, we get with probability at least $1 - \delta$:

$$\begin{aligned} \max_{\mathbf{u} \in \mathcal{N}} \mathbf{u}^\top (\widehat{\Sigma}_{\mathbf{x}} - \Sigma) \mathbf{u} &\leq \max_{\mathbf{u} \in \mathcal{N}} 2\mathbf{u}^\top \Sigma_{\mathbf{x}} \mathbf{u} \left(\sqrt{\frac{\log(|\mathcal{N}|/\delta)}{n}} + \frac{\log(|\mathcal{N}|/\delta)}{n} \right) \\ &\leq 2 \|\Sigma_{\mathbf{x}}\|_{\text{op}} \left(\sqrt{\frac{d \log(1 + \frac{2}{\varepsilon}) + \log(1/\delta)}{n}} + \frac{d \log(1 + \frac{2}{\varepsilon}) + \log(1/\delta)}{n} \right) \\ &\leq 4 \sqrt{\log\left(1 + \frac{2}{\varepsilon}\right)} \|\Sigma_{\mathbf{x}}\|_{\text{op}} \sqrt{\frac{d + \log(1/\delta)}{n}}, \end{aligned}$$

as long as $n \geq d \log(1 + \frac{2}{\varepsilon}) + \log(1/\delta)$. Chaining together inequalities, this yields with probability at least $1 - \delta$ under the same condition on n :

$$\|\widehat{\Sigma}_{\mathbf{x}} - \Sigma_{\mathbf{x}}\|_{\text{op}} \leq \|\Sigma_{\mathbf{x}}\|_{\text{op}} \frac{2}{1 - \varepsilon} \sqrt{\log\left(1 + \frac{2}{\varepsilon}\right)} \sqrt{\frac{d + \log(1/\delta)}{n}}.$$

Minimizing the RHS for $\varepsilon \approx 0.0605$ yields the result. \square

E.3. Extensions to subgaussianity

As previewed, many results can be extended from the Gaussian setting to subgaussian random vectors.

Definition E.4. A (scalar) random variable X is *subgaussian* with variance proxy σ^2 if the following holds on its moment-generating function:

$$\mathsf{E}[\exp(\lambda X)] \leq \exp\left(\frac{\lambda^2 \sigma^2}{2}\right).$$

A mean-zero random vector $\mathbf{x} \in \mathbb{R}^d$, $\mathsf{E}[\mathbf{x}] = \mathbf{0}$, is *subgaussian* with variance proxy σ^2 if every linear projection is a σ^2 -subgaussian random variable:

$$\mathsf{E}[\exp(\lambda \mathbf{v}^\top \mathbf{x})] \leq \exp\left(\frac{\lambda^2 \|\mathbf{v}\|^2 \sigma^2}{2}\right), \quad \text{for all } \mathbf{v} \in \mathbb{R}^d.$$

With this in hand, we may introduce the subgaussian variant of covariance concentration and the Hanson-Wright inequality.

E.4. Subgaussian Covariance Concentration

We state the subgaussian variant of Lemma E.2, whose proof is structurally the same, replacing the χ^2 random variables with a generic subexponential (c.f. Vershynin (2018, Chapter 2)) random variable, and using a generic Bernstein's inequality rather than the specific χ^2 concentration inequality. The result is qualitatively identical, sacrificing tight/explicit universal numerical constants. The result is relatively standard, and can be found in e.g., Vershynin (2018, Chapter 5) or Du et al. (2021, Lemma A.6).

1980 **Lemma E.5** (Subgaussian covariance concentration). *Let \mathbf{x}_i be i.i.d. zero-mean σ^2 -subgaussian random vectors for
1981 $i = 1, \dots, n$, where $\mathbf{x}_i \in \mathbb{R}^d$, and $\mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \Sigma_{\mathbf{x}}$. Defining the empirical covariance matrix $\widehat{\Sigma}_{\mathbf{x}} \triangleq \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, there
1982 exists a universal constant $C_1 > 0$ such that with probability at least $1 - 2\delta$:*

$$1984 \quad \left\| \widehat{\Sigma}_{\mathbf{x}} - \Sigma_{\mathbf{x}} \right\|_{\text{op}} \leq C\sigma^2 \|\Sigma_{\mathbf{x}}\|_{\text{op}} \left(\sqrt{\frac{d + \log(1/\delta)}{n}} + \frac{d + \log(1/\delta)}{n} \right).$$

1987 *Therefore, as long as $n \geq C_2 \frac{\sigma^2}{c^2} (d + \log(1/\delta))$, we have with probability at least $1 - \delta$,*

$$1989 \quad (1 - c)\Sigma_{\mathbf{x}} \preceq \widehat{\Sigma}_{\mathbf{x}} \preceq (1 + c)\Sigma_{\mathbf{x}}.$$

1991 **E.5. Hanson-Wright Inequality**

1992 We often use the following theorem to prove the concentration inequality for quadratic forms. A modern proof of this
1993 theorem can be found in [Rudelson & Vershynin \(2013\)](#).

1995 **Theorem E.6** (Hanson-Wright Inequality ([Hanson & Wright, 1971](#))). *Let $\mathbf{x} = (X_1, \dots, X_n) \in \mathbb{R}^d$ be a random vector
1996 with independent sub-gaussian components X_i with $\mathbb{E}X_i = 0$. Let \mathbf{D} be an $n \times n$ matrix. Then, for every $t \geq 0$, we have*

$$1998 \quad \mathbb{P} \left\{ |\mathbf{x}^\top \mathbf{D} \mathbf{x} - \mathbb{E}[\mathbf{x}^\top \mathbf{D} \mathbf{x}]| > t \right\} \leq 2 \exp \left[-c \min \left(\frac{t^2}{\|\mathbf{D}\|_F^2}, \frac{t}{\|\mathbf{D}\|_{\text{op}}} \right) \right],$$

2000 *where c is a constant that depends only on the subgaussian constants of X_i .*

2002 **E.6. Stein's Lemma**

2004 We use the following simple lemma which is an application of integration by parts for Gaussian integrals.

2005 **Lemma E.7** (Stein's Lemma). *Let X be a random variables drawn from $N(\mu, \sigma^2)$ and $g : \mathbb{R} \rightarrow \mathbb{R}$ be a differentiable
2006 function. We have $\mathbb{E}[g(X)(X - \mu)] = \sigma^2 \mathbb{E}[g'(X)]$.*

2008 **E.7. Woodbury Matrix Identity**

2009 In the proofs, we use the following elementary identity which states that the inverse of a rank- k correction of a matrix is
2010 equal to a rank- k correction to the inverse of the original matrix.

2012 **Theorem E.8** (Woodbury Matrix Identity ([Woodbury, 1950](#))). *Let $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{C} \in \mathbb{R}^{k \times k}$, $\mathbf{U} \in \mathbb{R}^{n \times k}$, and $\mathbf{V} \in \mathbb{R}^{k \times n}$.
2013 The following matrix identity holds:*

$$2014 \quad (\mathbf{A} + \mathbf{U}\mathbf{C}\mathbf{V})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{C}^{-1} + \mathbf{V}\mathbf{A}^{-1}\mathbf{U})^{-1}\mathbf{V}\mathbf{A}^{-1},$$

2016 *assuming that the inverse matrices in the expression exist.*

2018 **E.8. Stieltjes Transform of Empirical Eigenvalue Distribution**

2020 For a distribution μ over \mathbb{R} , its Stieltjes transform is defined as

$$2022 \quad m_\mu(z) = \int \frac{d\mu(x)}{x - z}.$$

2024 Let H_d be the (discrete) empirical eigenvalue distribution of $\Sigma_{\mathbf{x}} \in \mathbb{R}^{d \times d}$ and let F_d be the (discrete) empirical eigenvalue
2025 distribution of the sample covariance matrix $\widehat{\Sigma} \in \mathbb{R}^{d \times d}$. Consider the proportional limit where $d, n \rightarrow \infty$ with $d/n \rightarrow \phi > 0$.
2026 Suppose that the eigenvalue distribution H_d converges to a limit population spectral distribution $H_{\Sigma_{\mathbf{x}}}$; i.e., $H_d \Rightarrow H_{\Sigma_{\mathbf{x}}}$ in
2027 distribution. Given the definition of $m(z)$ from equation (33), we have $m(z) = m_F(z)$. The following theorem characterizes
2028 m_F in terms of $H_{\Sigma_{\mathbf{x}}}$.

2030 **Theorem E.9** (Silverstein Equation ([Silverstein & Choi, 1995](#))). *Let $\nu_F(z) = \phi(m_F(z) + 1/z) - 1/z$. The function ν_F is
2031 the solution of the following fixed-point equation:*

$$2032 \quad -\frac{1}{\nu_F(z)} = z - \phi \int \frac{t}{1 + t\nu_F(z)} dH_{\Sigma_{\mathbf{x}}}(t).$$

2035 Thus, using this theorem, given H_{Σ_x} , we can numerically compute ν_F (and hence, m_F) using fixed-point iteration. For
 2036 example, for $\Sigma_x^{(s)}$ from equation (14), we have $F = 1/2 \delta_{1-\varepsilon} + 1/2 \delta_{1+\varepsilon}$.
 2037

2038 F. Additional Numerical Results and Details

2039 F.1. Details of the Experiment Setups

2040 In the experiments, we generate $\mathbf{F}_0 \in \mathbb{R}^{d_y \times k}$ with i.i.d. $N(0, 1)$ entries. Then, for each task s we randomly draw a matrix
 2041 $\mathbf{B}_s \in \mathbb{R}^{d_y \times d_y}$ and set $\mathbf{F}_s^s = \exp(0.005(\mathbf{B}_s - \mathbf{B}_s^\top)) \mathbf{F}_0$, where $\exp(\cdot)$ is the matrix exponential. The shared representation
 2042 matrix $\mathbf{G}_* \in \mathbb{R}^{k \times d_x}$ is generated by sampling uniformly from the space of row-orthonormal matrices in $\mathbb{R}^{k \times d_x}$.
 2043

2044 We consider two settings for the covariance matrices $\Sigma_{x,s}$; the *low-anisotropic*, and the *high-anisotropic* settings. In the
 2045 low-anisotropic setting, we define $\mathbf{E} = 5 \mathbf{I}_{d_x} + \mathbf{N}$ where $\mathbf{N} \in \mathbb{R}^{d_x \times d_x}$ has i.i.d. $N(0, 1)$ entries, and set $\Sigma_{x,s} = 0.5 (\mathbf{E} + \mathbf{E}^\top)$.
 2046 For the high-anisotropic setting, we first sample uniformly a rotation matrix $\mathbf{O} \in \mathbb{R}^{d_x \times d_x}$ and set $\Sigma_{x,s} = \mathbf{O} \mathbf{D} \mathbf{O}^\top$ where
 2047 $\mathbf{D} = \text{diag}(\text{logspace}(0, 5, d_x))$. In the experiments for the main paper, we always consider the high-anisotropic setting.
 2048

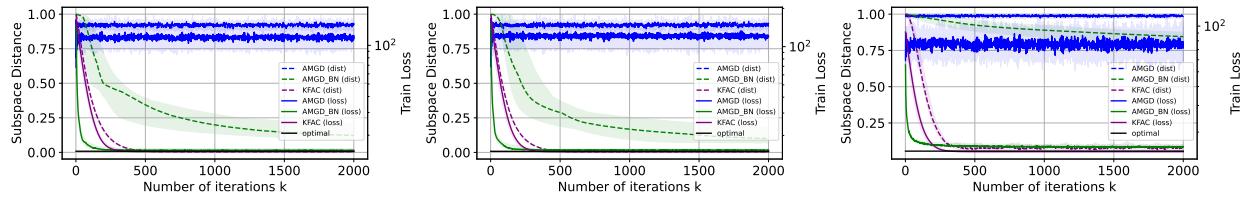
2049 In the following experiments, in addition to the data generation process in equation (13) used in the experiment in the main
 2050 paper, we also consider a Gaussian data setup where samples for task s are generated according to
 2051

$$2053 \quad \mathbf{y}_i^s = \mathbf{F}_s^s \mathbf{G}_* \mathbf{x}_i^s + \varepsilon_i^s, \quad \mathbf{x}_i^s \sim N(\mathbf{0}, \Sigma_{x,s}), \quad \varepsilon_i^s \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma_{\varepsilon,s} \mathbf{I}_{d_y}), \quad s \in \{\text{test, train}\}. \quad (38)$$

2056 F.2. Additional Experiments

2057 F.2.1. EFFECT OF BATCH NORMALIZATION

2058 We used the same experiment setting described in Section 4 to generate the plots in Figure 9. Explicitly, we use data
 2059 dimension $d_x = 100$, task dimension $d_y = 15$, and representation dimension $k = 8$. We use the same learning rate 10^{-2} for
 2060 each optimizer except for NGD, in which we used 10^{-4} . The batch size is 1024. In Figure 3 we considered the Uniform data
 2061 (13) with high anisotropy. Here, we consider the other three setting: Uniform data (13) with low anisotropy, Gaussian data
 2062 (38) with low anisotropy, and Gaussian data (38) with high anisotropy.
 2063



2073 *Figure 4.* The effect of batch normalization (on AMGD) vs. KFAC in our experiment settings (**Left**) Uniform with low anisotropy. (**Middle**)
 2074 Gaussian with low anisotropy. (**Right**) Gaussian with high anisotropy.

2075 As discussed in the main paper Section 3.1.1, we expect AMGD with batch-norm to converge in training loss but to perform
 2076 poorly with respect to the subspace distance from the optimal in settings in the case with high anisotropy (**Right**). However,
 2077 in the experiment settings with low anisotropy (**Left and Center**), we expect reasonable performance from this algorithm
 2078 because $\text{rowsp}(\mathbf{G}_* \Sigma_x)$ is close to the target $\text{rowsp}(\mathbf{G}_*)$.

2079 F.2.2. LEARNING RATE SWEEP

2080 We further test the performance of each learning algorithm at different learning rates from $10^{-6}, 10^{-5.5}, \dots, 10^{-0.5}, 10^0$,
 2081 with results shown in Figure 5, where we plot the subspace distance at 1000 iterations for different algorithms. If the
 2082 algorithm encounters numerical instability, then we report the subspace distance as the maximal value of 1.0. We observe
 2083 that KFAC and DFW coverage to a solution with small subspace distance to the true representation for a wide range of step
 2084 sizes, whereas the set of suitable learning rates for other algorithms is much narrower. Furthermore, we observe the poor
 2085 performance of various algorithms in Figure 1 and Figure 9 is not due to specific choice of learning rate.
 2086

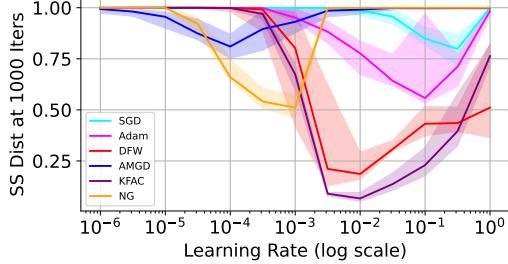


Figure 5. The subspace distance of representations learned by different algorithms after 1000 iterations and the true representation as a function of learning rate.

F.2.3. HEAD-TO-HEAD EXPERIMENTS

We again consider the same experimental setting used for Figure 1. In particular, we use data dimension $d_x = 100$, task dimension $d_y = 15$, and representation dimension $k = 8$. We use the same learning rate 10^{-2} for each optimizer except for NGD optimizer, in which we used 10^{-4} . The batch size is 1024. In Figure 1 we considered the Uniform data (13) with high anisotropy. Here, we consider the other three setting: Uniform data (13) with low anisotropy, Gaussian data (38) with low anisotropy, and Gaussian data (38) with high anisotropy. We plot the training loss, subspace distance to the ground truth shared representation, and the transfer loss obtained by different algorithms. See Figure 9. We observe that in all three settings, various algorithms converge in training loss. In the case with high anisotropy (second row), methods other than KFAC do not converge to the optimal representation in subspace distance and transfer loss. However, in the low anisotropy settings (first and third rows), the performance of other algorithms also improve, but are notably still suboptimal relative to KFAC, confirming the theoretical results showing that anisotropy is a root cause behind the sub-optimality of prior algorithms and analysis.

2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144

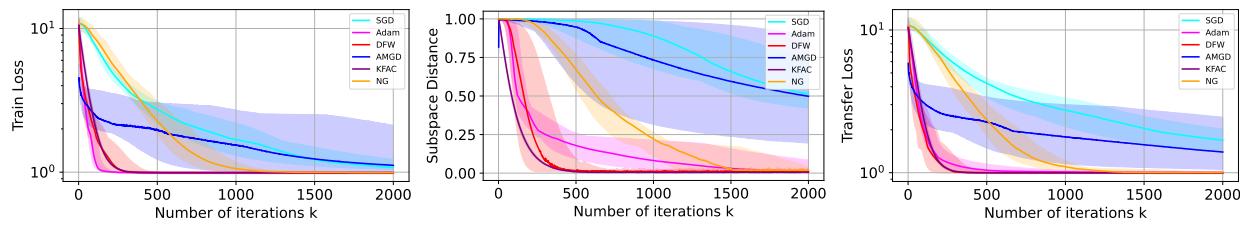


Figure 6. Gaussian with low anisotropy

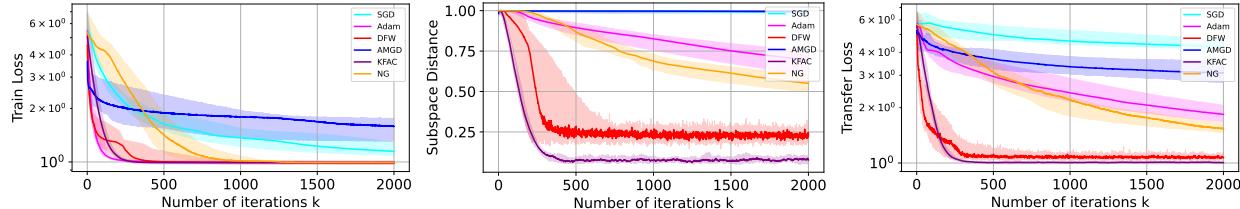


Figure 7. Gaussian with high anisotropy

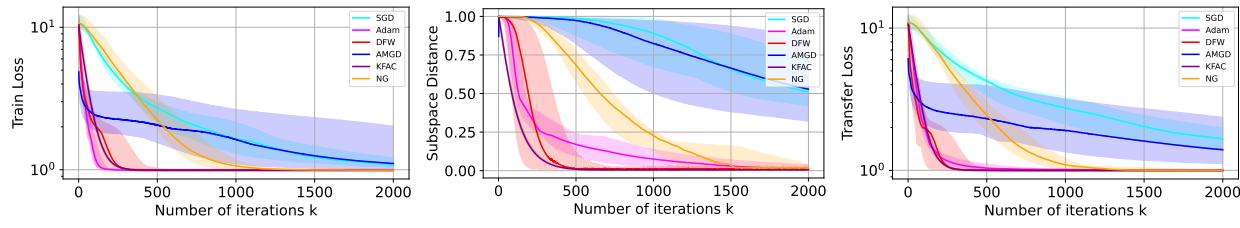


Figure 8. Bernoulli with low anisotropy

Figure 9. From **left to right**: the training loss, subspace distance, and transfer loss induced by various algorithms on a linear representation learning task.