

WHEN DOES PRECONDITIONING HELP OR HURT GENERALIZATION?

Anonymous authors

Paper under double-blind review

ABSTRACT

While second order optimizers such as natural gradient descent (NGD) often speed up optimization, their effect on generalization has been called into question. This work presents a more nuanced view on how the *implicit bias* of optimizers affects the comparison of generalization properties. We provide an exact bias-variance decomposition of the generalization error of overparameterized ridgeless regression under a general class of preconditioner \mathbf{P} , and consider the inverse population Fisher information matrix (used in NGD) as a particular example. We determine the optimal \mathbf{P} for both the bias and variance, and find that the relative generalization performance of different optimizers depends on label noise and “shape” of the signal (true parameters): when the labels are noisy, the model is misspecified, or the signal is misaligned, NGD can achieve lower risk; conversely, GD generalizes better under clean labels, a well-specified model, or aligned signal. Based on this analysis, we discuss approaches to manage the bias-variance tradeoff, and the benefit of interpolating between first- and second-order updates. We then extend our analysis to regression in the reproducing kernel Hilbert space and demonstrate that preconditioned GD can decrease the population risk faster than GD. Lastly, we empirically compare the generalization error of first- and second-order optimizers in neural network, and observe robust trends matching our theoretical analysis.

1 INTRODUCTION

We study the generalization property of an estimator $\hat{\theta}$ obtained by minimizing the empirical risk (or the training error) $L(f_{\theta})$ via a preconditioned gradient update with preconditioner \mathbf{P} :

$$\theta_{t+1} = \theta_t - \eta \mathbf{P}(t) \nabla_{\theta_t} L(f_{\theta_t}), \quad t = 0, 1, \dots \quad (1.1)$$

Setting $\mathbf{P} = \mathbf{I}$ recovers gradient descent (GD). Choices of \mathbf{P} which exploit second-order information include the inverse Fisher information matrix, which gives the natural gradient descent (NGD) (Amari, 1998); the inverse Hessian, which leads to Newton’s method; and diagonal matrices estimated from past gradients, corresponding to adaptive gradient methods (Duchi et al., 2011; Kingma & Ba, 2014). These preconditioners often alleviate the effect of pathological curvature and speed up *optimization*, but their *generalization* properties has been under debate. While several works reported that in neural network optimization, adaptive or second-order methods generalize worse compared to gradient descent (GD) (Wilson et al., 2017), other empirical studies suggested that second-order methods can achieve comparable, if not better generalization (Xu et al., 2020).

The generalization property of optimizers relates to the discussion of *implicit bias* (Gunasekar et al., 2018a), i.e. preconditioning may lead to a different converged solution (with same training loss), as shown in Figure 1. While many explanations have been proposed, our starting point is the well-known observation that GD implicitly regularizes the parameter ℓ_2 norm. For instance in overparameterized least squares regression, GD and many first-order methods find the minimum ℓ_2 norm solution from zero initialization (without explicit regularization), but preconditioned updates often do not. This being said, while the minimum norm solution may generalize well in the overparameterized regime (Bartlett et al., 2019), it is unclear whether preconditioning leads to inferior solutions – even in the simple setting of overparameterized linear regression, *quantitative* understanding of how preconditioning affects generalization is large lacking.

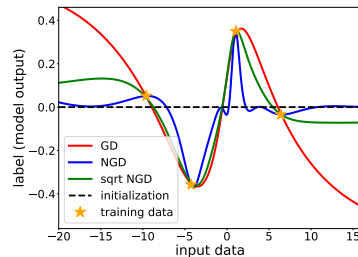


Figure 1: 1D illustration of different implicit biases: two-layer sigmoid network trained with preconditioned GD.

Motivated by the observations above, in Section 3 we start with overparameterized least squares regression (unregularized) and analyze the stationary solution ($t \rightarrow \infty$) of update (1.1) under time-invariant preconditioner. Extending previous analysis in the proportional limit (Hastie et al., 2019), we consider a more general random effects model and derive the exact population risk in its *bias-variance decomposition* via random matrix theory. We characterize the optimal \mathbf{P} within a general class of preconditioners for both the bias and variance, and focus on the comparison between GD, for which \mathbf{P} is identity, and NGD, for which \mathbf{P} is the inverse population Fisher information matrix¹. We find that the comparison of generalization performance is affected by the following factors:

1. **Label Noise:** Additive noise in the labels leads to the *variance* term in the risk. We show that NGD achieves the optimal variance among a general class of preconditioned updates.
2. **Model Misspecification:** Under misspecification, there does not exist f_θ that perfectly learns the true function (target). This contributes to the *bias* term, and we argue that its effect is similar to label noise, and thus NGD may also be beneficial when the model is misspecified.
3. **Data-Signal-Alignment:** Alignment describes how the target signal distributes among input features² and affects the *bias* term. GD achieves lower bias for isotropic signal, whereas NGD is preferred under “misalignment” — when the target function focuses on small feature directions.

Beyond the decomposition of stationary risk, our findings in Section 4 and 5 are summarized as:

- In Section 4.1 and 4.2 we discuss how the bias-variance tradeoff can be realized by different choices of preconditioner \mathbf{P} (e.g. interpolating between GD and NGD) or early stopping.
- In Section 4.3 we extend our analysis to regression in the RKHS and show that under early stopping, a preconditioned update interpolating between GD and NGD achieves minimax optimal convergence rate in much fewer steps, and thus reduces the population risk faster than GD.
- In Section 5 we empirically test how well our findings in linear model carry over to neural networks: under a student-teacher setup, we compare the generalization of GD with preconditioned updates and illustrate the influence of all aforementioned factors. The performance of neural networks under a variety of manipulations results in trends that align with our theoretical analysis.

2 BACKGROUND AND RELATED WORKS

Natural Gradient Descent. NGD is a second-order method originally proposed in Amari (1997). Consider a data distribution $p(\mathbf{x})$ on the space \mathcal{X} , a function $f_\theta : \mathcal{X} \rightarrow \mathcal{Z}$ parameterized by θ , and a loss function $L(\mathbf{X}, f_\theta) = \frac{1}{n} \sum_{i=1}^n l(y_i, f_\theta(\mathbf{x}_i))$, where $l : \mathcal{Y} \times \mathcal{Z} \rightarrow \mathbb{R}$. Also suppose a probability distribution $p(y|\mathbf{z}) = p(y|f_\theta(\mathbf{x}))$ is defined on the space of labels. Then, the natural gradient is defined as: $\tilde{\nabla}_\theta L(\mathbf{X}, f_\theta) = \mathbf{F}^{-1} \nabla_\theta L(\mathbf{X}, f_\theta)$, where $\mathbf{F} = \mathbb{E}[\nabla_\theta \log p(\mathbf{x}, y|\theta) \nabla_\theta \log p(\mathbf{x}, y|\theta)^\top]$ is the *Fisher information matrix*, or simply the (population) Fisher. Note that expectations in the Fisher are under the joint distribution of the model $p(\mathbf{x}, y|\theta) = p(\mathbf{x})p(y|f_\theta(\mathbf{x}))$. In the literature, the Fisher is sometimes defined under the empirical data distribution $\{\mathbf{x}_i\}_{i=1}^n$ (Amari et al., 2000). We instead refer to this quantity as the *sample Fisher*, the properties of which influence optimization and have been studied in Karakida et al. (2018); Kunstner et al. (2019); Thomas et al. (2020). We remark that in linear and kernel regression under squared loss, sample Fisher-based updates give the same stationary solution as GD (see Section 3), whereas population Fisher-based update may not.

While the population Fisher is typically difficult to obtain, extra unlabeled data can be used in its estimation, which empirically improves generalization (Pascanu & Bengio, 2013). Moreover, under structural assumptions, parametric approaches to estimate \mathbf{F} can be more sample-efficient (Martens & Grosse, 2015; Ollivier, 2015), and thus closing the gap between sample and population Fisher.

When the per-instance loss is the negative log-probability of an exponential family, the sample Fisher coincides with the *generalized Gauss-Newton matrix* (Martens, 2014). In least squares regression, which is the focus of this work, the quantity also coincides with the Hessian. We thus take NGD as a representative example of preconditioned update, and we expect our findings to also translate to other second-order methods (not including adaptive gradient methods) in regression problems.

¹From now on we use NGD to denote the *population* Fisher-based update, and we write “sample NGD” when \mathbf{P} is the inverse or pseudo-inverse of the sample Fisher; see Section 2 for discussion.

²Our notion of alignment relates to the *source condition* in RKHS (Cucker & Smale, 2002) (see Section 4.3).

Analysis of Preconditioned Gradient Descent. While Wilson et al. (2017) outlined one example under fixed training data where GD generalizes better than adaptive methods, in the online learning setting, for which optimization speed relates to generalization, several works have shown the advantage of preconditioning (Levy & Duchi, 2019; Zhang et al., 2019a). In addition, (Zhang et al., 2019b; Cai et al., 2019) established convergence and generalization guarantees of sample Fisher-based updates were derived for neural networks in the kernel regime. Lastly, the generalization of different optimizers also connects to the notion of “sharpness” (Keskar et al., 2016; Dinh et al., 2017), and it has been argued that second-order updates tend to find sharper minima (Wu et al., 2018).

We note that two concurrent works also discussed the generalization performance of preconditioned updates. Wadia et al. (2020) connected second-order methods with data whitening in linear models, and qualitatively showed that whitening (thus second-order update) harms generalization in certain cases. Vaswani et al. (2020) analyzed the complexity of the maximum P -margin solution in linear classification problems. We emphasize that instead of *upper bounding* the risk (e.g. Rademacher complexity), which may not decide the optimal P (for generalization), we compute the *exact risk* for least squares regression, which allows us to precisely compare different preconditioners.

3 ASYMPTOTIC RISK OF RIDGELESS INTERPOLANTS

In this section we consider the following setup: given n training samples $\{\mathbf{x}_i\}_{i=1}^n$ labeled by a teacher model (target function) $f^*: \mathbb{R}^d \rightarrow \mathbb{R}$ with additive noise: $y_i = f^*(\mathbf{x}_i) + \varepsilon_i$, we learn a linear student model f_θ by minimizing the squared loss: $L(\mathbf{X}, f_\theta) = \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \theta)^2$. We assume a random design: $\mathbf{x}_i = \Sigma_X^{1/2} \mathbf{z}_i$, where $\mathbf{z}_i \in \mathbb{R}^d$ is an i.i.d. vector with zero-mean, unit-variance, and finite 12th moment, and ε is i.i.d. noise independent to \mathbf{z} with mean 0 and variance σ^2 . Our goal is to compute the population risk $R(f) = \mathbb{E}_{P_X}[(f^*(\mathbf{x}) - f(\mathbf{x}))^2]$ in the proportional asymptotic limit:

- **(A1) Overparameterized Proportional Limit:** $n, d \rightarrow \infty, d/n \rightarrow \gamma \in (1, \infty)$.

(A1) entails that the number of features (or parameters) is larger than the number of samples. In this overparameterized setting, the population risk is equivalent to the generalization error, and there exist multiple empirical risk minimizers with potentially different generalization properties.

Denote $\mathbf{X} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top]^\top \in \mathbb{R}^{n \times d}$ the data matrix and $\mathbf{y} \in \mathbb{R}^n$ the corresponding label vector. We optimize the parameters θ via a preconditioned gradient flow with preconditioner $\mathbf{P}(t) \in \mathbb{R}^{d \times d}$,

$$\frac{\partial \theta(t)}{\partial t} = -\mathbf{P}(t) \frac{\partial L(\theta(t))}{\partial \theta(t)} = \frac{1}{n} \mathbf{P}(t) \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \theta(t)), \quad \theta(0) = 0. \quad (3.1)$$

In this linear setup, many common choices of preconditioner do not change through time: under Gaussian likelihood, the sample Fisher (and also Hessian) corresponds to the sample covariance $\mathbf{X}^\top \mathbf{X} / n$ up to variance scaling, whereas the population Fisher corresponds to the population covariance $\mathbf{F} = \Sigma_X$. We thus limit our analysis to fixed preconditioner of the form $\mathbf{P}(t) =: \mathbf{P}$.

Write parameters at time t under update (3.1) with fixed \mathbf{P} as $\theta_P(t)$. For positive definite \mathbf{P} , the stationary solution is given as: $\hat{\theta}_P := \lim_{t \rightarrow \infty} \theta_P(t) = \mathbf{P} \mathbf{X}^\top (\mathbf{X} \mathbf{P} \mathbf{X}^\top)^{-1} \mathbf{y}$. One may check that the discrete time gradient descent update (with appropriate step size) and other variants that do not alter the span of gradient (e.g. stochastic gradient or momentum) converge to the same solution as well.

Intuitively speaking, if the data distribution (blue contour in Figure 2) is not isotropic, then certain directions will be more “important” than others. In this case uniform ℓ_2 shrinkage (which GD implicitly provides) may not be the most desirable, and certain \mathbf{P} that takes data geometry into account may lead to better generalization instead. The above intuition will be made rigorous in this section.

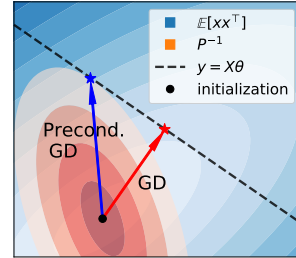


Figure 2: Geometric illustration (2D) of how the interpolating θ_P depends on the preconditioner.

Remark. $\hat{\theta}_P$ is the minimum $\|\theta\|_{\mathbf{P}^{-1}}$ norm interpolant: $\hat{\theta}_P = \arg \min_{\theta} \|\theta\|_{\mathbf{P}^{-1}}, \text{ s.t. } \mathbf{X} \theta = \mathbf{y}$ for positive definite \mathbf{P} . For GD this translates to the parameter ℓ_2 norm, whereas for NGD ($\mathbf{P} = \mathbf{F}^{-1} = \Sigma_X^{-1}$), the implicit bias is the $\|\theta\|_{\mathbf{F}}$ norm. Since $\mathbb{E}_{P_X}[f(\mathbf{x})^2] = \|\theta\|_{\Sigma_X}^2$, NGD finds an interpolating function with smallest norm under the data distribution. We empirically observe this divide between small parameter norm and function norm in neural networks as well (Figure 1 and Appendix A.1).

We highlight the following choices of \mathbf{P} and the corresponding stationary solution $\hat{\boldsymbol{\theta}}_{\mathbf{P}}$ as $t \rightarrow \infty$.

- **Identity:** $\mathbf{P} = \mathbf{I}_d$ recovers GD that converges to the min ℓ_2 norm interpolant (also true for momentum GD and SGD), which we write as $\hat{\boldsymbol{\theta}}_{\mathbf{I}} := \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{y}$ and refer to as the *GD solution*.
- **Population Fisher:** $\mathbf{P} = \mathbf{F}^{-1} = \boldsymbol{\Sigma}_{\mathbf{X}}^{-1}$ leads to the estimator $\hat{\boldsymbol{\theta}}_{\mathbf{F}^{-1}}$, which we refer to as the *NGD solution*.
- **Sample Fisher:** since the sample Fisher is rank-deficient, we may add a damping term $\mathbf{P} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1}$ or take the pseudo-inverse $\mathbf{P} = (\mathbf{X}^\top \mathbf{X})^\dagger$. In both cases, the gradient is still spanned by \mathbf{X} , and thus the update finds the same min ℓ_2 -norm solution $\hat{\boldsymbol{\theta}}_{\mathbf{I}}$ (also true for full-matrix Adagrad (Agarwal et al., 2018)), although the trajectory differs (see Figure 3).

Remark. The above choices reveal a gap between sample- and population-based \mathbf{P} : while the sample Fisher accelerates optimization (Zhang et al., 2019b), the following sections demonstrate generalization properties only possessed by the population Fisher.

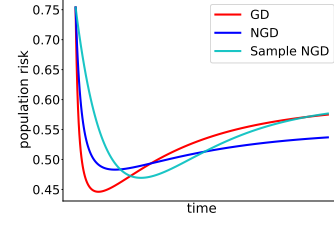


Figure 3: Population risk of preconditioned linear regression vs. time with the following \mathbf{P} : \mathbf{I} (red), $\boldsymbol{\Sigma}_{\mathbf{X}}^{-1}$ (blue) and $(\mathbf{X}^\top \mathbf{X})^\dagger$ (cyan). Time is rescaled differently for each curve (convergence speed is not comparable). Observe that GD and sample NGD give the same stationary risk.

We compare the population risk of the GD solution $\hat{\boldsymbol{\theta}}_{\mathbf{I}}$ and NGD solution $\hat{\boldsymbol{\theta}}_{\mathbf{F}^{-1}}$ in its bias-variance decomposition w.r.t. label noise (Hastie et al., 2019) and discuss the two components separately,

$$R(\boldsymbol{\theta}) = \underbrace{\mathbb{E}_{P_{\mathbf{X}}} [(f^*(\mathbf{x}) - \langle \mathbf{x}, \mathbb{E}_{P_{\varepsilon}}[\boldsymbol{\theta}] \rangle)^2]}_{B(\boldsymbol{\theta}), \text{ bias}} + \underbrace{\text{tr}(\text{Cov}(\boldsymbol{\theta}) \boldsymbol{\Sigma}_{\mathbf{X}})}_{V(\boldsymbol{\theta}), \text{ variance}}.$$

Note that the *bias* term does not depend on the label noise ε , whereas the *variance* term does not depend on the teacher f^* . In addition, given that f^* can be independently decomposed into a linear component on the features \mathbf{x} and a residual term: $f^*(\mathbf{x}) = \langle \mathbf{x}, \boldsymbol{\theta}^* \rangle + f_c^*(\mathbf{x})$, we can further decompose the bias into a *well-specified* component $\|\boldsymbol{\theta}^* - \mathbb{E}\boldsymbol{\theta}\|_{\boldsymbol{\Sigma}_{\mathbf{X}}}^2$, which captures the difficulty in learning the linear component of f^* , and a *misspecified* component (Hastie et al., 2019, Section 5), which corresponds to the error in fitting f_c^* (beyond the class of functions the student can represent).

3.1 THE VARIANCE TERM: NGD IS OPTIMAL

We first characterize the stationary variance which is independent to the teacher model f^* . We restrict ourselves to preconditioners satisfying the following assumption on the spectral distribution:

- **(A2) Converging Eigenvalues:** \mathbf{P} is positive definite and as $n, d \rightarrow \infty$, the spectral distribution of $\boldsymbol{\Sigma}_{\mathbf{X}\mathbf{P}} := \mathbf{P}^{1/2} \boldsymbol{\Sigma}_{\mathbf{X}} \mathbf{P}^{1/2}$ converges weakly to $\mathbf{H}_{\mathbf{X}\mathbf{P}}$ supported on $[c, C]$ for $c, C > 0$.

The following theorem characterizes the asymptotic variance and the corresponding optimal \mathbf{P} .

Theorem 1. Given (A1-2), the asymptotic variance is given as

$$V(\hat{\boldsymbol{\theta}}_{\mathbf{P}}) \rightarrow \sigma^2 \left(\lim_{\lambda \rightarrow 0+} m'(-\lambda) m^{-2}(-\lambda) - 1 \right), \quad (3.2)$$

where $m(z) > 0$ is the Stieltjes transform of the limiting distribution of eigenvalues of $\frac{1}{n} \mathbf{X} \mathbf{P} \mathbf{X}^\top$ (for z beyond its support) defined as the solution to $m^{-1}(z) = -z + \gamma \int \tau (1 + \tau m(z))^{-1} d\mathbf{H}_{\mathbf{X}\mathbf{P}}(\tau)$.

Furthermore, under (A1-2), $V(\hat{\boldsymbol{\theta}}_{\mathbf{P}}) \geq \sigma^2(\gamma - 1)^{-1}$, and the equality is achieved by $\mathbf{P} = \mathbf{F}^{-1}$.

Formula (3.2) is a direct extension of Hastie et al. (2019, Theorem 4), which can be obtained from Dobriban et al. (2018, Theorem 2.1) or Ledoit & P  ch   (2011, Theorem 1.2). We note that the eigenvalue condition in (A2) may also be relaxed as in Xu & Hsu (2019). Theorem 1 implies that preconditioning with the inverse population Fisher \mathbf{F} results in the optimal stationary variance, which is supported by Figure 4(a). In other words, when the labels are noisy so that the risk is dominated by the variance term, we expect NGD to generalize better upon convergence. We emphasize that this advantage is only present when the population Fisher is used, but not its sample-based counterpart (which converges to $\hat{\boldsymbol{\theta}}_{\mathbf{I}}$ as commented above). In Appendix A.3 we discuss the substitution error in replacing the population Fisher \mathbf{F} with a sample-based estimate based on unlabeled data.

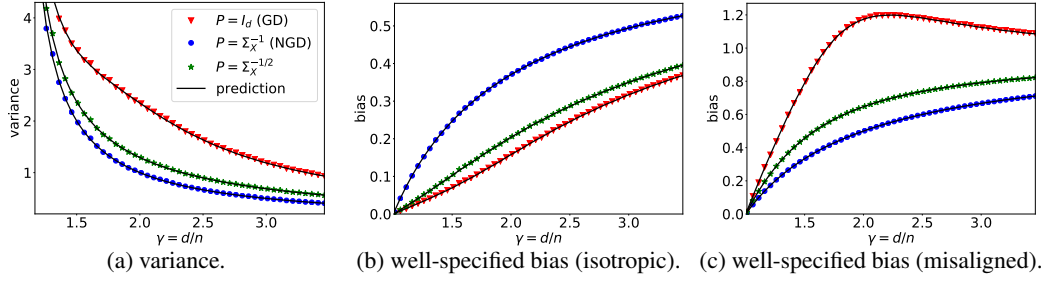


Figure 4: We set eigenvalues of Σ_X as two point masses with $\kappa_X = 20$ and $\|\Sigma_X\|_F^2 = d$; empirical values (dots) are computed with $n = 300$. (a) NGD (blue) achieves minimum variance. (b) GD (red) achieves lower bias under isotropic signal: $\Sigma_\theta = I_d$. (c) NGD achieves lower bias under “misalignment”: $\Sigma_X = \Sigma_\theta^{-1}$.

3.2 THE BIAS TERM: WELL-SPECIFIED CASE

We now analyze the bias term when the teacher model is linear on the input features (hence well-specified): $f^*(x) = x^\top \theta^*$. Extending the random effects hypothesis in Dobriban et al. (2018), we consider a more general prior on θ^* : $\mathbb{E}[\theta^* \theta^{*\top}] = d^{-1} \Sigma_\theta$, and assume the following joint relations³:

- **(A3) Joint Convergence:** Σ_X and P share the same eigenvector matrix U . The empirical distributions of elements of (e_x, e_θ, e_{xp}) jointly converge to random variables (v_x, v_θ, v_{xp}) supported on $[c', C']$ for $c', C' > 0$, where e_x, e_{xp} are eigenvalues of $\Sigma_X, \Sigma_X P$, and $e_\theta = \text{diag}(U^\top \Sigma_\theta U)$.

We remark that when $P = I_d$, previous works (Hastie et al., 2019; Xu & Hsu, 2019) considered the special case of isotropic prior $\Sigma_\theta = I_d$. Our assumption thus allows for analysis of the bias term under much more general Σ_θ , which gives rise to interesting phenomena that are not captured by simplified settings, such as non-monotonic bias and variance for $\gamma > 1$ (see Figure 15), and the epoch-wise double descent phenomenon (see Appendix A.2). Under this general setup, we have the following characterization of the asymptotic bias and the corresponding optimal preconditioner:

Theorem 2. Under (A1)(A3), the expected bias $B(\hat{\theta}_P) := \mathbb{E}_{\theta^*}[B(\hat{\theta}_P)]$ is given as

$$B(\hat{\theta}_P) \rightarrow \lim_{\lambda \rightarrow 0_+} m'(-\lambda) m^{-2}(-\lambda) \mathbb{E}[v_x v_\theta (1 + v_{xp} m(-\lambda))^{-2}], \quad (3.3)$$

where expectation is taken over v and $m(z)$ is the Stieltjes transform defined in Theorem 1.

Furthermore, among all P satisfying (A3), the optimal bias is achieved by $P = U \text{diag}(e_\theta) U^\top$.

Note that the optimal P depends on the “orientation” of the teacher model Σ_θ , which is usually not known in practice. This result can thus be interpreted as a *no-free-lunch* characterization in choosing an optimal preconditioner for the bias term *a priori*. As a consequence of the theorem, when the true parameters θ^* have roughly equal magnitude (isotropic), GD achieves lower bias (see Figure 4(b) where $\Sigma_\theta = I_d$). On the other hand, NGD leads to lower bias when Σ_X is “misaligned” with Σ_θ , i.e. when θ^* focus on the least varying directions of input features (see Figure 4(c) where $\Sigma_\theta = \Sigma_X^{-1}$), in which case learning is intuitively difficult since the features are not useful⁴.

In the analogy of source condition (similar to $\mathbb{E}[\Sigma_X^{-r/2} \theta^*] < \infty$), the coefficient r can be interpreted as a measure of “misalignment”, which relates to the hardness of learning (see Section 4.3); the above discussion thus suggests that GD may be beneficial in “easy” tasks, whereas NGD is preferable when the teacher is “difficult”. In Appendix A.4 we elaborate this connection by analyzing the setting where $\Sigma_\theta = \Sigma_X^{-r}$, and provide a more precise comparison between the well-specified bias of GD and NGD in certain special cases.

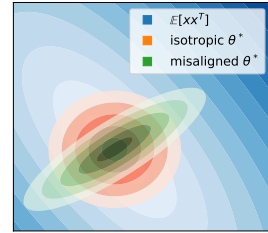


Figure 5: 2D illustration of isotropic (red) and misaligned (green) teacher θ^* .

³Note that (A2)(A3) covers many common choices of preconditioner, such as the population Fisher and variants of the sample Fisher (which is degenerate but leads to the same minimum ℓ_2 norm solution as GD).

⁴If θ^* is neither isotropic nor completely misaligned with x , then the optimal bias is not achieved by GD or NGD. Qualitatively speaking, we expect GD or NGD to be beneficial when Σ_θ is “closer” to the isotropic or fully misaligned case, respectively. See Appendix A.4 for more discussion.

3.3 MISSPECIFICATION \approx LABEL NOISE

Finally, we consider the additional bias introduced by model misspecification, under which there does not exist a linear student that perfectly recovers the teacher model f^* . In this case, we may decompose the teacher into a linear component and its residual: $f^*(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\theta}^* + f_c^*(\mathbf{x})$.

For simplicity, we first consider f_c^* to be a linear function on unobserved features (similar to [Hastie et al. \(2019\)](#)): $y_i = \mathbf{x}_i^\top \boldsymbol{\theta}^* + \mathbf{x}_{c,i}^\top \boldsymbol{\theta}^c + \varepsilon_i$, where $\mathbf{x}_{c,i} \in \mathbb{R}^{d_c}$ is independent to \mathbf{x}_i with covariance $\Sigma_{\mathbf{X}}^c$, and $\mathbb{E}[\boldsymbol{\theta}^c \boldsymbol{\theta}^{c\top}] = d_c^{-1} \Sigma_{\boldsymbol{\theta}}^c$. In this setting, the misspecified bias behaves the same as the variance:

Proposition 3. *For the above unobserved features model, given (A1-3), the bias can be written as $B(\hat{\boldsymbol{\theta}}) = B_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}_{\mathbf{P}}) + B_c(\hat{\boldsymbol{\theta}}_{\mathbf{P}})$, where $B_{\boldsymbol{\theta}}$ is the well-specified bias in Thm. 2, and $B_c = d_c^{-1} \text{tr}(\Sigma_{\mathbf{X}}^c \Sigma_{\boldsymbol{\theta}}^c) (V(\hat{\boldsymbol{\theta}}_{\mathbf{P}}) + 1)$, where $V(\hat{\boldsymbol{\theta}}_{\mathbf{P}})$ is the variance in Thm. 1.*

Misspecification can thus be interpreted as additional label noise, for which NGD is advantageous by Theorem 1. While Proposition 3 describes one specific example of misspecification, we expect such characterization to hold under broader settings. In particular, [Mei & Montanari \(2019, Remark 5\)](#) indicates that for many nonlinear f_c^* , the misspecified bias is the same as variance due to label noise. This result is only shown for isotropic data, but we empirically observe similar phenomenon under general covariances in Figure 6, in which f_c^* is a quadratic function. Observe that NGD leads to lower bias compared to GD as we further misspecify the teacher model.

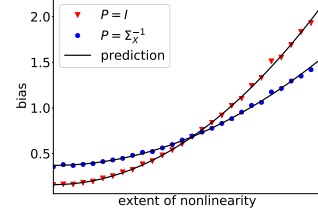


Figure 6: Misspecified bias: we take $\Sigma_{\boldsymbol{\theta}} = \mathbf{I}_d$ (favors GD) and $f_c^*(\mathbf{x}) = \alpha(\mathbf{x}^\top \mathbf{x} - \text{tr}(\Sigma_{\mathbf{X}}))$, where α controls the extent of nonlinearity. Predictions are generated by matching the noise level with 2nd moment of f_c^* .

4 BIAS-VARIANCE TRADEOFF

Our characterization of the stationary risk suggests that preconditioners that achieve the optimal bias and variance are in general different. This section discusses how the bias-variance tradeoff can be realized by interpolating between preconditioners or by early stopping. In addition, we analyze regression in the RKHS and show that by balancing the bias and variance, a preconditioned update that interpolates between GD and NGD also decreases the population risk faster than GD.

4.1 INTERPOLATING BETWEEN PRECONDITIONERS

Depending on the orientation of the teacher model, we may expect a bias-variance tradeoff in choosing \mathbf{P} . Intuitively, given \mathbf{P}_1 that minimizes the bias and \mathbf{P}_2 that minimizes the variance, it is possible that a preconditioner interpolating between \mathbf{P}_1 and \mathbf{P}_2 can balance the bias and variance and thus generalize better. The following proposition confirms this intuition in a setup of general $\Sigma_{\mathbf{X}}$ and isotropic $\Sigma_{\boldsymbol{\theta}}$ ⁵, for which GD ($\mathbf{P} = \mathbf{I}_d$) achieves optimal stationary bias and NGD ($\mathbf{P} = \mathbf{F}^{-1}$) achieves optimal variance.

Proposition 4 (Informal). *Let $\Sigma_{\mathbf{X}} \neq \mathbf{I}_d$ and $\Sigma_{\boldsymbol{\theta}} = \mathbf{I}_d$. Consider the following choices of interpolation scheme: (i) $\mathbf{P}_\alpha = \alpha \Sigma_{\mathbf{X}}^{-1} + (1-\alpha)\mathbf{I}_d$, (ii) $\mathbf{P}_\alpha = (\alpha \Sigma_{\mathbf{X}} + (1-\alpha)\mathbf{I}_d)^{-1}$, (iii) $\mathbf{P}_\alpha = \Sigma_{\mathbf{X}}^{-\alpha}$. The stationary variance monotonically decreases with $\alpha \in [0, 1]$ for all three choices. For (i), the stationary bias monotonically increases with $\alpha \in [0, 1]$, whereas for (ii) and (iii), the bias monotonically increases with α in certain range that depends on $\Sigma_{\mathbf{X}}$.*

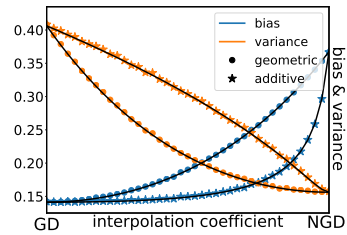


Figure 7: Bias-variance tradeoff with $\kappa_{\mathbf{X}} = 25$, $\Sigma_{\boldsymbol{\theta}} = \mathbf{I}_d$ and $\text{SNR}=32/5$. As we additively (ii) or geometrically (iii) interpolate from GD to NGD (left to right), the stationary bias (blue) increases and the stationary variance (orange) decreases.

In other words, as the signal-to-noise ratio (SNR) decreases (i.e. more label noise), one can increase α , which makes the update closer to NGD, to improve generalization, and vice versa⁶ (small α entails GD-like update). This intuition is supported by Figure 7 and 16(c): for certain SNR, interpolating between $\Sigma_{\mathbf{X}}^{-1}$ and $\Sigma_{\boldsymbol{\theta}}$ can lead to lower stationary risk than both GD and NGD.

⁵Note that this reduces to the random effects model studied in [Dobriban et al. \(2018\)](#).

⁶In Appendix D.5 we empirically verify the monotonicity bias over all $\alpha \in [0, 1]$ beyond the proposition.

Remark. Two of the aforementioned interpolation schemes are analogous to common choices in practice: the additive interpolation (ii) corresponds to the damping term to stably invert the Fisher, whereas the geometric interpolation (iii) includes the “conservative” square-root scaling in adaptive gradient methods (Duchi et al., 2011; Kingma & Ba, 2014).

4.2 THE ROLE OF EARLY STOPPING

Thus far we considered the stationary solution of the unregularized objective. It is known that the bias-variance tradeoff can also be controlled by either explicit or algorithmic regularization. We briefly comment on the effect of early stopping, starting from the monotonicity of the variance term.

Proposition 5. For all P satisfying (A2), the variance $V(\theta_P(t))$ monotonically increases with time.

The proposition confirms the intuition that early stopping reduces overfitting. Variance reduction can benefit GD in its comparison to NGD, which achieves lowest stationary variance: indeed, Figure 3 and 19 show that under early stopping, GD may be favored even if NGD has lower stationary risk.

On the other hand, early stopping may not always improve the well-specified bias. While a complete analysis is difficult partially due to the non-monotonicity of the bias term (see Appendix A.2), we speculate that previous observations on the stationary bias also translate to early stopping. As a concrete example, we consider well-specified settings in which either GD or NGD achieves the optimal stationary bias, and demonstrate that such optimality is also preserved under early stopping:

Proposition 6. Given (A1) and denote the optimal early stopping bias as $B^{\text{opt}}(\theta) = \inf_{t \geq 0} B(\theta(t))$. When $\Sigma_\theta = \Sigma_X^{-1}$, we have $B^{\text{opt}}(\theta_P) \geq B^{\text{opt}}(\theta_{F^{-1}})$ for all P satisfying (A3). Whereas when $\Sigma_\theta = I_d$, we have $B^{\text{opt}}(\theta_{F^{-1}}) \geq B^{\text{opt}}(\theta_I)$.

Figure 19 illustrates that the observed trend in the stationary bias (well-specified) is indeed preserved under optimal early stopping: GD or NGD achieves lower early stopping bias under isotropic or misaligned teacher model, respectively. We leave a more precise characterization as future work.

4.3 FAST DECAY OF POPULATION RISK

Our previous analysis suggests that certain preconditioners can achieve lower population risk, but does not address which method decreases the risk more efficiently. Knowing that preconditioned updates may accelerate optimization, one natural question to ask is, is this speedup also present for generalization under fixed dataset? We answer this question in the affirmative in a slightly different model: we study least squares regression in the RKHS, and show that a preconditioned update that interpolates between GD and NGD achieves the minimax optimal rate in much fewer steps than GD.

We provide a brief outline and defer the detailed setup to Appendix D.8. Let \mathcal{H} be an RKHS included in $L_2(P_X)$ equipped with a bounded kernel function k , and $K_x \in \mathcal{H}$ be the Riesz representation of the kernel. Define S as the canonical operator from \mathcal{H} to $L_2(P_X)$, and write $\Sigma = S^*S$ and $L = SS^*$. We aim to learn the teacher model f^* under the following standard regularity assumptions:

- **(A4) Source Condition:** $\exists r \in (0, \infty)$, $M > 0$ s.t. $f^* = L^r h^*$ for $h^* \in L_2(P_X)$ and $\|f^*\|_\infty \leq M$.
- **(A5) Capacity Condition:** $\exists s > 1$ s.t. $\text{tr}(\Sigma^{1/s}) < \infty$ and $2r + s^{-1} > 1$.
- **(A6) Regularity of RKHS:** $\exists \mu \in [s^{-1}, 1]$, $C_\mu > 0$ s.t. $\sup_{x \in \text{supp}(P_X)} \|\Sigma^{1/2-1/\mu} K_x\|_{\mathcal{H}} \leq C_\mu$.

Note that in the source condition (A4), the coefficient r controls the complexity of the teacher model and relates to the notions of model misalignment in Section 3.2: large r indicates a smoother teacher model which is “easier” to learn, and vice versa (Steinwart et al., 2009). Given training points $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, we consider the following preconditioned update on the student model $f_t \in \mathcal{H}$:

$$f_t = f_{t-1} - \eta(\Sigma + \alpha I)^{-1}(\hat{\Sigma} f_{t-1} - \hat{S}^* Y), \quad f_0 = 0, \quad (4.1)$$

where $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{x}_i} \otimes K_{\mathbf{x}_i}$, $\hat{S}^* Y = \frac{1}{n} \sum_{i=1}^n y_i K_{\mathbf{x}_i}$. In this setting, the population Fisher corresponds to covariance operator Σ , and thus (4.1) can be interpreted as *additive* interpolation between GD and NGD: update with large α behaves like GD, and small α like NGD. Related to our update is the FALKON algorithm (Rudi et al., 2017), which is a preconditioned gradient method for kernel ridge regression (under a different preconditioner). The key distinction is that we consider optimizing the original objective (instead of a regularized version as in FALKON) under early stopping. This

is important since we aim to understand *how preconditioning affects generalization*, and therefore explicit regularization should not be taken into account (for more discussion see Appendix D.8.1).

The following theorem shows that with appropriately chosen α , the preconditioned update (4.1) leads to more efficient decrease in the population risk, due to faster decay of the bias term.

Theorem 7 (Informal). *Under (A4-6), the population risk of f_t can be written as $R(f_t) = \|Sf_t - f^*\|_{L_2(P_X)}^2 \leq B(t) + V(t)$ defined in Appendix D.8. Given $r \geq 1/2$ or $\mu \leq 2r$, preconditioned update (4.1) with $\alpha = n^{-\frac{2s}{2rs+1}}$ achieves minimax optimal convergence rate $R(f_t) = \tilde{O}\left(n^{-\frac{2rs}{2rs+1}}\right)$ in $t = \Theta(\log n)$ steps, whereas ordinary gradient descent requires $t = \Theta\left(n^{\frac{2rs}{2rs+1}}\right)$ steps.*

We comment that the optimal interpolation coefficient α and stopping time t are chosen to balance the bias $B(t)$ and variance $V(t)$. Note that α depends on the teacher model in the following way: for $n > 1$, α decreases as r becomes smaller, which corresponds to non-smooth and “difficult” f^* , and vice versa. This agrees with our previous observation that NGD is advantageous when the teacher model is difficult to learn. We defer empirical verification of this result to Appendix C.

5 NEURAL NETWORK EXPERIMENTS

Protocol. We compare the generalization performance of GD and NGD in neural network settings and illustrate the influence of the following three factors: (i) label noise; (ii) misspecification; (iii) signal misalignment. We also verify the potential advantage of interpolating between GD and NGD.

We consider the MNIST and CIFAR-10 datasets. To create a student-teacher setup, we split the training set into two halves, one of which (*pretrain* split) along with the original labels is used to pretrain the teacher, and the other (*distill* split) along with the teacher’s labels is used to distill (Hinton et al., 2015) the student. We take the teacher to be either a two-layer fully-connected ReLU network or a ResNet (He et al., 2016), and the student is a two-layer ReLU network. We normalize the teacher’s labels logits following Ba & Caruana (2014) before potentially adding label noise, and fit the student by minimizing the L2 loss. Student models are trained on a subset of the distill split with full-batch updates. We implement NGD using Hessian-free optimization (Martens, 2010). We use 100k unlabeled data (possibly applying data augmentation) to estimate the population Fisher. We report the test error when the training error is below 0.2% of its initial value as a proxy for the stationary risk. We defer detailed setup to Appendix E and additional results to Appendix C.

5.1 EMPIRICAL FINDINGS

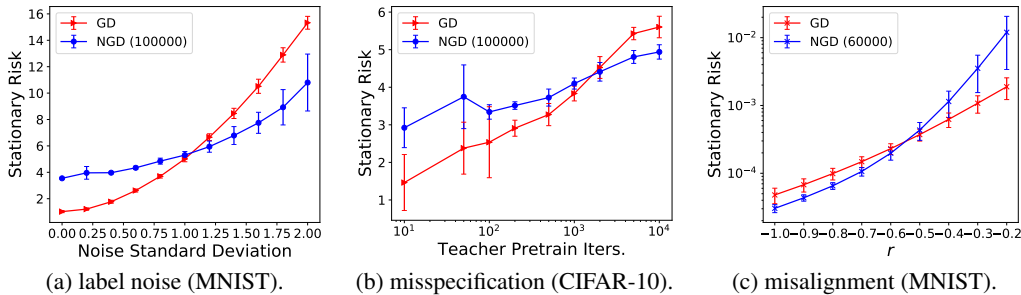


Figure 8: Comparison between NGD and GD. Error bar is one standard deviation away from mean over five independent runs. Numbers in parentheses denote amount of unlabeled examples for estimating the Fisher.

Label Noise. We pretrain the teacher with the full pretrain split and use 1024 examples from the distill split to fit the student. For both the student and teacher, we use a two-layer ReLU net with 80 hidden units. We corrupt the labels with isotropic Gaussian noise of varying magnitude. Figure 8(a) shows that as the noise level increases (the variance term begins to dominate), the stationary risk of both NGD and GD worsen, with GD worsening faster, which aligns with our observation in Figure 4.

Misspecification. We use a ResNet-20 teacher and the same two-layer ReLU student from the label noise experiment. We control the misspecification level by varying amount of pretraining

of the teacher. Intuitively, large teacher models that are trained longer should be more complex and thus likely to be outside of functions that a small two-layer student can represent (i.e. more misspecified). Indeed, Figure 8(b) shows that NGD eventually achieves better generalization as the number of training steps for the teacher increases. In Appendix A.5 we report a heuristic measure of model misspecification that relates to the student’s NTK matrix (Jacot et al., 2018), and confirm that the quantity increases as more label noise is added or as the teacher model is trained longer.

Misalignment. We set the student and teacher models to be the same two-layer ReLU network. We construct the teacher by perturbing the student’s initialization, the direction of which is given by \mathbf{F}^r , where \mathbf{F} is the student’s Fisher (estimated from extra unlabeled data) and $r \in [-1, 0]$. Intuitively, as r approaches -1 , the important parameters of the teacher (i.e. larger update directions) becomes misaligned with the student’s Hessian, and thus learning becomes more “difficult”. While this analogy is rather superficial due to non-convexity, Figure 8(c) shows that as r becomes smaller (setup is more misaligned), NGD begins to generalize better than GD in terms of stationary risk.

Interpolating between Preconditioners. We validate our observations in Section 3 and 4 on the difference between the sample Fisher and population Fisher, and the potential benefit of interpolating between GD and NGD, in neural network experiments. Figure 9(a) shows that as we decrease the number of unlabeled data in estimating the Fisher, which renders the preconditioner closer to the sample Fisher, the stationary risk becomes more akin to that of GD, especially under large noise. This agrees with our remark on sample vs. population Fisher in Section 3 and Appendix A.1.

Figure 9(b)(c) supports the finding in Section 4.1 on the bias-variance tradeoff in neural network settings. In particular, we interpret the left end to correspond to a bias-dominant regime (due to the same architecture of two-layer MLP for the student and teacher), and the right end to correspond to the variance-dominant regime (due to the added label noise). Observe that at a certain SNR, a preconditioner that interpolates between GD and NGD achieves lower stationary risk.

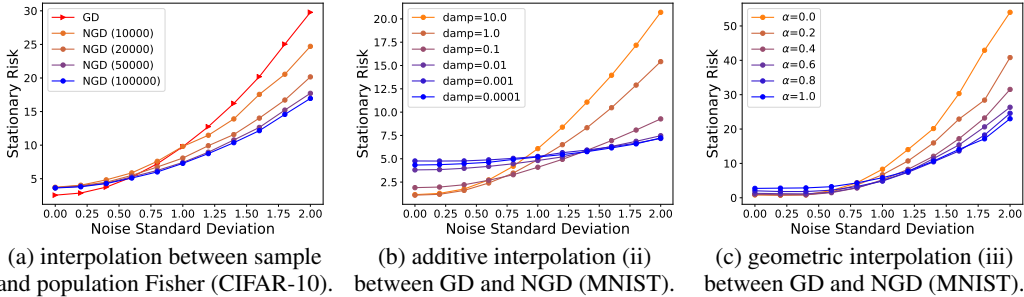


Figure 9: (a) numbers in parentheses indicate the amount of unlabeled data used in estimating the Fisher \mathbf{F} ; we expect the estimated Fisher to be closer to the sample Fisher when the number of unlabeled data is small. (a) additive interpolation $\mathbf{P} = (\mathbf{F} + \alpha \mathbf{I}_d)^{-1}$; larger damping parameter yields update closer to GD. (b) geometric interpolation $\mathbf{P} = \mathbf{F}^{-\alpha}$; larger α parameter yields update closer to that of NGD (blue).

6 DISCUSSION AND CONCLUSION

We analyzed the generalization properties of a general class of preconditioned gradient descent in overparameterized least squares regression, with particular emphasis on natural gradient descent. We identified three factors that affect the comparison of generalization performance of different optimizers, the influence of which we also empirically observed in neural network⁷. We then determined the optimal preconditioner for each factor. While the optimal \mathbf{P} is often not known in practice, we provided justification for common algorithmic choices by discussing the bias-variance tradeoff. Note that our current setup is limited to fixed preconditioner or those that do not alter the span of the gradients, which does not cover many adaptive gradient methods; understanding these optimizers in similar setting would be an interesting future direction. Another important problem is to further characterize the interplay between preconditioning and explicit (e.g. weight decay) or algorithmic regularization (e.g. large step size and gradient noise).

⁷We however note that observations in linear or kernel models do not always translate to neural networks: many works have demonstrated such a gap (e.g. Suzuki (2018); Ghorbani et al. (2019); Allen-Zhu & Li (2019)).

REFERENCES

- Naman Agarwal, Brian Bullins, Xinyi Chen, Elad Hazan, Karan Singh, Cyril Zhang, and Yi Zhang. The case for full-matrix adaptive regularization. *arXiv preprint arXiv:1806.02958*, 2018.
- Alnur Ali, J Zico Kolter, and Ryan J Tibshirani. A continuous-time view of early stopping for least squares. In *International Conference on Artificial Intelligence and Statistics*, volume 22, 2019.
- Zeyuan Allen-Zhu and Yuanzhi Li. What can resnet learn efficiently, going beyond kernels? *arXiv preprint arXiv:1905.10337*, 2019.
- Shun-ichi Amari. Neural learning in structured parameter spaces-natural riemannian gradient. In *Advances in neural information processing systems*, pp. 127–133, 1997.
- Shun-Ichi Amari. Natural gradient works efficiently in learning. *Neural computation*, 10(2):251–276, 1998.
- Shun-Ichi Amari, Hyeyoung Park, and Kenji Fukumizu. Adaptive method of realizing natural gradient learning for multilayer perceptrons. *Neural computation*, 12(6):1399–1409, 2000.
- Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. In *Advances in Neural Information Processing Systems*, pp. 7411–7422, 2019a.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. *arXiv preprint arXiv:1901.08584*, 2019b.
- Navid Azizan and Babak Hassibi. Stochastic gradient/mirror descent: Minimax optimality and implicit regularization. *arXiv preprint arXiv:1806.00952*, 2018.
- Navid Azizan, Sahin Lale, and Babak Hassibi. Stochastic mirror descent on overparameterized nonlinear models: Convergence, implicit regularization, and generalization. *arXiv preprint arXiv:1906.03830*, 2019.
- Jimmy Ba and Rich Caruana. Do deep nets really need to be deep? In *Advances in neural information processing systems*, pp. 2654–2662, 2014.
- Jimmy Ba, Murat A. Erdogdu, Taiji Suzuki, Denny Wu, and Tianzong Zhang. Generalization of two-layer neural networks: An asymptotic viewpoint. *International Conference on Learning Representations*, 2020.
- Zhi-Dong Bai and Yong-Qua Yin. Limit of the smallest eigenvalue of a large dimensional sample covariance matrix. In *Advances In Statistics*, pp. 108–127. World Scientific, 2008.
- Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *arXiv preprint arXiv:1906.11300*, 2019.
- Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine learning and the bias-variance trade-off. *arXiv preprint arXiv:1812.11118*, 2018a.
- Mikhail Belkin, Daniel J Hsu, and Partha Mitra. Overfitting or perfect fitting? risk bounds for classification and regression rules that interpolate. In *Advances in neural information processing systems*, pp. 2300–2311, 2018b.
- Mikhail Belkin, Alexander Rakhlin, and Alexandre B Tsybakov. Does data interpolation contradict statistical optimality? *arXiv preprint arXiv:1806.09471*, 2018c.
- Alberto Bietti and Julien Mairal. On the inductive bias of neural tangent kernels. In *Advances in Neural Information Processing Systems*, pp. 12873–12884, 2019.
- Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- Tianle Cai, Ruiqi Gao, Jikai Hou, Siyu Chen, Dong Wang, Di He, Zhihua Zhang, and Liwei Wang. A gram-gauss-newton method learning overparameterized deep neural networks for regression problems. *arXiv preprint arXiv:1905.11675*, 2019.

- Yuan Cao, Zhiying Fang, Yue Wu, Ding-Xuan Zhou, and Quanquan Gu. Towards understanding the spectral bias of deep learning. *arXiv preprint arXiv:1912.01198*, 2019.
- Andrea Caponnetto and Ernesto De Vito. Optimal rates for the regularized least-squares algorithm. *Foundations of Computational Mathematics*, 7(3):331–368, 2007.
- Lenaïc Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. *arXiv preprint arXiv:2002.04486*, 2020.
- Felipe Cucker and Steve Smale. On the mathematical foundations of learning. *Bulletin of the American mathematical society*, 39(1):1–49, 2002.
- Zeyu Deng, Abba Kammoun, and Christos Thrampoulidis. A model of double descent for high-dimensional binary linear classification. *arXiv preprint arXiv:1911.05822*, 2019.
- Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize for deep nets. *arXiv preprint arXiv:1703.04933*, 2017.
- Edgar Dobriban, Stefan Wager, et al. High-dimensional asymptotics of prediction: Ridge regression and classification. *The Annals of Statistics*, 46(1):247–279, 2018.
- Bin Dong, Jikai Hou, Yiping Lu, and Zhihua Zhang. Distillation \approx early stopping? harvesting dark knowledge utilizing anisotropic information retrieval for overparameterized neural network. *arXiv preprint arXiv:1910.01255*, 2019.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(Jul):2121–2159, 2011.
- Federica Gerace, Bruno Loureiro, Florent Krzakala, Marc Mézard, and Lenka Zdeborová. Generalisation error in learning with random features and the hidden manifold model. *arXiv preprint arXiv:2002.09339*, 2020.
- Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Limitations of lazy training of two-layers neural networks. *arXiv preprint arXiv:1906.08899*, 2019.
- Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient dynamics in deep linear neural networks. *arXiv preprint arXiv:1904.13262*, 2019.
- Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Implicit regularization in matrix factorization. In *Advances in Neural Information Processing Systems*, pp. 6151–6159, 2017.
- Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. *arXiv preprint arXiv:1802.08246*, 2018a.
- Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. In *Advances in Neural Information Processing Systems*, pp. 9461–9471, 2018b.
- Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *arXiv preprint arXiv:1903.08560*, 2019.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Hong Hu and Yue M Lu. Universality laws for high-dimensional learning with random features. *arXiv preprint arXiv:2009.07669*, 2020.
- Takashi Ishida, Ikko Yamane, Tomoya Sakai, Gang Niu, and Masashi Sugiyama. Do we need zero training loss after achieving zero training error? *arXiv preprint arXiv:2002.08709*, 2020.

- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pp. 8571–8580, 2018.
- Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. *arXiv preprint arXiv:1810.02032*, 2018.
- Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In *Conference on Learning Theory*, pp. 1772–1798, 2019.
- Ryo Karakida and Kazuki Osawa. Understanding approximate fisher information for fast convergence of natural gradient descent in wide neural networks. *Advances in Neural Information Processing Systems*, 33, 2020.
- Ryo Karakida, Shotaro Akaho, and Shun-ichi Amari. Universal statistics of fisher information in deep neural networks: Mean field approach. *arXiv preprint arXiv:1806.01316*, 2018.
- Noureddine El Karoui. Asymptotic behavior of unregularized and ridge-regularized high-dimensional robust regression estimators: rigorous results. *arXiv preprint arXiv:1311.2445*, 2013.
- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Frederik Kunstner, Lukas Balles, and Philipp Hennig. Limitations of the empirical fisher approximation. *arXiv preprint arXiv:1905.12558*, 2019.
- Olivier Ledoit and Sandrine Pécché. Eigenvectors of some large sample covariance matrix ensembles. *Probability Theory and Related Fields*, 151(1-2):233–264, 2011.
- Daniel Levy and John C Duchi. Necessary and sufficient geometries for gradient methods. In *Advances in Neural Information Processing Systems*, pp. 11491–11501, 2019.
- Mingchen Li, Mahdi Soltanolkotabi, and Samet Oymak. Gradient descent with early stopping is provably robust to label noise for overparameterized neural networks. *arXiv preprint arXiv:1903.11680*, 2019.
- Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. *arXiv preprint arXiv:1712.09203*, 2017.
- Tengyuan Liang and Alexander Rakhlin. Just interpolate: Kernel” ridgeless” regression can generalize. *arXiv preprint arXiv:1808.00387*, 2018.
- Zhenyu Liao and Romain Couillet. The dynamics of learning: a random matrix approach. *arXiv preprint arXiv:1805.11917*, 2018.
- Junhong Lin and Lorenzo Rosasco. Optimal rates for multi-pass stochastic gradient methods. *The Journal of Machine Learning Research*, 18(1):3375–3421, 2017.
- Cosme Louart, Zhenyu Liao, Romain Couillet, et al. A random matrix approach to neural networks. *The Annals of Applied Probability*, 28(2):1190–1248, 2018.
- Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. *arXiv preprint arXiv:1906.05890*, 2019.
- James Martens. Deep learning via hessian-free optimization. In *ICML*, volume 27, pp. 735–742, 2010.

- James Martens. New insights and perspectives on the natural gradient method. *arXiv preprint arXiv:1412.1193*, 2014.
- James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pp. 2408–2417, 2015.
- James Martens and Ilya Sutskever. Training deep and recurrent networks with hessian-free optimization. In *Neural networks: Tricks of the trade*, pp. 479–535. Springer, 2012.
- Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and double descent curve. *arXiv preprint arXiv:1908.05355*, 2019.
- Stanislav Minsker. On some extensions of Bernstein’s inequality for self-adjoint operators. *Statistics & Probability Letters*, 127:111–119, 2017.
- Andrea Montanari, Feng Ruan, Youngtak Sohn, and Jun Yan. The generalization error of max-margin linear classifiers: High-dimensional asymptotics in the overparametrized regime. *arXiv preprint arXiv:1911.01544*, 2019.
- Tomoya Murata and Taiji Suzuki. Gradient descent in rkhs with importance labeling. *arXiv preprint arXiv:2006.10925*, 2020.
- Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. *arXiv preprint arXiv:1912.02292*, 2019.
- Jorge Nocedal and Stephen Wright. *Numerical optimization*. Springer Science & Business Media, 2006.
- Yann Ollivier. Riemannian metrics for neural networks i: feedforward networks. *Information and Inference: A Journal of the IMA*, 4(2):108–153, 2015.
- Razvan Pascanu and Yoshua Bengio. Revisiting natural gradient for deep networks. *arXiv preprint arXiv:1301.3584*, 2013.
- Loucas Pillaud-Vivien, Alessandro Rudi, and Francis Bach. Statistical optimality of stochastic gradient descent on hard learning problems through multiple passes. In *Advances in Neural Information Processing Systems*, pp. 8114–8124, 2018.
- Qian Qian and Xiaoyuan Qian. The implicit bias of adagrad on separable data. In *Advances in Neural Information Processing Systems*, pp. 7759–7767, 2019.
- Noam Razin and Nadav Cohen. Implicit regularization in deep learning may not be explainable by norms. *arXiv preprint arXiv:2005.06398*, 2020.
- Francisco Rubio and Xavier Mestre. Spectral convergence for a general class of random matrices. *Statistics & probability letters*, 81(5):592–602, 2011.
- Alessandro Rudi and Lorenzo Rosasco. Generalization properties of learning with random features. In *Advances in Neural Information Processing Systems*, pp. 3215–3225, 2017.
- Alessandro Rudi, Luigi Carratino, and Lorenzo Rosasco. Falkon: An optimal large scale kernel method. In *Advances in neural information processing systems*, pp. 3888–3898, 2017.
- Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.
- Ingo Steinwart, Don R Hush, and Clint Scovel. Optimal rates for regularized least squares regression. In *COLT*, pp. 79–93, 2009.

- Lili Su and Pengkun Yang. On learning over-parameterized neural networks: A functional approximation perspective. In *Advances in Neural Information Processing Systems*, pp. 2637–2646, 2019.
- Arun Suggala, Adarsh Prasad, and Pradeep K Ravikumar. Connecting optimization and regularization paths. In *Advances in Neural Information Processing Systems*, pp. 10608–10619, 2018.
- Taiji Suzuki. Adaptivity of deep relu network for learning in besov and mixed smooth besov spaces: optimal rate and curse of dimensionality. *arXiv preprint arXiv:1810.08033*, 2018.
- Valentin Thomas, Fabian Pedregosa, Bart Merriënboer, Pierre-Antoine Manzagol, Yoshua Bengio, and Nicolas Le Roux. On the interplay between noise and curvature and its effect on optimization and generalization. In *International Conference on Artificial Intelligence and Statistics*, pp. 3503–3513, 2020.
- Sharan Vaswani, Reza Babanezhad, Jose Gallego, Aaron Mishkin, Simon Lacoste-Julien, and Nicolas Le Roux. To each optimizer a norm, to each norm its generalization. *arXiv preprint arXiv:2006.06821*, 2020.
- Neha S Wadia, Daniel Duckworth, Samuel S Schoenholz, Ethan Dyer, and Jascha Sohl-Dickstein. Whitening and second order optimization both destroy information about the dataset, and can make generalization impossible. *arXiv preprint arXiv:2008.07545*, 2020.
- Per-Åke Wedin. Perturbation theory for pseudo-inverses. *BIT Numerical Mathematics*, 13(2):217–232, 1973.
- Francis Williams, Matthew Trager, Claudio Silva, Daniele Panozzo, Denis Zorin, and Joan Bruna. Gradient dynamics of shallow univariate relu networks. *arXiv preprint arXiv:1906.07842*, 2019.
- Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. The marginal value of adaptive gradient methods in machine learning. In *Advances in Neural Information Processing Systems*, pp. 4148–4158, 2017.
- Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. *arXiv preprint arXiv:2002.09277*, 2020.
- Lei Wu, Chao Ma, and E Weinan. How sgd selects the global minima in over-parameterized learning: A dynamical stability perspective. In *Advances in Neural Information Processing Systems*, pp. 8279–8288, 2018.
- Bo Xie, Yingyu Liang, and Le Song. Diverse neural network learns true target functions. *arXiv preprint arXiv:1611.03131*, 2016.
- Ji Xu and Daniel Hsu. How many variables should be entered in a principal component regression equation? *arXiv preprint arXiv:1906.01139*, 2019.
- Peng Xu, Fred Roosta, and Michael W Mahoney. Second-order optimization for non-convex machine learning: An empirical study. In *Proceedings of the 2020 SIAM International Conference on Data Mining*, pp. 199–207. SIAM, 2020.
- Guodong Zhang, Lala Li, Zachary Nado, James Martens, Sushant Sachdeva, George Dahl, Chris Shallue, and Roger B Grosse. Which algorithmic choices matter at which batch sizes? insights from a noisy quadratic model. In *Advances in Neural Information Processing Systems*, pp. 8194–8205, 2019a.
- Guodong Zhang, James Martens, and Roger B Grosse. Fast convergence of natural gradient descent for over-parameterized neural networks. In *Advances in Neural Information Processing Systems*, pp. 8080–8091, 2019b.

TABLE OF CONTENTS

A	Discussion on Additional Results	16
A.1	Implicit Bias of GD vs. NGD	16
A.2	Non-monotonicity of the Bias Term w.r.t. Time	17
A.3	Estimating the Population Fisher	17
A.4	Bias Term Under Specific Source Condition	18
A.5	Interpretation of $\sqrt{\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}/n}$	18
B	Additional Related Works	19
C	Additional Figures	20
C.1	Overparameterized Linear Regression	20
C.2	RKHS Regression	21
C.3	Neural Networks	22
D	Proofs and Derivations	23
D.1	Missing Derivations in Section 3	23
D.2	Proof of Theorem 1	23
D.3	Proof of Theorem 2	24
D.4	Proof of Proposition 3	26
D.5	Proof of Proposition 4	26
D.6	Proof of Proposition 5	29
D.7	Proof of Proposition 6	29
D.8	Proof of Theorem 7	30
D.9	Proof of Proposition 8	37
D.10	Proof of Proposition 9	37
D.11	Proof of Corollary 10	38
E	Experiment Setup	39
E.1	Processing the Datasets	39
E.2	Setup and Implementation for Optimizers	39
E.3	Other Details	39

A DISCUSSION ON ADDITIONAL RESULTS

A.1 IMPLICIT BIAS OF GD VS. NGD

It is known that gradient descent is the steepest descent with respect to the ℓ_2 norm, i.e. the update direction is constructed to decrease the loss under small changes in the parameters measured by the ℓ_2 norm (Gunasekar et al., 2018a). Following this analogy, NGD is the steepest descent with respect to the KL divergence on the predictive distributions (Martens, 2014); this can be interpreted as a proximal update which penalizes how much the predictions change on the data distribution.

Intuitively, the above discussion suggests GD tend to find solution that is close to the initialization in the Euclidean distance between parameters, whereas NGD prefers solution close to the initialization in terms of the function outputs on P_X . This observation turns out to be exact in the case of ridgeless interpolant under the squared loss, as remarked in Section 3. Moreover, Figure 1 and 10 confirms the same trend in neural network optimization. In particular, we observe that

- GD results in small changes in the parameters (left), and NGD results in small changes in the learned function (right).
- preconditioning with the pseudo-inverse of the sample Fisher, i.e., $\mathbf{P} = (\mathbf{J}^\top \mathbf{J})^\dagger$, leads to implicit bias similar to that of GD, but not NGD with the population Fisher.
- interpolating between GD and NGD ($\mathbf{P} = \mathbf{F}^{-1/2}$) results in properties in between GD and NGD.

Remark. *Qualitatively speaking, the small change in the function output is the essential reason that NGD performs well under noisy labels in the interpolation setting: NGD seeks to interpolate the training data by changing the function only “locally”, so that memorizing the noisy labels has small impact on the “global” shape of the learned function.*

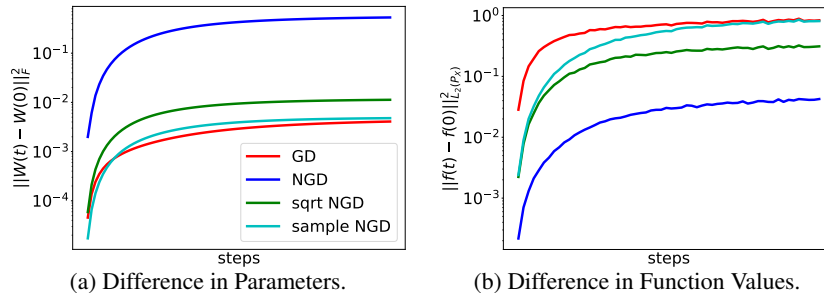


Figure 10: Illustration of implicit bias of GD and NGD. We set $n = 100$, $d = 50$, and regress a two-layer ReLU network with 50 hidden units towards a teacher model of the same architecture on Gaussian input. The x-axis is rescaled for each optimizer such that the final training error is below 10^{-3} . GD finds solution with small changes in the parameters, whereas NGD finds solution with small changes in the function. Note that the sample Fisher (cyan) has implicit bias similar to GD and does not resemble NGD (population Fisher).

We note that the above observation also implies that wide neural networks trained with NGD (population Fisher) is less likely to stay in the kernel regime: the distance traveled from initialization can be large (see Figure 10(a)) and thus the Taylor expansion around the initialization is no longer accurate. In other words, the analogy between wide neural net and its linearized kernel model (which we partially employed in Section 5) may not be valid in models trained with NGD⁸.

Implicit Bias of Interpolating Preconditioners. We also expect that as we interpolate from GD to NGD, the distance traveled by the parameter space would gradually increase, and distance traveled in the function space would decrease. Figure 11 demonstrate that this is indeed the case for linear model as well as neural network: we use the same two-layer MLP setup on MNIST as in Section 5. Note that updates that are closer to GD result in smaller change in the parameters, whereas ones close to NGD lead to smaller change in the function outputs.

⁸Note that this gap is only present when the population Fisher is used; previous works have shown NTK-type global convergence for sample Fisher-related update (Zhang et al., 2019b; Cai et al., 2019).

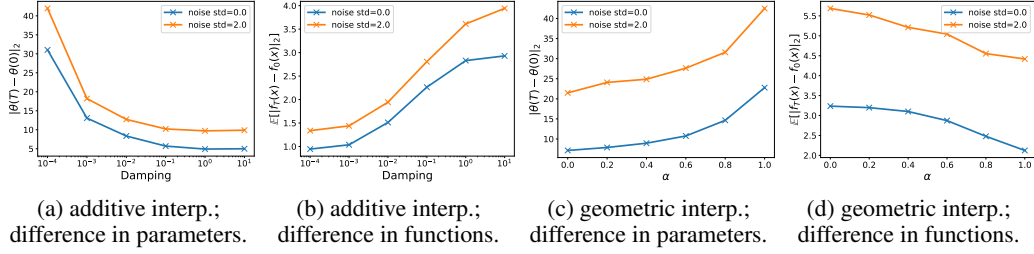


Figure 11: Illustration of the implicit bias of preconditioned gradient descent that interpolates between GD and NGD on MNIST. Note that as the update becomes more similar to NGD (smaller damping or larger α), the distance traveled in the parameter space increases, where as the distance traveled on the output space decreases.

A.2 NON-MONOTONICITY OF THE BIAS TERM W.R.T. TIME

Many previous works on the high-dimensional characterization of linear regression assumed a random effects model with an isotropic prior on the true parameters (Dobriban et al., 2018; Hastie et al., 2019; Xu & Hsu, 2019), which may not be realistic. As an example of the limitation of this assumption, we note that when $\Sigma_\theta = \mathbf{I}_d$, it can be shown that the expected bias $B(\hat{\theta}(t))$ monotonically decreases through time (see proof of Proposition 6 for details). In contrast, when the target parameters do not follow an isotropic prior, the bias of GD can exhibit non-monotonicity, which gives rise to the surprising “epoch-wise double descent” phenomenon also observed in deep learning (Nakki-ran et al., 2019; Ishida et al., 2020).

We empirically demonstrate this non-monotonicity when the model is close to the interpolation threshold in Figure 12. We set eigenvalues of Σ_X to be two equally-weighted point masses with $\kappa_X = 32$, $\Sigma_\theta = \Sigma_X^{-1}$ and $\gamma = 16/15$. Note that the GD trajectory (red) exhibits non-monotonicity in the bias term, whereas for NGD the bias is monotonically decreasing through time (which we confirm in the proof of Proposition 6). We remark that this mechanism of epoch-wise double descent may not be related to the empirical findings in deep neural networks (the robustness of which is also unknown), in which it is typically speculated that the variance term exhibits non-monotonicity.

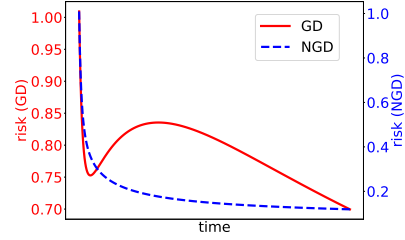


Figure 12: Epoch-wise double descent. Note that non-monotonicity of the bias term is present in GD but not NGD.

A.3 ESTIMATING THE POPULATION FISHER

Our analysis on linear model considers the idealized setup with access to the exact population Fisher, which can be estimated using additional unlabeled data. In this section we discuss how our result in Section 3 and Section 4 are affected when the population covariance is approximated from N i.i.d. (unlabeled) samples $\mathbf{X}_u \in \mathbb{R}^{N \times d}$. For the ridgeless interpolant we have the following result on the substitution error in replacing the true population covariance with the sample estimate.

Proposition 8. *Given (A1)(A3) and $N/d \rightarrow \psi > 1$ as $d \rightarrow \infty$, let $\hat{\Sigma}_X = \mathbf{X}_u^\top \mathbf{X}_u / N$, we have*

- (a) $\|\Sigma_X - \hat{\Sigma}_X\|_2 = O(\psi^{-1/2})$ almost surely.
- (b) Denote the stationary bias and variance of NGD (with the exact population Fisher) as B^* and V^* , respectively, and the bias and variance of the preconditioned update using the approximate Fisher $\hat{\mathbf{F}} = \hat{\Sigma}_X$ as \hat{B} and \hat{V} , respectively. Let $0 < \epsilon < 1$ be the desired accuracy. Then $\psi = \Theta(\epsilon^{-2})$ suffices to achieve $|B^* - \hat{B}| < \epsilon$ and $|V^* - \hat{V}| < \epsilon$.

Proposition 8 entails that when the preconditioner is a sample estimate $\hat{\mathbf{F}}$, we can approach the stationary bias and variance of the population Fisher at a rate of $\psi^{-1/2}$ as we increase the number of unlabeled data N linearly in the dimensionality d . In other words, any accuracy ϵ such that $1/\epsilon$ is bounded can be achieved with finite ψ (to push $\epsilon \rightarrow 0$, additional logarithmic factor is required, e.g. $N = O(d \log d)$, which is beyond the proportional limit).

On the other hand, for our result in Section 4.3, Murata & Suzuki (2020, Lemma A.5) directly implies that setting $N = \Theta(n^s \log n)$ is sufficient to approximate the population covariance operator (i.e., so that $\|\Sigma^{1/2}\Sigma_{N,\lambda}^{-1/2}\| = O(1)$). Finally, note that our analysis above does not impose any structural assumptions on the estimated matrix. When the Fisher exhibits certain structures (e.g. Kronecker factorization (Martens & Grosse, 2015)), then estimation can be more sample-efficient. For more discussion on such approximation of the Fisher matrix see Karakida & Osawa (2020).

A.4 BIAS TERM UNDER SPECIFIC SOURCE CONDITION

Motivated by the connection between the notion of “alignment” and the source condition (Cucker & Smale, 2002) in Section 3.2, we now consider a specific case of θ^* : $\Sigma_\theta = \Sigma_X^{-r}$, where r controls the extent of misalignment, and Theorem 2 implies that the optimal preconditioner for the bias term (well-specified) is $P = \Sigma_X^r$. Note that larger r corresponds to more misaligned and thus “difficult” problem, and vice versa. In this setup we have the following comparison between GD and NGD.

Proposition 9. *Under the setting of Theorem 2 and $\Sigma_\theta = \Sigma_X^{-r}$, for all $r \geq 1$ we have $B(\hat{\theta}_{F^{-1}}) \leq B(\hat{\theta}_I)$, whereas for $r \leq 0$, we have $B(\hat{\theta}_{F^{-1}}) \geq B(\hat{\theta}_I)$; the equality is achieved when the features are isotropic ($\Sigma_X = cI_d$).*

The proposition confirms the intuition that when parameters of the teacher model θ^* are more “aligned” with x than the isotropic case ($r \leq 0$), then GD achieves lower bias than NGD; on the other hand, when Σ_θ is more “misaligned” than Σ_X^{-1} , then NGD is advantageous for the bias term. We therefore expect a transition from the GD-dominated to NGD-dominated regime for some $r \in (0, 1)$. The exact value of r for such transition depends on the spectral distribution of Σ_X (and one would need to specifically evaluate (D.26)). To give a concrete example, in the case where Σ_X has a simple block structure, we can explicitly determine the transition point $r^* \in (0, 1)$, as shown by the following corollary.

Corollary 10. *Assume $\Sigma_\theta = \Sigma_X^{-r}$, and eigenvalues of Σ_X come from two equally-weighted point masses with $\kappa \triangleq \lambda_{\max}(\Sigma_X)/\lambda_{\min}(\Sigma_X)$. WLOG we take $\text{tr}(\Sigma_X)/d = 1$. Then given $r^* = \ln c_{\kappa,\gamma}/\ln \kappa$ (see Appendix D.11 for definition), we have $B(\hat{\theta}_I) \geq B(\hat{\theta}_{F^{-1}})$ if and only if $r \geq r^*$.*

Remark. When $\gamma = 2$, the transition happens at $r^* = 1/2$ which is independent of the condition number κ , as indicated by the dashed line in Figure 13. However for other $\gamma > 1$, r^* generally relates to both γ and κ .

Our characterization above is supported by Figure 13, in which we plot the bias term under varying extent of misalignment (controlled by r) in the setting of Corollary 10. Observe that as we construct the teacher model to be more “misaligned” (and thus difficult to learn) by increasing r , NGD (blue) achieves lower bias compared to GD (red), and vice versa.

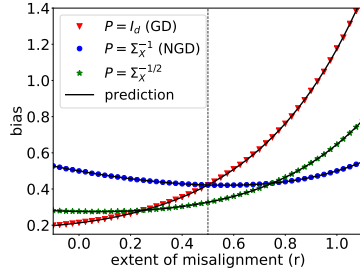


Figure 13: We set $\Sigma_\theta = \Sigma_X^{-r}$, $\gamma = 2$, $\kappa = 20$ and plot the bias (well-specified) under varying r .

A.5 INTERPRETATION OF $\sqrt{y^\top K^{-1}y}/n$

As a heuristic measure of model misspecification, in Figure 14 we report $\sqrt{y^\top K^{-1}y}/n$ studied in Arora et al. (2019b), where y is the label vector and K is the student’s NTK matrix (Jacot et al., 2018). This quantity relates to generalization of wide neural networks in the kernel regime, and can be interpreted as a proxy for measuring how much signal and noise are distributed along the eigendirections of the NTK (see Li et al. (2019); Dong et al. (2019); Su & Yang (2019); Cao et al. (2019) for detailed discussion). Roughly speaking, large $\sqrt{y^\top K^{-1}y}/n$ implies that the problem is difficult to learn by GD, vice versa.

Here we give a heuristic argument on how this quantity relates to label noise and misspecification. For the ridgeless regression

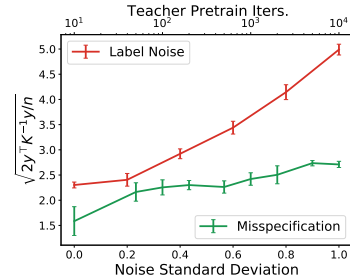


Figure 14: A heuristic measure of misspecification and label noise: $\sqrt{y^\top K^{-1}y}/n$ on CIFAR-10.

C ADDITIONAL FIGURES

C.1 OVERPARAMETERIZED LINEAR REGRESSION

Non-monotonicity of the Risk. Under our generalized (anisotropic) assumption on the covariance of the features and the target, both the bias and the variance term can exhibit non-monotonicity w.r.t. the overparameterization level $\gamma > 1$: in Figure 15 we observe two peaks in the bias term and three peaks in the variance term. In contrast, it is known that when $\Sigma_X = I_d$, both the bias and variance are *monotonic* for when $\gamma > 1$.

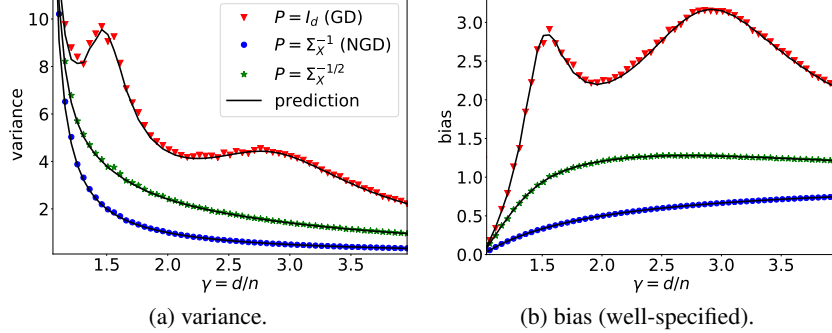


Figure 15: Illustration of the “multiple-descent” curve of the risk for $\gamma > 1$. We take $n = 300$, eigenvalues of Σ_X as three equally-spaced point masses with $\kappa_X = 5000$ and $\|\Sigma_X\|_F^2 = d$, and $\Sigma_\theta = \Sigma_X^{-1}$ (misaligned). Note that for GD, both the bias and the variance are highly non-monotonic for $\gamma > 1$.

Additional Figures for Section 3 and 4. We include additional figures on (a) well-specified bias when $\Sigma_\theta = I_d$ (GD is optimal); (b) misspecified bias under unobserved features (predicted by Proposition 3); (c) bias-variance tradeoff by interpolating between preconditioners (SNR=5). As shown in Figure 16 and 17, in all cases the experimental values match the theoretical predictions.

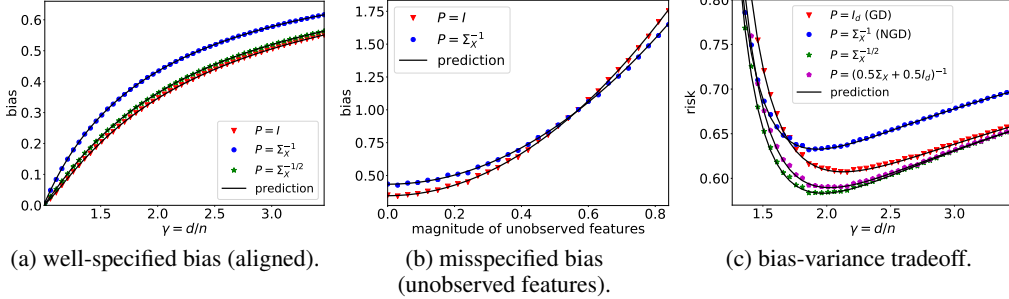


Figure 16: We set eigenvalues of Σ_X as a uniform distribution with $\kappa_X = 20$ and $\|\Sigma_X\|_F^2 = d$.

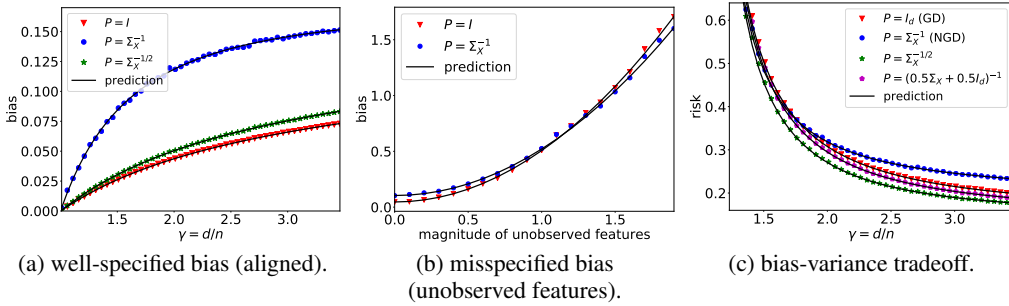


Figure 17: We construct eigenvalues of Σ_X with a polynomial decay: $\lambda_i(\Sigma_X) = i^{-1}$ and then rescale the eigenvalues such that $\kappa_X = 500$ and $\|\Sigma_X\|_F^2 = d$.

Early Stopping Risk. Figure 18 compares the stationary risk with the optimal early stopping risk under varying misalignment level. To increase the extent of misalignment, we set $\Sigma_\theta = \Sigma_X^{-\alpha}$ and vary α from 0 to 1: larger α entails more “misaligned” teacher, and vice versa. Note that as the problem becomes more misaligned, NGD achieves lower stationary and early stopping risk.

Figure 19 reports the optimal early stopping risk under misspecification (same trend can be obtained when the x-axis is label noise). In contrast to the stationary risk (Figure 6), GD can be advantageous under early stopping even with large extent of misspecification (for isotropic teacher). This aligns with our finding in Section 4.2 that early stopping reduces the variance and the misspecified bias.

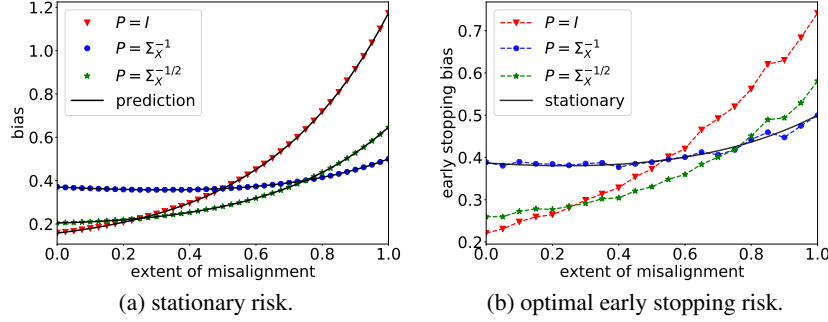


Figure 18: Well-specified bias against different extent of “alignment”. We set $n = 300$, eigenvalues of Σ_X as two point masses with $\kappa_X = 20$, and take $\Sigma_\theta = \Sigma_X^{-\alpha}$ and vary α from 0 to 1. (a) GD achieves lower bias when Σ_θ is isotropic, whereas NGD dominates when $\Sigma_X = \Sigma_\theta^{-1}$; $P = \Sigma_X^{-1/2}$ (interpolates between GD and NGD) is advantageous in between. (b) optimal early stopping bias follows similar trend as stationary bias.

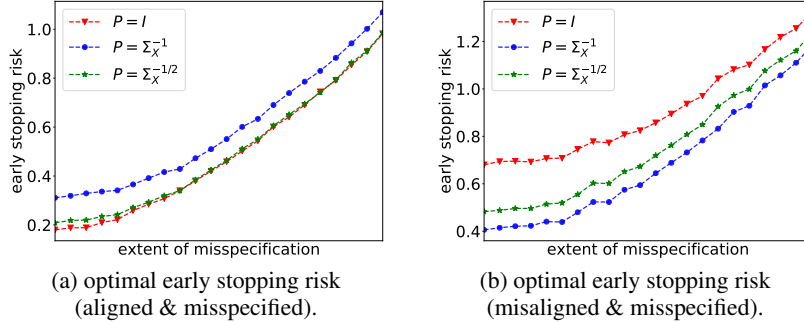


Figure 19: Optimal early stopping risk vs. increasing model misspecification. We follow the same setup as Figure 4(c). (a) $\Sigma_\theta = I_d$ (favors GD); unlike Figure 4(c), GD has lower early stopping risk even under large extent of misspecification. (b) $\Sigma_\theta = \Sigma_X^{-1}$ (favors NGD); NGD is also advantageous under early stopping.

C.2 RKHS REGRESSION

We simulate the optimization in the coordinates of RKHS via a finite-dimensional approximation (using extra unlabeled data). In particular, we consider the teacher model in the form of $f^*(\mathbf{x}) = \sum_{i=1}^N h_i \mu_i^r \phi_i(\mathbf{x})$ for square summable $\{h_i\}_{i=1}^N$, in which r controls the “difficulty” of the learning problem. We find $\{\mu_i\}_{i=1}^N$ and $\{\phi_i\}_{i=1}^N$ by solving the eigenfunction problem for some kernel k . The student model takes the form of $f(\mathbf{x}) = \sum_{i=1}^N \frac{a_i}{\sqrt{\mu_i}} \phi_i(\mathbf{x})$ and we optimize the coefficients $\{a_i\}_{i=1}^N$ via the preconditioned update (4.1). We set $n = 1000$, $d = 5$, $N = 2500$ and consider the inverse multiquadratic (IMQ) kernel: $k(\mathbf{x}, \mathbf{y}) = \frac{1}{\sqrt{1 + \|\mathbf{x} - \mathbf{y}\|_2^2}}$.

Recall that Theorem 7 suggests that for small r , i.e. “difficult” problem, the damping coefficient λ would need to be small (which makes the update NGD-like), and vice versa. This result is (qualitatively) supported by Figure 20, from which we can see that small λ is beneficial when r is small, and vice versa. We remark that this observed trend is rather fragile and sensitive to various hyperparameters, and we leave a comprehensive characterization of this observation as future work.

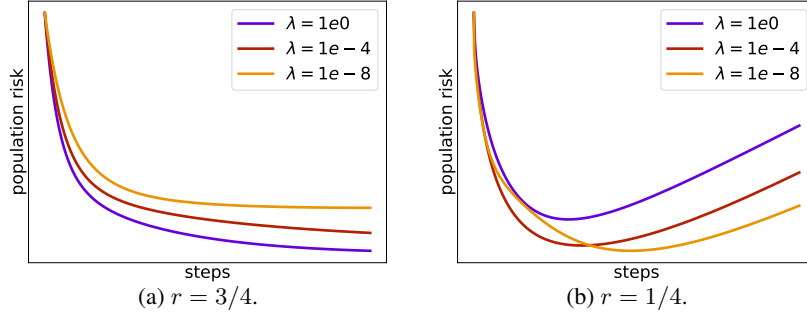


Figure 20: Population risk of the preconditioned update in RKHS that interpolates between GD and NGD. We use the IMQ kernel and set $n = 1000$, $d = 5$, $N = 2500$, $\sigma^2 = 5 \times 10^{-4}$. The x-axis has been rescaled for each curve and thus convergence speed is not directly comparable. Note that (a) large λ (i.e., GD-like update) is beneficial when r is large, and (b) small λ (i.e., NGD-like update) is beneficial when r is small.

C.3 NEURAL NETWORKS

Label Noise. In Figure 21, (a) we observe the same phenomenon on CIFAR-10 that NGD generalizes better as more label noise is added to the training data, and vice versa. Figure 21 (b) shows that in all cases with varying amounts of label noise, the early stopping risk is however worse than that of GD. This agrees with the observation in Section 4 and Figure 19(a) that early stopping can potentially favor GD due to the reduced variance.

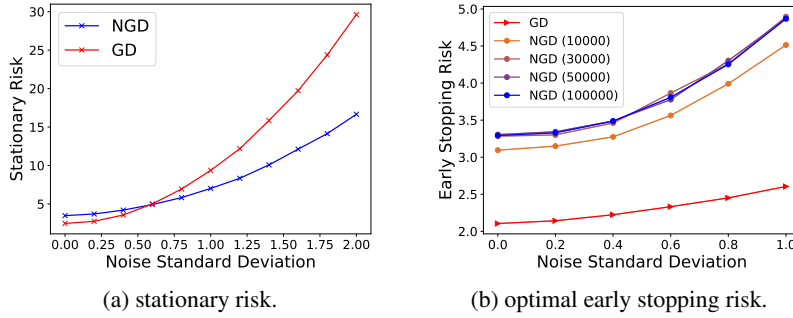


Figure 21: Additional label noise experiment on CIFAR-10.

Misalignment. We illustrate the finding in Proposition 6 and Figure 18(b) in neural networks under synthetic data: we consider 50-dimensional Gaussian input, and both the teacher and the student model are two-layer ReLU networks with 50 hidden units. We construct the teacher by perturbing the initialization of the student as described in Section 5. Figure 22 shows that as r approaches -1 (more “misaligned”), NGD achieves lower early stopping risk (b), whereas GD dominates the early stopping risk in less misaligned setting (a). We note that this phenomenon is difficult to observe in practical neural network training on real-world data, which may be partially due to the fragility of the analogy between neural nets and linear models, especially under NGD (discussed in Appendix A).

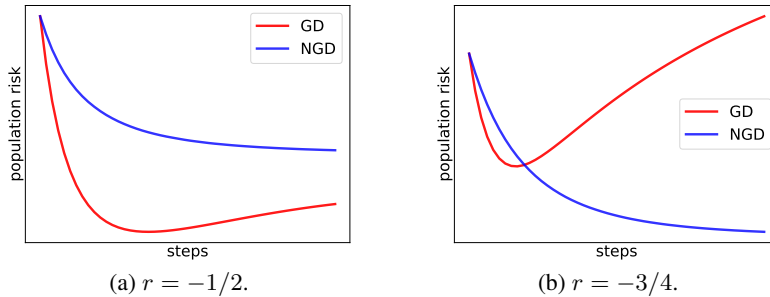


Figure 22: Population risk of two-layer neural networks in the misalignment setup (noiseless) with synthetic Gaussian data. We set $n = 200$, $d = 50$, the damping coefficient $\lambda = 10^{-6}$, and both the student and the teacher are two-layer ReLU networks with 50 hidden units. The x-axis and the learning rate have been rescaled for each curve. When r is sufficiently small, NGD achieves lower early stopping risk, and vice versa.

D PROOFS AND DERIVATIONS

D.1 MISSING DERIVATIONS IN SECTION 3

Gradient Flow of Preconditioned Updates. Given positive definite \mathbf{P} and $\gamma > 1$, it is clear that the gradient flow solution at time t can be written as

$$\boldsymbol{\theta}_P(t) = \mathbf{P}\mathbf{X}^\top \left[\mathbf{I}_n - \exp\left(-\frac{t}{n}\mathbf{X}\mathbf{P}\mathbf{X}^\top\right) \right] (\mathbf{X}\mathbf{P}\mathbf{X}^\top)^{-1}\mathbf{y}.$$

Taking $t \rightarrow \infty$ yields the stationary solution $\hat{\boldsymbol{\theta}}_P = \mathbf{P}\mathbf{X}^\top (\mathbf{X}\mathbf{P}\mathbf{X}^\top)^{-1}\mathbf{y}$. We remark that the damped inverse of the sample Fisher $\mathbf{P} = (\mathbf{X}\mathbf{X}^\top + \lambda\mathbf{I}_d)^{-1}$ leads to the same minimum-norm solution as GD $\hat{\boldsymbol{\theta}}_I = \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{y}$ since $\mathbf{P}\mathbf{X}^\top$ and \mathbf{X} share the same eigenvectors. On the other hand, when \mathbf{P} is the pseudo-inverse of the sample Fisher $(\mathbf{X}\mathbf{X}^\top)^\dagger$ which is not full-rank, the trajectory can be obtained via the variation of constants formula:

$$\boldsymbol{\theta}(t) = \left[\frac{t}{n} \sum_{k=0}^{\infty} \frac{1}{(k+1)!} \left(-\frac{t}{n} \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{X} \right)^k \right] \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{y},$$

for which taking the large t limit also yields the minimum-norm solution $\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{y}$.

Minimum $\|\boldsymbol{\theta}\|_{\mathbf{P}^{-1}}$ Norm Interpolant. For positive definite \mathbf{P} and the corresponding stationary solution $\hat{\boldsymbol{\theta}}_P = \mathbf{P}\mathbf{X}^\top (\mathbf{X}\mathbf{P}\mathbf{X}^\top)^{-1}\mathbf{y}$, note that given any other interpolant $\hat{\boldsymbol{\theta}}'$, we have $(\hat{\boldsymbol{\theta}}_P - \hat{\boldsymbol{\theta}}')\mathbf{P}^{-1}\hat{\boldsymbol{\theta}}_P = 0$ because both $\hat{\boldsymbol{\theta}}_P$ and $\hat{\boldsymbol{\theta}}'$ achieves zero empirical risk. Therefore, $\|\hat{\boldsymbol{\theta}}'\|_{\mathbf{P}^{-1}}^2 - \|\hat{\boldsymbol{\theta}}_P\|_{\mathbf{P}^{-1}}^2 = \|\hat{\boldsymbol{\theta}}' - \hat{\boldsymbol{\theta}}_P\|_{\mathbf{P}^{-1}}^2 \geq 0$. This confirms that $\hat{\boldsymbol{\theta}}_P$ is the unique minimum $\|\boldsymbol{\theta}\|_{\mathbf{P}^{-1}}$ norm solution.

D.2 PROOF OF THEOREM 1

Proof. By the definition of the variance term and the stationary $\hat{\boldsymbol{\theta}}$,

$$V(\hat{\boldsymbol{\theta}}) = \text{tr}\left(\text{Cov}(\hat{\boldsymbol{\theta}})\boldsymbol{\Sigma}_X\right) = \sigma^2 \text{tr}\left(\mathbf{P}\mathbf{X}^\top (\mathbf{X}\mathbf{P}\mathbf{X}^\top)^{-2} \mathbf{X}\mathbf{P}\boldsymbol{\Sigma}_X\right).$$

Write $\bar{\mathbf{X}} = \mathbf{X}\mathbf{P}^{1/2}$. Similarly, we define $\boldsymbol{\Sigma}_{\mathbf{X}\mathbf{P}} = \mathbf{P}^{1/2}\boldsymbol{\Sigma}_X\mathbf{P}^{1/2}$. The equation above thus simplifies to

$$V(\hat{\boldsymbol{\theta}}_P) = \sigma^2 \text{tr}\left(\bar{\mathbf{X}}^\top (\bar{\mathbf{X}}\bar{\mathbf{X}}^\top)^{-2} \bar{\mathbf{X}}\boldsymbol{\Sigma}_{\mathbf{X}\mathbf{P}}\right).$$

The analytic expression of the variance term follows from a direct application of [Hastie et al. \(2019, Theorem 4\)](#), in which conditions on the population covariance are satisfied by (A2).

Taking the derivative of $m(-\lambda)$ yields

$$m'(-\lambda) = \left(\frac{1}{m^2(-\lambda)} - \gamma \int \frac{\tau^2}{(1 + \tau m(-\lambda))^2} d\mathbf{H}_{\mathbf{X}\mathbf{P}}(\tau) \right)^{-1}.$$

Plugging the quantity into the expression of the variance (omitting the scaling σ^2 and constant shift),

$$\frac{m'(-\lambda)}{m^2(-\lambda)} = \left(1 - \gamma m^2(-\lambda) \int \frac{\tau^2}{(1 + \tau m(-\lambda))^2} d\mathbf{H}_{\mathbf{X}\mathbf{P}}(\tau) \right)^{-1}.$$

From the monotonicity of $\frac{x}{1+x}$ on $x > 0$ or the Jensen's inequality we know that

$$1 - \gamma \int \left(\frac{\tau m(-\lambda)}{1 + \tau m(-\lambda)} \right)^2 d\mathbf{H}_{\mathbf{X}\mathbf{P}}(\tau) \leq 1 - \gamma \left(\int \frac{\tau m(-\lambda)}{1 + \tau m(-\lambda)} d\mathbf{H}_{\mathbf{X}\mathbf{P}}(\tau) \right)^2.$$

Taking $\lambda \rightarrow 0$ and omitting the scalar σ^2 , the RHS evaluates to $1 - 1/\gamma$. We thus arrive at the lower bound $V \geq (\gamma - 1)^{-1}$. Note that the equality is only achieved when $\mathbf{H}_{\mathbf{X}\mathbf{P}}$ is a point mass, i.e. $\mathbf{P} = \boldsymbol{\Sigma}_X^{-1}$. In other words, the minimum variance is achieved by NGD. As a verification, the variance of the NGD solution $\hat{\boldsymbol{\theta}}_{\mathbf{P}^{-1}}$ agrees with the calculation in [Hastie et al. \(2019, A.3\)](#). \square

D.3 PROOF OF THEOREM 2

Proof. By the definition of the bias term (note that $\Sigma_X, \Sigma_\theta, P$ are all positive semi-definite),

$$\begin{aligned}
B(\hat{\theta}_P) &= \mathbb{E}_{\theta^*} \left[\left\| P X^\top (X P X^\top)^{-1} X \theta_* - \theta^* \right\|_{\Sigma_X}^2 \right] \\
&= \frac{1}{d} \text{tr} \left(\Sigma_\theta \left(I_d - P X^\top (X P X^\top)^{-1} X \right)^\top \Sigma_X \left(I_d - P X^\top (X P X^\top)^{-1} X \right) \right) \\
&\stackrel{(i)}{=} \frac{1}{d} \text{tr} \left(\Sigma_{\theta/P} \left(I_d - \bar{X}^\top (\bar{X} \bar{X}^\top)^{-1} \bar{X} \right)^\top \Sigma_{XP} \left(I_d - \bar{X}^\top (\bar{X} \bar{X}^\top)^{-1} \bar{X} \right) \right) \\
&\stackrel{(ii)}{=} \lim_{\lambda \rightarrow 0^+} \frac{\lambda^2}{d} \text{tr} \left(\Sigma_{\theta/P} \left(\frac{1}{n} \bar{X}^\top \bar{X} + \lambda I_d \right)^{-1} \Sigma_{XP} \left(\frac{1}{n} \bar{X}^\top \bar{X} + \lambda I_d \right)^{-1} \right) \\
&\stackrel{(iii)}{=} \lim_{\lambda \rightarrow 0^+} \frac{\lambda^2}{d} \text{tr} \left(\left(\frac{1}{n} \hat{X}^\top \hat{X} + \lambda \Sigma_{\theta/P}^{-1} \right)^{-2} \Sigma_{\theta/P}^{-1/2} \Sigma_{XP} \Sigma_{\theta/P}^{-1/2} \right),
\end{aligned}$$

where we utilized (A3) and defined $\bar{X} = X P^{1/2}$, $\Sigma_{XP} = P^{1/2} \Sigma_X P^{1/2}$, $\Sigma_{\theta/P} = P^{-1/2} \Sigma_\theta P^{-1/2}$ in (i), applied the equality $(A A^\top)^\dagger A = \lim_{\lambda \rightarrow 0} (A^\top A + \lambda I)^{-1} A$ in (ii), and defined $\hat{X} = X P^{1/2} \Sigma_\theta^{-1/2}$ in (iii). To proceed, we first assume that $\Sigma_{\theta/P}$ is invertible (i.e. $\lambda_{\min}(\Sigma_{\theta/P})$ is bounded away from 0) and observe the following relation via a leave-one-out argument similar to that in Xu & Hsu (2019),

$$\frac{1}{d} \text{tr} \left(\frac{1}{n} \hat{X}^\top \hat{X} \left(\frac{1}{n} \hat{X}^\top \hat{X} + \lambda \Sigma_{\theta/P}^{-1} \right)^{-2} \right) \quad (\text{D.1})$$

$$\begin{aligned}
&\stackrel{(i)}{=} \frac{1}{d} \sum_{i=1}^n \frac{\frac{1}{n} \hat{x}_i^\top \left(\frac{1}{n} \hat{X}^\top \hat{X} + \lambda \Sigma_{\theta/P}^{-1} \right)^{-2} \hat{x}_i}{\left(1 + \frac{1}{n} \hat{x}_i^\top \left(\frac{1}{n} \hat{X}^\top \hat{X} + \lambda \Sigma_{\theta/P}^{-1} \right)^{-1} \hat{x}_i \right)^2} \\
&\stackrel{(ii)}{\xrightarrow{p}} \frac{\frac{1}{d} \text{tr} \left(\left(\frac{1}{n} \hat{X}^\top \hat{X} + \lambda \Sigma_{\theta/P}^{-1} \right)^{-2} \Sigma_{\theta/P}^{-1/2} \Sigma_{XP} \Sigma_{\theta/P}^{-1/2} \right)}{\left(1 + \frac{1}{n} \text{tr} \left(\left(\frac{1}{n} \bar{X}^\top \bar{X} + \lambda I_d \right)^{-1} \Sigma_{XP} \right) \right)^2}, \quad (\text{D.2})
\end{aligned}$$

where (i) is due to the Woodbury identity and we defined $\left(\frac{1}{n} \hat{X}^\top \hat{X} + \lambda \Sigma_{\theta/P}^{-1} \right)_{-i} = \frac{1}{n} \hat{X}^\top \hat{X} - \frac{1}{n} \hat{x}_i \hat{x}_i^\top + \lambda \Sigma_{\theta/P}^{-1}$ which is independent to \hat{x}_i (see Xu & Hsu (2019, Eq. 58) for details), and in (ii) we used (A3), the convergence to trace (Ledoit & Péché, 2011, Lemma 2.1) and its stability under low-rank perturbation (e.g., see Ledoit & Péché (2011, Eq. 18)) which we elaborate below. In particular, denote $\hat{\Sigma} = \frac{1}{n} \hat{X}^\top \hat{X} + \lambda \Sigma_{\theta/P}^{-1}$, for the denominator we have

$$\begin{aligned}
&\sup_i \left| \frac{\lambda}{n} \text{tr} \left(\hat{\Sigma}^{-1} \Sigma_{\theta/P}^{-1/2} \Sigma_{XP} \Sigma_{\theta/P}^{-1/2} \right) - \frac{\lambda}{n} \text{tr} \left(\hat{\Sigma}_{-i}^{-1} \Sigma_{\theta/P}^{-1/2} \Sigma_{XP} \Sigma_{\theta/P}^{-1/2} \right) \right| \\
&\leq \frac{\lambda}{n} \left\| \Sigma_{\theta/P}^{-1/2} \Sigma_{XP} \Sigma_{\theta/P}^{-1/2} \right\|_2 \sup_i \left| \text{tr} \left(\hat{\Sigma}^{-1} (\hat{\Sigma} - \hat{\Sigma}_{-i}) \hat{\Sigma}_{-i}^{-1} \right) \right| \\
&\leq \frac{\lambda}{n} \left\| \Sigma_{\theta/P}^{-1/2} \Sigma_{XP} \Sigma_{\theta/P}^{-1/2} \right\|_2 \left\| \hat{\Sigma}^{-1} \right\|_2 \sup_i \left\| \hat{\Sigma}_{-i}^{-1} \right\|_2 \text{tr} \left(\hat{\Sigma} - \hat{\Sigma}_{-i} \right) \stackrel{(i)}{\rightarrow} O_p \left(\frac{1}{n} \right),
\end{aligned}$$

where (i) is due to the definition of $\hat{\Sigma}_{-i}$ and (A1)(A3). The result on the numerator can be obtained via a similar calculation, the details of which we omit.

Note that the denominator can be evaluated by previous results (e.g. Dobriban et al. (2018, Theorem 2.1)) as follows,

$$\frac{1}{n} \text{tr} \left(\left(\frac{1}{n} \bar{X}^\top \bar{X} + \lambda I_d \right)^{-1} \Sigma_{XP} \right) \xrightarrow{a.s.} \frac{1}{\lambda m(-\lambda)} - 1. \quad (\text{D.3})$$

$$\begin{aligned}
& + 4\mathbb{E}\left[\frac{h_\alpha(g_\alpha - h_\alpha)}{(1 + h_\alpha)^3}\right]\mathbb{E}\left[\frac{g_\alpha - h_\alpha}{(1 + h_\alpha)^2}\right]\mathbb{E}\left[\frac{h_\alpha}{(1 + h_\alpha)^2}\right] \\
& \stackrel{(i)}{\leq} -4\sqrt{\mathbb{E}\left[\frac{(g_\alpha - h_\alpha)^2}{(1 + h_\alpha)^3}\right]\mathbb{E}\left[\frac{h_\alpha^2}{(1 + h_\alpha)^3}\right]\left(\mathbb{E}\left[\frac{g_\alpha - h_\alpha}{(1 + h_\alpha)^2}\right]\right)^2\left(\mathbb{E}\left[\frac{h_\alpha}{(1 + h_\alpha)^2}\right]\right)^2} \\
& + 4\mathbb{E}\left[\frac{h_\alpha(g_\alpha - h_\alpha)}{(1 + h_\alpha)^3}\right]\mathbb{E}\left[\frac{g_\alpha - h_\alpha}{(1 + h_\alpha)^2}\right]\mathbb{E}\left[\frac{h_\alpha}{(1 + h_\alpha)^2}\right] \stackrel{(ii)}{\leq} 0,
\end{aligned}$$

where (i) is due to AM-GM and (ii) due to Cauchy-Schwarz on the first term. Note that the two equalities hold when $g_\alpha = h_\alpha$, from which one can easily deduce that the optimum is achieved when $v_{xp} \stackrel{a.s.}{=} v_x v_\theta$, and thus we know that the proposed P is the optimal preconditioner that is codiagonalizable with Σ_X . \square

D.4 PROOF OF PROPOSITION 3

Proof. Via calculation similar to [Hastie et al. \(2019, Section 5\)](#), the bias can be decomposed as

$$\begin{aligned}
\mathbb{E}[B(\hat{\theta}_P)] &= \mathbb{E}_{\mathbf{x}, \hat{\mathbf{x}}, \theta^*, \theta^c} \left[\left(\mathbf{x}^\top P \mathbf{X}^\top (\mathbf{X} P \mathbf{X}^\top)^{-1} (\mathbf{X} \theta^* + \mathbf{X}^c \theta^c) - (\mathbf{x}^\top \theta^* + \hat{\mathbf{x}}^\top \theta^c) \right)^2 \right] \\
&\stackrel{(i)}{=} \mathbb{E}_{\mathbf{x}, \theta^*} \left[\left(\mathbf{x}^\top P \mathbf{X}^\top (\mathbf{X} P \mathbf{X}^\top)^{-1} \mathbf{X} \theta^* - \mathbf{x}^\top \theta^* \right)^2 \right] + \mathbb{E}_{\mathbf{x}^c, \theta^c} \left[(\hat{\mathbf{x}}^\top \theta^c)^2 \right] \\
&\quad + \mathbb{E}_{\mathbf{x}, \theta^c} \left[\left(\mathbf{x}^\top P \mathbf{X}^\top (\mathbf{X} P \mathbf{X}^\top)^{-1} \mathbf{X}^c \theta^c \theta^{c\top} \mathbf{X}^{c\top} (\mathbf{X} P \mathbf{X}^\top)^{-1} \mathbf{X} P \mathbf{x} \right)^2 \right] \\
&\stackrel{(ii)}{\rightarrow} B_\theta(\hat{\theta}_P) + \frac{1}{d^c} \text{tr}(\Sigma_X^c \Sigma_\theta^c) (1 + V(\hat{\theta}_P)),
\end{aligned}$$

where we used the independence of $\mathbf{x}, \hat{\mathbf{x}}$ and θ^*, θ^c in (i), and (A1-3) as well as the definition of the well-specified bias $B_\theta(\hat{\theta}_P)$ and variance $V(\hat{\theta}_P)$ in (ii). \square

D.5 PROOF OF PROPOSITION 4

Proof. We first outline a more general setup where $P_\alpha = f(\Sigma_X; \alpha)$ for continuous and differentiable function of α and f applied to eigenvalues of Σ_x . For any interval $\mathcal{I} \subseteq [0, 1]$, we claim

- Suppose all four functions $\frac{1}{xf(x; \alpha)}$, $f(x; \alpha)$, $\frac{\partial f(x; \alpha)}{\partial \alpha} / f(x; \alpha)$ and $x \frac{\partial f(x; \alpha)}{\partial \alpha}$ are decreasing functions of x on the support of v_x for all $\alpha \in \mathcal{I}$. In addition, $\frac{\partial f(x; \alpha)}{\partial \alpha} \geq 0$ on the support of v_x for all $\alpha \in \mathcal{I}$. Then the stationary bias is an increasing function of α on \mathcal{I} .
- For all $\alpha \in \mathcal{I}$, suppose $xf(x; \alpha)$ is a monotonic function of x on the support of v_x and $\frac{\partial f(x; \alpha)}{\partial \alpha} / f(x; \alpha)$ is a decreasing function of x on the support of v_x . Then the stationary variance is a decreasing function of α on \mathcal{I} .

Let us verify the three choices of P_α in Proposition 4 one by one.

- When $P_\alpha = (1 - \alpha)\mathbf{I}_d + \alpha(\Sigma_X)^{-1}$, the corresponding $f(x; \alpha)$ is $(1 - \alpha) + \alpha x$. It is clear that it satisfies all conditions in (a) and (b) for all $\alpha \in [0, 1]$. Hence, the stationary variance is a decreasing function and the stationary bias is an increasing function of $\alpha \in [0, 1]$.
- When $P_\alpha = (\Sigma_X)^{-\alpha}$, the corresponding $f(x; \alpha)$ is $x^{-\alpha}$. It is clear that it satisfies all conditions in (a) and (b) for all $\alpha \in [0, 1]$ except for the condition that $x \frac{\partial f(x; \alpha)}{\partial \alpha} = -x^{1-\alpha} \ln x$ is a decreasing function of x . Note that $x \frac{\partial f(x; \alpha)}{\partial \alpha} = -x^{1-\alpha} \ln x$ is a decreasing function of x on the support of v_x only for $\alpha \geq \frac{\ln(\kappa)-1}{\ln(\kappa)}$ where $\kappa = \sup v_x / \inf v_x$. Hence, the stationary variance is a decreasing function of $\alpha \in [0, 1]$ and the stationary bias is an increasing function of $\alpha \in [\max(0, \frac{\ln(\kappa)-1}{\ln(\kappa)}), 1]$.

- When $\mathbf{P}_\alpha = (\alpha \Sigma_{\mathbf{X}} + (1-\alpha) \mathbf{I}_d)^{-1}$, the corresponding $f(x; \alpha)$ is $1/(\alpha x + (1-\alpha))$. It is clear that it satisfies all conditions in (a) and (b) for all $\alpha \in [0, 1]$ except for the condition that $x \frac{\partial f(x; \alpha)}{\partial \alpha} = \frac{x(1-x)}{(\alpha x + (1-\alpha))^2}$ is a decreasing function of x . Note that $x \frac{\partial f(x; \alpha)}{\partial \alpha} = \frac{x(1-x)}{(\alpha x + (1-\alpha))^2}$ is a decreasing function of x on the support of v_x only for $\alpha \geq \frac{\kappa-2}{\kappa-1}$. Hence, the stationary variance is a decreasing function of $\alpha \in [0, 1]$ and the stationary bias is an increasing function of $\alpha \in [\max(0, \frac{\kappa-2}{\kappa-1}), 1]$.

To show (a) and (b), note that under the conditions on $\Sigma_{\mathbf{x}}$ and $\Sigma_{\boldsymbol{\theta}}$ assumed in Proposition 4, the stationary bias $B(\hat{\boldsymbol{\theta}}_{\mathbf{P}_\alpha})$ and the stationary variance $V(\hat{\boldsymbol{\theta}}_{\mathbf{P}_\alpha})$ can be simplified to

$$B(\hat{\boldsymbol{\theta}}_{\mathbf{P}_\alpha}) = \frac{m'_\alpha(0)}{m_\alpha^2(0)} \mathbb{E} \frac{v_x}{(1 + v_x f(v_x; \alpha) m_\alpha(0))^2} \quad \text{and} \quad V(\hat{\boldsymbol{\theta}}_{\mathbf{P}_\alpha}) = \sigma^2 \cdot \left(\frac{m'_\alpha(0)}{m_\alpha^2(0)} - 1 \right),$$

where $m_\alpha(z)$ and $m'_\alpha(z)$ satisfy

$$1 = -z m_\alpha(z) + \gamma \mathbb{E} \frac{v_x f(v_x; \alpha) m_\alpha(z)}{1 + v_x f(v_x; \alpha) m_\alpha(z)} \quad (\text{D.7})$$

$$\frac{m'_\alpha(z)}{m_\alpha^2(z)} = \frac{1}{1 - \gamma \mathbb{E} \left(\frac{f(v_x; \alpha) m_\alpha(z)}{1 + f(v_x; \alpha) m_\alpha(z)} \right)^2}. \quad (\text{D.8})$$

For notation convenience, let $f_\alpha := v_x f(v_x; \alpha)$. From (D.8), we have the following equivalent expressions.

$$B(\hat{\boldsymbol{\theta}}_{\mathbf{P}_\alpha}) = \frac{\mathbb{E} \frac{v_x}{(1 + f_\alpha m_\alpha(0))^2}}{1 - \gamma \mathbb{E} \left(\frac{f_\alpha m_\alpha(0)}{1 + f_\alpha m_\alpha(0)} \right)^2}, \quad (\text{D.9})$$

$$V(\hat{\boldsymbol{\theta}}_{\mathbf{P}_\alpha}) = \sigma^2 \left(\frac{1}{1 - \gamma \mathbb{E} \left(\frac{f_\alpha m_\alpha(0)}{1 + f_\alpha m_\alpha(0)} \right)^2} - 1 \right). \quad (\text{D.10})$$

We first show that (b) holds. Note that from (D.10), we have

$$\frac{\partial V(\hat{\boldsymbol{\theta}}_{\mathbf{P}_\alpha})}{\partial \alpha} = \gamma \sigma^2 \left(\frac{1}{1 - \gamma \mathbb{E} \left(\frac{f_\alpha m_\alpha(0)}{1 + f_\alpha m_\alpha(0)} \right)^2} \right)^2 \mathbb{E} \left[\frac{2 f_\alpha m_\alpha(0)}{(1 + f_\alpha m_\alpha(0))^3} \left(f_\alpha \frac{\partial m_\alpha(z)}{\partial \alpha} \Big|_{z=0} + \frac{\partial f_\alpha}{\partial \alpha} m_\alpha(0) \right) \right]. \quad (\text{D.11})$$

To calculate $\frac{\partial m_\alpha(z)}{\partial \alpha} \Big|_{z=0}$, we take derivatives with respect to α on both sides of (D.7),

$$0 = \gamma \mathbb{E} \left[\frac{1}{(1 + f_\alpha m_\alpha(0))^2} \cdot \left(f_\alpha \frac{\partial m_\alpha(z)}{\partial \alpha} \Big|_{z=0} + \frac{\partial f_\alpha}{\partial \alpha} m_\alpha(0) \right) \right]. \quad (\text{D.12})$$

Therefore, plugging (D.12) into (D.11) yields

$$\begin{aligned} \frac{\partial V(\hat{\boldsymbol{\theta}}_{\mathbf{P}_\alpha})}{\partial \alpha} &= 2\gamma \sigma^2 \left(\frac{m_\alpha(0)}{1 - \gamma \mathbb{E} \left(\frac{f_\alpha m_\alpha(0)}{1 + f_\alpha m_\alpha(0)} \right)^2} \right)^2 \left(\mathbb{E} \frac{f_\alpha}{(1 + f_\alpha m_\alpha(0))^2} \right)^{-1} \\ &\quad \times \left(\mathbb{E} \frac{f_\alpha \frac{\partial f_\alpha}{\partial \alpha}}{(1 + f_\alpha m_\alpha(0))^3} \mathbb{E} \frac{f_\alpha}{(1 + f_\alpha m_\alpha(0))^2} - \mathbb{E} \frac{f_\alpha^2}{(1 + f_\alpha m_\alpha(0))^3} \mathbb{E} \frac{\frac{\partial f_\alpha}{\partial \alpha}}{(1 + f_\alpha m_\alpha(0))^2} \right) \end{aligned}$$

Thus showing $V(\hat{\boldsymbol{\theta}}_{\mathbf{P}_\alpha})$ is a decreasing function of α is equivalent to showing that

$$\mathbb{E} \frac{f_\alpha^2}{(1 + f_\alpha m_\alpha(0))^3} \mathbb{E} \frac{\frac{\partial f_\alpha}{\partial \alpha}}{(1 + f_\alpha m_\alpha(0))^2} \geq \mathbb{E} \frac{f_\alpha \frac{\partial f_\alpha}{\partial \alpha}}{(1 + f_\alpha m_\alpha(0))^3} \mathbb{E} \frac{f_\alpha}{(1 + f_\alpha m_\alpha(0))^2}. \quad (\text{D.13})$$

Let μ_x be the probability measure of v_x . We define a new measure $\tilde{\mu}_x = \frac{f_\alpha \mu_x}{(1+f_\alpha m_\alpha(0))^2}$, and let \tilde{v}_x follow the new measure. Since $\frac{\partial f(x;\alpha)}{\partial \alpha} / f(x;\alpha)$ is a decreasing function of x and $xf(x;\alpha)$ is a monotonic function of x ,

$$\mathbb{E} \frac{\tilde{v}_x f(\tilde{v}_x; \alpha)}{1 + \tilde{v}_x f(\tilde{v}_x; \alpha) m_\alpha(0)} \mathbb{E} \frac{\frac{\partial \tilde{v}_x f(\tilde{v}_x; \alpha)}{\partial \alpha}}{\tilde{v}_x f(\tilde{v}_x; \alpha)} \geq \mathbb{E} \frac{\frac{\partial \tilde{v}_x f(\tilde{v}_x; \alpha)}{\partial \alpha}}{1 + \tilde{v}_x f(\tilde{v}_x; \alpha) m_\alpha(0)}.$$

Changing \tilde{v}_x back to v_x , we arrive at (D.13) and thus (b).

For the bias term $B(\hat{\theta}_{P_\alpha})$, note that from (D.7) and (D.9), we have

$$\begin{aligned} \frac{\partial B(\hat{\theta}_{P_\alpha})}{\partial \alpha} &= \frac{1}{\gamma} \left(\frac{1}{\gamma} - \mathbb{E} \left(\frac{f_\alpha m_\alpha(0)}{1 + f_\alpha m_\alpha(0)} \right)^2 \right)^{-2} \\ &\quad \times \left(-\mathbb{E} \left[2 \frac{v_x}{(1 + f_\alpha m_\alpha(0))^3} \cdot \left(f_\alpha \frac{\partial m_\alpha(z)}{\partial \alpha} \Big|_{z=0} + \frac{\partial f_\alpha}{\partial \alpha} m_\alpha(0) \right) \right] \mathbb{E} \frac{f_\alpha m_\alpha(0)}{(1 + f_\alpha m_\alpha(0))^2} \right. \\ &\quad \left. + \mathbb{E} \frac{v_x}{(1 + f_\alpha m_\alpha(0))^2} \mathbb{E} \left[2 \frac{f_\alpha m_\alpha(0)}{(1 + f_\alpha m_\alpha(0))^3} \cdot \left(f_\alpha \frac{\partial m_\alpha(z)}{\partial \alpha} \Big|_{z=0} + \frac{\partial f_\alpha}{\partial \alpha} m_\alpha(0) \right) \right] \right). \end{aligned} \quad (\text{D.14})$$

Similarly, we combine (D.12) and (D.14) and simplify the expression. To verify $B(\hat{\theta}_{P_\alpha})$ is an increasing function of α , we need to show that

$$\begin{aligned} 0 &\leq \left(\mathbb{E} \frac{v_x f_\alpha m_\alpha(0)}{(1 + f_\alpha m_\alpha(0))^3} \mathbb{E} \frac{\frac{\partial f_\alpha}{\partial \alpha}}{(1 + f_\alpha m_\alpha(0))^2} - \mathbb{E} \frac{v_x \frac{\partial f_\alpha}{\partial \alpha}}{(1 + f_\alpha m_\alpha(0))^3} \mathbb{E} \frac{f_\alpha m_\alpha(0)}{(1 + f_\alpha m_\alpha(0))^2} \right) \mathbb{E} \frac{f_\alpha m_\alpha(0)}{(1 + f_\alpha m_\alpha(0))^2} \\ &\quad - \mathbb{E} \frac{v_x}{(1 + f_\alpha m_\alpha(0))^2} \left(\mathbb{E} \frac{(f_\alpha m_\alpha(0))^2}{(1 + f_\alpha m_\alpha(0))^3} \mathbb{E} \frac{\frac{\partial f_\alpha}{\partial \alpha}}{(1 + f_\alpha m_\alpha(0))^2} - \mathbb{E} \frac{f_\alpha m_\alpha(0) \frac{\partial f_\alpha}{\partial \alpha}}{(1 + f_\alpha m_\alpha(0))^3} \mathbb{E} \frac{f_\alpha m_\alpha(0)}{(1 + f_\alpha m_\alpha(0))^2} \right), \end{aligned} \quad (\text{D.15})$$

Let $h_\alpha \triangleq f_\alpha m_\alpha(0) = v_x f(v_x; \alpha) m_\alpha(0)$ and $g_\alpha \triangleq \frac{\partial f_\alpha}{\partial \alpha} = v_x \frac{\partial f(v_x; \alpha)}{\partial \alpha}$. Then (D.15) can be further simplified to the following equation

$$\begin{aligned} 0 &\leq \underbrace{\mathbb{E} \frac{v_x h_\alpha}{(1 + h_\alpha)^3} \mathbb{E} \frac{g_\alpha}{(1 + h_\alpha)^3} \mathbb{E} \frac{h_\alpha}{(1 + h_\alpha)^3} - \mathbb{E} \frac{v_x}{(1 + h_\alpha)^3} \mathbb{E} \frac{g_\alpha}{(1 + h_\alpha)^3} \mathbb{E} \frac{h_\alpha^2}{(1 + h_\alpha)^3}}_{\text{part 1}} \\ &\quad + \underbrace{\mathbb{E} \frac{v_x}{(1 + h_\alpha)^3} \mathbb{E} \frac{g_\alpha h_\alpha}{(1 + h_\alpha)^3} \mathbb{E} \frac{h_\alpha}{(1 + h_\alpha)^3} - \mathbb{E} \frac{v_x g_\alpha}{(1 + h_\alpha)^3} \mathbb{E} \frac{h_\alpha}{(1 + h_\alpha)^3} \mathbb{E} \frac{h_\alpha}{(1 + h_\alpha)^3}}_{\text{part 2}} \\ &\quad + \underbrace{2 \mathbb{E} \frac{v_x h_\alpha}{(1 + h_\alpha)^3} \mathbb{E} \frac{g_\alpha h_\alpha}{(1 + h_\alpha)^3} \mathbb{E} \frac{h_\alpha}{(1 + h_\alpha)^3} - 2 \mathbb{E} \frac{v_x g_\alpha}{(1 + h_\alpha)^3} \mathbb{E} \frac{h_\alpha^2}{(1 + h_\alpha)^3} \mathbb{E} \frac{h_\alpha}{(1 + h_\alpha)^3}}_{\text{part 3}} \\ &\quad + \underbrace{\mathbb{E} \frac{v_x h_\alpha}{(1 + h_\alpha)^3} \mathbb{E} \frac{g_\alpha h_\alpha}{(1 + h_\alpha)^3} \mathbb{E} \frac{h_\alpha^2}{(1 + h_\alpha)^3} - \mathbb{E} \frac{v_x g_\alpha}{(1 + h_\alpha)^3} \mathbb{E} \frac{h_\alpha^2}{(1 + h_\alpha)^3} \mathbb{E} \frac{h_\alpha^2}{(1 + h_\alpha)^3}}_{\text{part 4}}. \end{aligned} \quad (\text{D.16})$$

Note that under condition of (a), we know that both h_α and v_x/h_α are increasing functions of v_x ; and both g_α/h_α and g_α are decreasing functions of v_x . Hence, with calculation similar to (D.13), we know part 1,2,3,4 in (D.16) are all non-negative, and therefore (D.16) holds. \square

Remark. The above characterization provides sufficient but not necessary conditions for the monotonicity of the bias term. In general, the expression of the bias is rather opaque, and determining the sign of its derivative can be tedious, except for certain special cases (e.g. $\gamma = 2$ and the eigenvalues of $\Sigma_{\mathbf{X}}$ are two equally weighted point masses, for which m_α has a simple form and one may analytically check the monotonicity). We conjecture that the bias is monotone for $\alpha \in [0, 1]$ for a much wider class of $\Sigma_{\mathbf{X}}$, as shown in Figure 23.

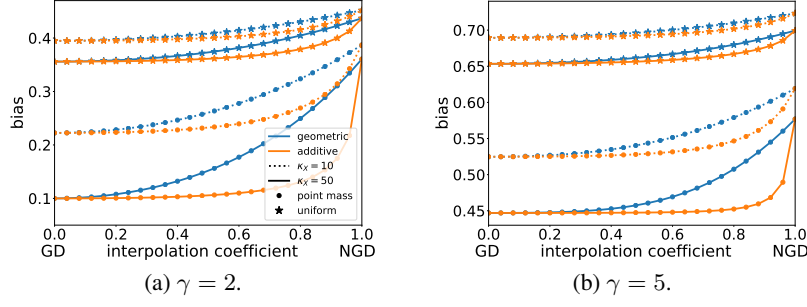


Figure 23: Illustration of the monotonicity of the bias term under $\Sigma_\theta = I_d$. We consider two distributions of eigenvalues for Σ_X : two equally weighted point masses (circle) and a uniform distribution (star), and vary the condition number κ_X and overparameterization level γ . In all cases the bias is monotone in $\alpha \in [0, 1]$.

D.6 PROOF OF PROPOSITION 5

Proof. Taking the derivative of $V(\theta_P(t))$ w.r.t. time yields (omitting the scalar σ^2),

$$\begin{aligned} \frac{dV(\theta_P(t))}{dt} &= \frac{d}{dt} \left\| \Sigma_X^{1/2} P X^\top \left(I_n - \exp\left(-\frac{t}{n} X P X^\top\right) \right) (X P X^\top)^{-1} \right\|_F^2 \\ &\stackrel{(i)}{=} \frac{1}{n} \text{tr} \left(\underbrace{\Sigma_{XP} \bar{X}^\top S_P \exp\left(-\frac{t}{n} S_P\right) S_P^{-2} \left(I_n - \exp\left(-\frac{t}{n} S_P\right) \right) \bar{X}}_{p.s.d.} \right) \stackrel{(ii)}{>} 0, \end{aligned}$$

where we defined $\bar{X} = X P^{1/2}$ and $S_P = X P X^\top$ in (i), and (ii) is due to (A2-3) the inequality $\text{tr}(AB) \geq \lambda_{\min}(A) \text{tr}(B)$ for positive semi-definite A and B . \square

D.7 PROOF OF PROPOSITION 6

Proof. Recall the definition of the bias (well-specified) of $\hat{\theta}_P(t)$,

$$\begin{aligned} B(\theta_P(t)) &\stackrel{(i)}{=} \frac{1}{d} \text{tr} \left(\Sigma_\theta \left(I_d - P X^\top W_P(t) S_P^{-1} X \right)^\top \Sigma_X \left(I_d - P X^\top W_P(t) S_P^{-1} X \right) \right) \\ &\stackrel{(ii)}{=} \frac{1}{d} \text{tr} \left(\Sigma_{\theta/P} \left(I_d - \bar{X}^\top W_P(t) S_P^{-1} \bar{X} \right)^\top \Sigma_{XP} \left(I_d - \bar{X}^\top W_P(t) S_P^{-1} \bar{X} \right) \right) \\ &\stackrel{(iii)}{\geq} \frac{1}{d} \text{tr} \left(\left(\Sigma_{XP}^{1/2} \left(I_d - \bar{X}^\top W_P(t) S_P^{-1} \bar{X} \right) \Sigma_{\theta/P}^{1/2} \right)^2 \right), \end{aligned} \quad (\text{D.17})$$

where we defined $S_P = X P X^\top$, $W_P(t) = I_n - \exp(-\frac{t}{n} S_P)$ in (i), $\bar{X} = X P^{1/2}$ in (ii), and (iii) is due to the inequality $\text{tr}(A^\top A) \geq \text{tr}(A^2)$.

When $\Sigma_X = \Sigma_\theta^{-1}$, i.e. NGD achieves lowest stationary bias, (D.17) simplifies to

$$B(\theta_P(t)) \geq \frac{1}{d} \text{tr} \left(\left(I_d - \bar{X}^\top W_P(t) S_P^{-1} \bar{X} \right)^2 \right) = \left(1 - \frac{1}{\gamma} \right) + \frac{1}{d} \sum_{i=1}^n \exp\left(-\frac{t}{n} \bar{\lambda}_i\right)^2, \quad (\text{D.18})$$

where $\bar{\lambda}$ is the eigenvalue of S_P . On the other hand, since $F = \Sigma_X$, for the NGD iterate $\hat{\theta}_{F^{-1}}(t)$ we have

$$B(\theta_{F^{-1}}(t)) = \frac{1}{d} \text{tr} \left(\left(I_d - \hat{X}^\top W_{F^{-1}}(t) S_{F^{-1}}^{-1} \hat{X} \right)^2 \right) = \left(1 - \frac{1}{\gamma} \right) + \frac{1}{d} \sum_{i=1}^n \exp\left(-\frac{t}{n} \hat{\lambda}_i\right)^2 \quad (\text{D.19})$$

where $\hat{X} = X \Sigma_X^{-1/2}$ and $\hat{\lambda}$ is the eigenvalue of $S_{F^{-1}} = \hat{X} \hat{X}^\top$. Comparing (D.18)(D.19), we see that given $\hat{\theta}_P(t)$ at a fixed t , if we run NGD for time $T > \frac{\bar{\lambda}_{\max}}{\bar{\lambda}_{\min}} t$ (note that $T/t = O(1)$ by

(A2-3)), then we have $B(\boldsymbol{\theta}_P(t)) \geq B(\boldsymbol{\theta}_{F^{-1}}(T))$ for any P satisfying (A3). This thus implies that $B^{\text{opt}}(\boldsymbol{\theta}_P) \geq B^{\text{opt}}(\boldsymbol{\theta}_{F^{-1}})$.

On the other hand, when $\boldsymbol{\Sigma}_\theta = \mathbf{I}_d$, we can show that the bias term of GD is monotonically decreasing through time by taking its derivative,

$$\begin{aligned} \frac{d}{dt} B(\boldsymbol{\theta}_I(t)) &= \frac{1}{d} \frac{d}{dt} \text{tr} \left(\left(\mathbf{I}_d - \mathbf{X}^\top \mathbf{W}_I(t) \mathbf{S}_I^{-1} \mathbf{X} \right)^\top \boldsymbol{\Sigma}_X \left(\mathbf{I}_d - \mathbf{X}^\top \mathbf{W}_I(t) \mathbf{S}_I^{-1} \mathbf{X} \right) \right) \\ &= -\frac{1}{nd} \text{tr} \left(\underbrace{\boldsymbol{\Sigma}_X \mathbf{X}^\top \mathbf{S} \exp\left(-\frac{t}{n} \mathbf{S}\right) \mathbf{S}^{-1} \mathbf{X} \left(\mathbf{I}_d - \mathbf{X}^\top \mathbf{W}_I(t) \mathbf{S}_I^{-1} \mathbf{X} \right)}_{p.s.d.} \right) < 0. \end{aligned} \quad (\text{D.20})$$

Similarly, one can verify that the expected bias of NGD is monotonically decreasing for all choices of $\boldsymbol{\Sigma}_X$ and $\boldsymbol{\Sigma}_\theta$ satisfying (A2-4),

$$\begin{aligned} &\frac{d}{dt} \text{tr} \left(\boldsymbol{\Sigma}_\theta \left(\mathbf{I}_d - \mathbf{F}^{-1} \mathbf{X}^\top \mathbf{W}_{F^{-1}}(t) \mathbf{S}_{F^{-1}}^{-1} \mathbf{X} \right)^\top \boldsymbol{\Sigma}_X \left(\mathbf{I}_d - \mathbf{F}^{-1} \mathbf{X}^\top \mathbf{W}_{F^{-1}}(t) \mathbf{S}_{F^{-1}}^{-1} \mathbf{X} \right) \right) \\ &= \frac{d}{dt} \text{tr} \left(\boldsymbol{\Sigma}_{X\theta} \left(\mathbf{I}_d - \hat{\mathbf{X}}^\top \mathbf{W}_{F^{-1}}(t) \mathbf{S}_{F^{-1}}^{-1} \hat{\mathbf{X}} \right)^\top \left(\mathbf{I}_d - \hat{\mathbf{X}}^\top \mathbf{W}_{F^{-1}}(t) \mathbf{S}_{F^{-1}}^{-1} \hat{\mathbf{X}} \right) \right) \stackrel{(i)}{<} 0, \end{aligned}$$

where (i) follows from calculation similar to (D.20). Since the expected bias is decreasing through time for both GD and NGD when $\boldsymbol{\Sigma}_\theta = \mathbf{I}_d$, and from Theorem 2 we know that $B(\hat{\boldsymbol{\theta}}_I) \leq B(\hat{\boldsymbol{\theta}}_{F^{-1}})$, we conclude that $B^{\text{opt}}(\boldsymbol{\theta}_I) \leq B^{\text{opt}}(\boldsymbol{\theta}_{F^{-1}})$. \square

D.8 PROOF OF THEOREM 7

D.8.1 SETUP AND MAIN RESULT

We restate the setting and assumptions. \mathcal{H} is an RKHS included in $L_2(P_X)$ equipped with a bounded kernel function k satisfying $\sup_{\text{supp}(P_X)} k(\mathbf{x}, \mathbf{x}) \leq 1$. $K_{\mathbf{x}} \in \mathcal{H}$ is the Riesz representation of the kernel function $k(\mathbf{x}, \cdot)$, that is, $k(\mathbf{x}, \mathbf{y}) = \langle K_{\mathbf{x}}, K_{\mathbf{y}} \rangle_{\mathcal{H}}$. S is the canonical embedding operator from \mathcal{H} to $L_2(P_X)$. We write $\Sigma = S^* S : \mathcal{H} \rightarrow \mathcal{H}$ and $L = S S^*$. Note that the boundedness of the kernel gives $\|Sf\|_{L_2(P_X)} \leq \sup_{\mathbf{x}} |f(\mathbf{x})| = \sup_{\mathbf{x}} |\langle K_{\mathbf{x}}, f \rangle| \leq \|K_{\mathbf{x}}\|_{\mathcal{H}} \|f\|_{\mathcal{H}} \leq \|f\|_{\mathcal{H}}$. Hence we know $\|\Sigma\| \leq 1$ and $\|L\| \leq 1$. Our analysis will be made under the following assumptions.

- There exist $r \in (0, \infty)$ and $M > 0$ such that $f^* = L^r h^*$ for some $h^* \in L_2(P_X)$ and $\|f^*\|_{\infty} \leq M$.
- There exists $s > 1$ s.t. $\text{tr}(\Sigma^{1/s}) < \infty$ and $2r + s^{-1} > 1$.
- There exist $\mu \in [s^{-1}, 1]$ and $C_\mu > 0$ such that $\sup_{\mathbf{x} \in \text{supp}(P_X)} \|\Sigma^{1/2-1/\mu} K_{\mathbf{x}}\|_{\mathcal{H}} \leq C_\mu$.

(A5)(A6) are standard regularity assumptions in the literature that provide capacity control of the RKHS (e.g., see Caponnetto & De Vito (2007); Pillaud-Vivien et al. (2018)). It is also worth noting that most previous works (including Rudi et al. (2017)) consider $r \geq 1/2$ which implies that $f^* \in \mathcal{H}$.

The training data is generated as $y_i = f^*(\mathbf{x}_i) + \varepsilon_i$, where ε_i is an i.i.d. noise satisfying $|\varepsilon_i| \leq \sigma$ almost surely. Let $\mathbf{y} \in \mathbb{R}^n$ be the label vector. We identify \mathbb{R}^n with $L_2(P_n)$ and define

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{x}_i} \otimes K_{\mathbf{x}_i} : \mathcal{H} \rightarrow \mathcal{H}, \quad \hat{S}^* Y = \frac{1}{n} \sum_{i=1}^n Y_i K_{\mathbf{x}_i}, \quad (Y \in L_2(P_n)).$$

We consider the following preconditioned update on $f_t \in \mathcal{H}$:

$$f_t = f_{t-1} - \eta(\Sigma + \lambda I)^{-1}(\hat{\Sigma} f_{t-1} - \hat{S}^* Y), \quad f_0 = 0.$$

We briefly comment on how our analysis differs from Rudi et al. (2017), which showed that a preconditioned update (the FALKON algorithm) for kernel ridge regression can also achieve accelerated convergence in the population risk. We emphasize the following differences.

- The two algorithms optimize different objectives, as highlighted by the different role of the “ridge” coefficient (α in our update, λ in FALKON). In FALKON, λ turns the objective into kernel ridge regression; whereas in our (4.1), α controls the interpolation between GD and NGD. As we aim to study how the preconditioner affects generalization, it is important that we look at the objective in its original (instead of regularized) form.
- To elaborate on the first point, since FALKON minimizes a regularized objective, it would not overfit even after large number of gradient steps, yet it is unclear how preconditioning impacts the generalization error (i.e., any preconditioner may generalize well with proper regularization). In contrast, we consider the “unregularized” objective, and thus early stopping plays a crucial role in avoiding overfitting; this is different from most standard analysis on gradient descent.
- Algorithm-wise, the two updates employ different preconditioners. FALKON involves inverting the kernel matrix K defined on the training points, whereas we consider the population covariance operator Σ , which is consistent with our earlier discussion on the population Fisher in Section 3.
- In terms of the theoretical setup, our analysis allows for $r < 1/2$, whereas Rudi et al. (2017) and many other previous works assumed $r \in [1/2, 1]$, as commented above.

We aim to show the following theorem:

Theorem (Restatement of Theorem 7). *Given (A4-6), if the sample size n is sufficiently large so that $1/(n\lambda) \ll 1$, then for $\eta < \|\Sigma\|$ with $\eta t \geq 1$ and $0 < \delta < 1$ and $0 < \lambda < 1$, it holds that*

$$\|Sf_t - f^*\|_{L_2(P_X)}^2 \leq C(B(t) + V(t)),$$

with probability $1 - 3\delta$, where C is a constant and

$$B(t) := \exp(-\eta t) \vee \left(\frac{\lambda}{\eta t}\right)^{2r},$$

$$V(t) := V_1(t) + (1 + \eta t) \left(\frac{\lambda^{-1} B(t) + \sigma^2 \text{tr}(\Sigma^{\frac{1}{s}}) \lambda^{-\frac{1}{s}}}{n} + \frac{\lambda^{-1} (\sigma + M + (1 + t\eta) \lambda^{-(\frac{1}{2}-r)_+})^2}{n^2} \right) \log(1/\delta)^2,$$

in which

$$V_1(t) := \left[\exp(-\eta t) \vee \left(\frac{\lambda}{\eta t}\right)^{2r} + (t\eta)^2 \left(\frac{\beta' (1 \vee \lambda^{2r-\mu}) \text{tr}(\Sigma^{\frac{1}{s}}) \lambda^{-\frac{1}{s}}}{n} + \frac{\beta'^2 (1 + \lambda^{-\mu} (1 \vee \lambda^{2r-\mu}))}{n^2} \right) \right] (1 + t\eta)^2,$$

for $\beta' = \log\left(\frac{28C_\mu^2 (2^{2r-\mu} \vee \lambda^{-\mu+2r}) \text{tr}(\Sigma^{1/s}) \lambda^{-1/s}}{\delta}\right)$. When $r \geq 1/2$, if we set $\lambda = n^{-\frac{s}{2rs+1}} =: \lambda^*$ and $t = \Theta(\log(n))$, then the overall convergence rate becomes

$$\|Sg_t - f^*\|_{L_2(P_X)}^2 = \tilde{O}_p\left(n^{-\frac{2rs}{2rs+1}}\right),$$

which is the minimax optimal rate ($\tilde{O}_p(\cdot)$ hides a poly-log(n) factor). On the other hand, when $r < 1/2$, the bound is also $\tilde{O}_p\left(n^{-\frac{2rs}{2rs+1}}\right)$ except the term $V_1(t)$. In this case, if $2r \geq \mu$ holds additionally, we have $V_t(t) = \tilde{O}_p\left(n^{-\frac{2rs}{2rs+1}}\right)$, which again recovers the optimal rate.

Note that if the GD (with iterates \tilde{f}_t) is employed, from previous work (Lin & Rosasco, 2017) we know that the bias term $\left(\frac{\lambda}{\eta t}\right)^{2r}$ is replaced by $\left(\frac{1}{\eta t}\right)^{2r}$, and therefore the upper bound translates to

$$\|S\tilde{f}_t - f^*\|_{L_2(P_X)}^2 \leq C \left\{ (\eta t)^{-2r} + \frac{1}{n} \left(\text{tr}(\Sigma^{1/s}) (\eta t)^{1/s} + \frac{\eta t}{n} \right) \left(\sigma^2 + \left(\frac{1}{\eta t}\right)^{2r} + \frac{M^2 + (\eta t)^{-(2r-1)}}{n} \right) \right\},$$

with high probability. In other words, by the condition $\eta = O(1)$, we need $t = \Theta(n^{\frac{2rs}{2rs+1}})$ steps to sufficiently diminish the bias term. In contrast, the preconditioned update that interpolates between GD and NGD (4.1) only require $t = O(\log(n))$ steps to make the bias term negligible. This is because the NGD amplifies the high frequency component and rapidly captures the detailed “shape” of the target function f^* .

D.8.2 PROOF OF MAIN RESULT

Proof. We follow the proof strategy of [Lin & Rosasco \(2017\)](#). First we define a reference optimization problem with iterates \bar{f}_t that directly minimize the population risk:

$$\bar{f}_t = \bar{f}_{t-1} - \eta(\Sigma + \lambda I)^{-1}(\Sigma \bar{f}_{t-1} - S^* f^*), \quad \bar{f}_0 = 0.$$

Note that $\mathbb{E}[f_t] = \bar{f}_t$. In addition, we define the degrees of freedom and its related quantity as

$$\mathcal{N}_\infty(\lambda) := \mathbb{E}_{\mathbf{x}}[\langle K_{\mathbf{x}}, \Sigma_\lambda^{-1} K_{\mathbf{x}} \rangle_{\mathcal{H}}] = \text{tr}(\Sigma \Sigma_\lambda^{-1}), \quad \mathcal{F}_\infty(\lambda) := \sup_{\mathbf{x} \in \text{supp}(P_X)} \|\Sigma_\lambda^{-1/2} K_{\mathbf{x}}\|_{\mathcal{H}}^2.$$

We can see that the risk admits the following bias-variance decomposition

$$\|Sf_t - f^*\|_{L_2(P_X)}^2 \leq \underbrace{2(\|Sf_t - S\bar{f}_t\|_{L_2(P_X)}^2)}_{V(t), \text{ variance}} + \underbrace{\|\bar{f}_t - f^*\|_{L_2(P_X)}^2}_{B(t), \text{ bias}}.$$

We upper bound the bias and variance separately.

Bounding the bias term $B(t)$: Note that by the update rule (4.1), it holds that

$$\begin{aligned} Sf_t - f^* &= S\bar{f}_{t-1} - f^* - \eta S(\Sigma + \lambda I)^{-1}(\Sigma \bar{f}_{t-1} - S^* f^*) \\ \Leftrightarrow S\bar{f}_t - f^* &= (I - \eta S(\Sigma + \lambda I)^{-1} S^*)(S\bar{f}_{t-1} - f^*). \end{aligned}$$

Therefore, unrolling the recursion gives $S\bar{f}_t - f^* = (I - \eta S(\Sigma + \lambda I)^{-1} S^*)^t(S\bar{f}_0 - f^*) = (I - \eta S(\Sigma + \lambda I)^{-1} S^*)^t(-f^*) = -(I - \eta S(\Sigma + \lambda I)^{-1} S^*)^t L^r h^*$. Write the spectral decomposition of L as $L = \sum_{j=1}^\infty \sigma_j \phi_j \phi_j^*$ for $\phi_j \in L_2(P_X)$ for $\sigma_j \geq 0$. We have $\|(I - \eta S(\Sigma + \lambda I)^{-1} S^*)^t L^r h^*\|_{L_2(P_X)} = \sum_{j=1}^\infty (1 - \eta \frac{\sigma_j}{\sigma_j + \lambda})^{2t} \sigma_j^{2r} h_j^2$, where $h = \sum_{j=1}^\infty h_j \phi_j$. We then apply Lemma 11 to obtain

$$B(t) \leq \exp(-\eta t) \sum_{j: \sigma_j \geq \lambda} h_j^2 + \left(\frac{2r}{e} \frac{\lambda}{\eta t}\right)^{2r} \sum_{j: \sigma_j < \lambda} h_j^2 \leq C \left[\exp(-\eta t) \vee \left(\frac{\lambda}{\eta t}\right)^{2r} \right] \|h^*\|_{L_2(P_X)}^2,$$

where C is a constant depending only on r .

Bounding the variance term $V(t)$: We now handle the variance term $V(t)$. For notational convenience, we write $A_\lambda := A + \lambda I$ for a linear operator A from a Hilbert space H to H . By the definition of f_t , we know

$$\begin{aligned} f_t &= (I - \eta(\Sigma + \lambda I)^{-1} \hat{\Sigma}) f_{t-1} + \eta(\Sigma + \lambda I)^{-1} \hat{S}^* Y \\ &= \sum_{j=0}^{t-1} (I - \eta(\Sigma + \lambda I)^{-1} \hat{\Sigma})^j \eta(\Sigma + \lambda I)^{-1} \hat{S}^* Y \\ &= \Sigma_\lambda^{-1/2} \eta \left[\sum_{j=0}^{t-1} (I - \eta \Sigma_\lambda^{-1/2} \hat{\Sigma} \Sigma_\lambda^{-1/2})^j \right] \Sigma_\lambda^{-1/2} \hat{S}^* Y =: \Sigma_\lambda^{-1/2} G_t \Sigma_\lambda^{-1/2} \hat{S}^* Y, \end{aligned}$$

where we defined $G_t := \eta \left[\sum_{j=0}^{t-1} (I - \eta \Sigma_\lambda^{-1/2} \hat{\Sigma} \Sigma_\lambda^{-1/2})^j \right]$. Accordingly, we decompose $V(t)$ as

$$\begin{aligned} \|Sf_t - S\bar{f}_t\|_{L_2(P_X)}^2 &\leq 2 \underbrace{\|S(f_t - \Sigma_\lambda^{-1/2} G_t \Sigma_\lambda^{-1/2} \hat{\Sigma} \bar{f}_t)\|_{L_2(P_X)}^2}_{(a)} \\ &\quad + \underbrace{\|S(\Sigma_\lambda^{-1/2} G_t \Sigma_\lambda^{-1/2} \hat{\Sigma} \bar{f}_t - \bar{f}_t)\|_{L_2(P_X)}^2}_{(b)}. \end{aligned}$$

We bound (a) and (b) separately.

Step 1. Bounding (a). Decompose (a) as

$$\|S(f_t - \Sigma_\lambda^{-1/2} G_t \Sigma_\lambda^{-1/2} \hat{\Sigma} \bar{f}_t)\|_{L_2(P_X)}^2 = \|S \Sigma_\lambda^{-1/2} G_t \Sigma_\lambda^{-1/2} (\hat{S}^* Y - \hat{\Sigma} \bar{f}_t)\|_{L_2(P_X)}^2$$

$$\leq \|S\Sigma_\lambda^{-1/2}\|^2 \|G_t\Sigma_\lambda^{-1/2}\hat{\Sigma}_\lambda\Sigma_\lambda^{-1/2}\|^2 \|\Sigma_\lambda^{1/2}\hat{\Sigma}_\lambda^{-1}\Sigma_\lambda^{1/2}\|^2 \|\Sigma_\lambda^{-1/2}(\hat{S}^*Y - \hat{\Sigma}\bar{f}_t)\|_{\mathcal{H}}^2.$$

We bound the terms in the RHS individually.

(i) $\|S\Sigma_\lambda^{-1/2}\|^2 = \|\Sigma_\lambda^{-1/2}\Sigma\Sigma_\lambda^{-1/2}\| \leq 1.$

(ii) Note that $\Sigma_\lambda^{-1/2}\hat{\Sigma}_\lambda\Sigma_\lambda^{-1/2} = I - \Sigma_\lambda^{-1/2}(\Sigma - \hat{\Sigma})\Sigma_\lambda^{-1/2} \succeq (1 - \|\Sigma_\lambda^{-1/2}(\Sigma - \hat{\Sigma})\Sigma_\lambda^{-1/2}\|)I.$

Proposition 6 of [Rudi & Rosasco \(2017\)](#) and its proof implies that for $\lambda \leq \|\Sigma\|$ and $0 < \delta < 1$, it holds that

$$\|\Sigma_\lambda^{-1/2}(\Sigma - \hat{\Sigma})\Sigma_\lambda^{-1/2}\| \leq \sqrt{\frac{2\beta\mathcal{F}_\infty(\lambda)}{n}} + \frac{2\beta(1 + \mathcal{F}_\infty(\lambda))}{3n} =: \Xi_n, \quad (\text{D.21})$$

with probability $1 - \delta$, where $\beta = \log\left(\frac{4\text{tr}(\Sigma\Sigma_\lambda^{-1})}{\delta}\right) = \log\left(\frac{4\mathcal{N}_\infty(\lambda)}{\delta}\right)$. By Lemma 14, $\beta \leq \log\left(\frac{4\text{tr}(\Sigma^{1/s})\lambda^{-1/s}}{\delta}\right)$ and $\mathcal{F}_\infty(\lambda) \leq \lambda^{-1}$. Therefore, if $\lambda = o(n^{-1}\log(n))$ and $\lambda = \Omega(n^{-1/s})$, the RHS can be smaller than $1/2$ for sufficiently large n , i.e. $\Xi_n = O(\sqrt{\log(n)/(n\lambda)}) \leq 1/2$. In this case we have,

$$\Sigma_\lambda^{-1/2}\hat{\Sigma}_\lambda\Sigma_\lambda^{-1/2} \succeq \frac{1}{2}I.$$

We denote this event as \mathcal{E}_1 .

(iii) Note that

$$\begin{aligned} G_t\Sigma_\lambda^{-1/2}\hat{\Sigma}_\lambda\Sigma_\lambda^{-1/2} &= \eta \left[\sum_{j=0}^{t-1} (I - \eta\Sigma_\lambda^{-1/2}\hat{\Sigma}_\lambda\Sigma_\lambda^{-1/2})^j \right] \Sigma_\lambda^{-1/2}\hat{\Sigma}_\lambda\Sigma_\lambda^{-1/2} \\ &= \eta \left[\sum_{j=0}^{t-1} (I - \eta\Sigma_\lambda^{-1/2}\hat{\Sigma}_\lambda\Sigma_\lambda^{-1/2})^j \right] (\Sigma_\lambda^{-1/2}\hat{\Sigma}_\lambda\Sigma_\lambda^{-1/2} + \lambda\Sigma_\lambda^{-1}). \end{aligned}$$

Thus, by Lemma 12 we have

$$\begin{aligned} &\|G_t\Sigma_\lambda^{-1/2}\hat{\Sigma}_\lambda\Sigma_\lambda^{-1/2}\| \\ &\leq \underbrace{\left\| \eta \left[\sum_{j=0}^{t-1} (I - \eta\Sigma_\lambda^{-1/2}\hat{\Sigma}_\lambda\Sigma_\lambda^{-1/2})^j \right] \Sigma_\lambda^{-1/2}\hat{\Sigma}_\lambda\Sigma_\lambda^{-1/2} \right\|}_{\leq 1 \text{ (due to Lemma 12)}} + \left\| \eta \left[\sum_{j=0}^{t-1} (I - \eta\Sigma_\lambda^{-1/2}\hat{\Sigma}_\lambda\Sigma_\lambda^{-1/2})^j \right] \lambda\Sigma_\lambda^{-1} \right\| \\ &\leq 1 + \eta \sum_{j=0}^{t-1} \|(I - \eta\Sigma_\lambda^{-1/2}\hat{\Sigma}_\lambda\Sigma_\lambda^{-1/2})^j\| \|\lambda\Sigma_\lambda^{-1}\| \leq 1 + \eta t. \end{aligned}$$

(iv) Note that

$$\|\Sigma_\lambda^{-1/2}(\hat{S}^*Y - \hat{\Sigma}\bar{f}_t)\|_{\mathcal{H}}^2 \leq 2(\|\Sigma_\lambda^{-1/2}[(\hat{S}^*Y - \hat{\Sigma}\bar{f}_t) - (S^*f^* - \Sigma\bar{f}_t)]\|_{\mathcal{H}}^2 + \|\Sigma_\lambda^{-1/2}(S^*f^* - \Sigma\bar{f}_t)\|_{\mathcal{H}}^2).$$

First we bound the first term of the RHS. Let $\xi_i = \Sigma_\lambda^{-1/2}[K_{\mathbf{x}_i}y_i - K_{\mathbf{x}_i}\bar{f}_t(\mathbf{x}_i) - (S^*f^* - \Sigma\bar{f}_t)]$. Then, $\{\xi_i\}_{i=1}^n$ is an i.i.d. sequence of zero-centered random variables taking value in \mathcal{H} and thus we have

$$\|\Sigma_\lambda^{-1/2}[(\hat{S}^*Y - \hat{\Sigma}\bar{f}_t) - (S^*f^* - \Sigma\bar{f}_t)]\|_{\mathcal{H}}^2 = \left\| \frac{1}{n} \sum_{i=1}^n \xi_i \right\|_{\mathcal{H}}^2.$$

The RHS can be bounded by using Bernstein's inequality in Hilbert space ([Caponnetto & De Vito, 2007](#)). To apply the inequality, we need to bound the variance and sup-norm of the random variable. The variance can be bounded as

$$\mathbb{E}[\|\xi_i\|_{\mathcal{H}}^2] \leq \mathbb{E}_{(\mathbf{x}, y)} \left[\|\Sigma_\lambda^{-1/2}(K_{\mathbf{x}}(f^*(\mathbf{x}) - \bar{f}_t(\mathbf{x})) + K_{\xi}\epsilon)\|_{\mathcal{H}}^2 \right]$$

$$\begin{aligned}
&\leq 2 \left\{ \mathbb{E}_{(\mathbf{x}, y)} \left[\|\Sigma_\lambda^{-1/2}(K_{\mathbf{x}}(f^*(\mathbf{x}) - \bar{f}_t(x))\|_{\mathcal{H}}^2 + \|\Sigma_\lambda^{-1/2}(K_{\mathbf{x}}\epsilon)\|_{\mathcal{H}}^2 \right] \right\} \\
&\leq 2 \left\{ \sup_{\mathbf{x} \in \text{supp}(P_X)} \|\Sigma_\lambda^{-1/2} K_{\mathbf{x}}\|^2 \|f^* - S\bar{f}_t\|_{L_2(P_X)}^2 + \sigma^2 \text{tr}(\Sigma_\lambda^{-1} \Sigma) \right\} \\
&\leq 2 \{ \mathcal{F}_\infty(\lambda) B(t) + \sigma^2 \text{tr}(\Sigma_\lambda^{-1} \Sigma) \} \\
&\leq 2 \{ \lambda^{-1} B(t) + \sigma^2 \text{tr}(\Sigma_\lambda^{-1} \Sigma) \},
\end{aligned}$$

The sup-norm can be bounded as follows. Observe that $\|\bar{f}_t\|_\infty \leq \|\bar{f}_t\|_{\mathcal{H}}$, and thus by Lemma 13,

$$\begin{aligned}
\|\xi_i\|_{\mathcal{H}} &\leq 2 \sup_{\mathbf{x} \in \text{supp}(P_X)} \|\Sigma_\lambda^{-1/2} K_{\mathbf{x}}\|_{\mathcal{H}} (\sigma + \|f^*\|_\infty + \|\bar{f}_t\|_\infty) \\
&\lesssim \mathcal{F}_\infty^{1/2}(\lambda) (\sigma + M + (1 + t\eta) \lambda^{-(1/2-r)_+}) \\
&\lesssim \lambda^{-1/2} (\sigma + M + (1 + t\eta) \lambda^{-(1/2-r)_+}).
\end{aligned}$$

Therefore, for $0 < \delta < 1$, Bernstein's inequality (see Proposition 2 of [Caponnetto & De Vito \(2007\)](#)) yields that

$$\left\| \frac{1}{n} \sum_{i=1}^n \xi_i \right\|_{\mathcal{H}}^2 \leq C \left(\sqrt{\frac{\lambda^{-1} B(t) + \sigma^2 \text{tr}(\Sigma_\lambda^{-1} \Sigma)}{n}} + \frac{\lambda^{-1/2} (\sigma + M + (1 + t\eta) \lambda^{-(1/2-r)_+})}{n} \right)^2 \log(1/\delta)^2$$

with probability $1 - \delta$ where C is a universal constant. We define this event as \mathcal{E}_2 .

For the second term $\|\Sigma_\lambda^{-1/2}(S^* f^* - \Sigma \bar{f}_t)\|_{\mathcal{H}}^2$ we have

$$\|\Sigma_\lambda^{-1/2}(S^* f^* - \Sigma \bar{f}_t)\|_{\mathcal{H}}^2 \leq \|\Sigma_\lambda^{-1/2}(f^* - S\bar{f}_t)\|_{\mathcal{H}}^2 = \|f^* - S\bar{f}_t\|_{L_2(P_X)}^2 \leq B(t).$$

Combining these evaluations, on the event \mathcal{E}_2 where $P(\mathcal{E}_2) \geq 1 - \delta$ for $0 < \delta < 1$ we have

$$\begin{aligned}
&\|\Sigma_\lambda^{-1/2}(\hat{S}^* Y - \hat{\Sigma} \bar{f}_t)\|_{\mathcal{H}}^2 \\
&\stackrel{(i)}{\leq} C \left(\sqrt{\frac{\lambda^{-1} B(t) + \sigma^2 \text{tr}(\Sigma_\lambda^{-1} \Sigma)}{n}} + \frac{\lambda^{-1/2} (\sigma + M + (1 + t\eta) \lambda^{-(1/2-r)_+})}{n} \right)^2 \log(1/\delta)^2 + B(t).
\end{aligned}$$

where we used Lemma 14 in (i).

Step 2. Bounding (b). On the event \mathcal{E}_1 , the term (b) can be evaluated as

$$\begin{aligned}
&\|S(\Sigma_\lambda^{-1/2} G_t \Sigma_\lambda^{-1/2} \hat{\Sigma} \bar{f}_t - \bar{f}_t)\|_{L_2(P_X)}^2 \\
&\leq \|\Sigma^{1/2}(\Sigma_\lambda^{-1/2} G_t \Sigma_\lambda^{-1/2} \hat{\Sigma} \bar{f}_t - \bar{f}_t)\|_{\mathcal{H}}^2 \\
&\leq \|\Sigma^{1/2} \Sigma_\lambda^{-1/2} (G_t \Sigma_\lambda^{-1/2} \hat{\Sigma} \Sigma_\lambda^{-1/2} - I) \Sigma_\lambda^{1/2} \bar{f}_t\|_{\mathcal{H}}^2 \\
&\leq \|\Sigma^{1/2} \Sigma_\lambda^{-1/2} \| (G_t \Sigma_\lambda^{-1/2} \hat{\Sigma} \Sigma_\lambda^{-1/2} - I) \Sigma_\lambda^{1/2} \bar{f}_t\|_{\mathcal{H}}^2 \\
&\leq \|(G_t \Sigma_\lambda^{-1/2} \hat{\Sigma} \Sigma_\lambda^{-1/2} - I) \Sigma_\lambda^{1/2} \bar{f}_t\|_{\mathcal{H}}^2.
\end{aligned} \tag{D.22}$$

where we used Lemma 13 in the last inequality. The term $\|(G_t \Sigma_\lambda^{-1/2} \hat{\Sigma} \Sigma_\lambda^{-1/2} - I) \Sigma_\lambda^{1/2} \bar{f}_t\|_{\mathcal{H}}$ can be bounded as follows. First, note that

$$\begin{aligned}
(G_t \Sigma_\lambda^{-1/2} \hat{\Sigma} \Sigma_\lambda^{-1/2} - I) \Sigma_\lambda^{1/2} &= \left\{ \eta \left[\sum_{j=0}^{t-1} (I - \eta \Sigma_\lambda^{-1/2} \hat{\Sigma} \Sigma_\lambda^{-1/2})^j \right] \Sigma_\lambda^{-1/2} \hat{\Sigma} \Sigma_\lambda^{-1/2} - I \right\} \Sigma_\lambda^{1/2} \\
&= (I - \eta \Sigma_\lambda^{-1/2} \hat{\Sigma} \Sigma_\lambda^{-1/2})^t \Sigma_\lambda^{1/2}.
\end{aligned}$$

Therefore, the RHS of (D.22) can be further bounded by

$$\begin{aligned}
&\|(I - \eta \Sigma_\lambda^{-1/2} \hat{\Sigma} \Sigma_\lambda^{-1/2})^t \Sigma_\lambda^{1/2} \bar{f}_t\|_{\mathcal{H}} \\
&= \|(I - \eta \Sigma_\lambda^{-1/2} \Sigma \Sigma_\lambda^{-1/2} + \eta \Sigma_\lambda^{-1/2} (\Sigma - \hat{\Sigma}) \Sigma_\lambda^{-1/2})^t \Sigma_\lambda^{1/2} \bar{f}_t\|_{\mathcal{H}}
\end{aligned}$$

$$\begin{aligned}
&= \left\| \sum_{k=0}^{t-1} (I - \eta \Sigma_\lambda^{-1/2} \hat{\Sigma} \Sigma_\lambda^{-1/2})^k (\eta \Sigma_\lambda^{-1/2} (\Sigma - \hat{\Sigma}) \Sigma_\lambda^{-1/2}) (I - \eta \Sigma_\lambda^{-1} \Sigma)^{t-k-1} \Sigma_\lambda^{1/2} \bar{f}_t - (I - \eta \Sigma_\lambda^{-1} \Sigma)^t \Sigma_\lambda^{1/2} \bar{f}_t \right\|_{\mathcal{H}} \\
&\stackrel{(i)}{\leq} \left\| (I - \eta \Sigma_\lambda^{-1} \Sigma)^t \Sigma_\lambda^{1/2} \bar{f}_t \right\|_{\mathcal{H}} \\
&\quad + \eta \sum_{k=0}^{t-1} \left\| (I - \eta \Sigma_\lambda^{-1/2} \hat{\Sigma} \Sigma_\lambda^{-1/2})^k \Sigma_\lambda^{-1/2} (\Sigma - \hat{\Sigma}) \Sigma_\lambda^{-1/2+r} (I - \eta \Sigma_\lambda^{-1} \Sigma)^{t-k-1} \Sigma_\lambda^{1/2-r} \bar{f}_t \right\|_{\mathcal{H}} \\
&\leq \left\| (I - \eta \Sigma_\lambda^{-1} \Sigma)^t \Sigma_\lambda^{1/2} \bar{f}_t \right\|_{\mathcal{H}} + t\eta \left\| \Sigma_\lambda^{-1/2} (\Sigma - \hat{\Sigma}) \Sigma_\lambda^{-1/2+r} \right\| \left\| \Sigma_\lambda^{1/2-r} \bar{f}_t \right\|_{\mathcal{H}} \\
&= \left\| (I - \eta \Sigma_\lambda^{-1} \Sigma)^t \Sigma_\lambda^r \right\| \left\| \Sigma_\lambda^{1/2-r} \bar{f}_t \right\|_{\mathcal{H}} + t\eta \left\| \Sigma_\lambda^{-1/2} (\Sigma - \hat{\Sigma}) \Sigma_\lambda^{-1/2+r} \right\| \left\| \Sigma_\lambda^{1/2-r} \bar{f}_t \right\|_{\mathcal{H}} \\
&\lesssim \left\| (I - \eta \Sigma_\lambda^{-1} \Sigma)^t \Sigma_\lambda^r \right\| + t\eta \left\| \Sigma_\lambda^{-1/2} (\Sigma - \hat{\Sigma}) \Sigma_\lambda^{-1/2+r} \right\| (1 + t\eta) \|h^*\|_{L_2(P_X)}, \tag{D.23}
\end{aligned}$$

where (i) is due to exchangeability of Σ_λ and Σ . By Lemma 11, for the RHS we have

$$\left\| (I - \eta \Sigma_\lambda^{-1} \Sigma)^t \Sigma_\lambda^r \right\| \leq \exp(-\eta t/2) \vee \left(\frac{1}{e} \frac{\lambda}{\eta t} \right)^r.$$

Next, as in the (D.21), by applying the Bernstein inequality for asymmetric operators (Corollary 3.1 of Minsker (2017) with the argument in its Section 3.2), it holds that

$$\begin{aligned}
&\left\| \Sigma_\lambda^{-1/2} (\Sigma - \hat{\Sigma}) \Sigma_\lambda^{-1/2+r} \right\| \\
&\leq C' \left(\sqrt{\frac{\beta' C_\mu^2 (2^{2r-\mu} \vee \lambda^{2r-\mu}) \mathcal{N}_\infty(\lambda)}{n}} + \frac{\beta' ((1+\lambda)^r + C_\mu^2 \lambda^{-\mu/2} (2^{2r-\mu} \vee \lambda^{r-\mu/2}))}{n} \right) =: \Xi'_n,
\end{aligned}$$

with probability $1 - \delta$, where C' is a universal constant and $\beta' \leq \log \left(\frac{28 C_\mu^2 (2^{2r-\mu} \vee \lambda^{-\mu+2r}) \text{tr}(\Sigma^{1/s}) \lambda^{-1/s}}{\delta} \right)$. We also used the following bounds on the sup-norm and the second order moments:

$$\begin{aligned}
&\text{(sup-norm)} \quad \left\| \Sigma_\lambda^{-1/2} (K_{\mathbf{x}} K_{\mathbf{x}}^* - \Sigma) \Sigma_\lambda^{-1/2+r} \right\| \\
&\quad \leq \left\| \Sigma_\lambda^{-1/2} K_{\mathbf{x}} K_{\mathbf{x}}^* \Sigma_\lambda^{-1/2+r} \right\| + \left\| \Sigma_\lambda^r \right\| \\
&\quad \leq \left\| \Sigma_\lambda^{-\mu/2} \Sigma_\lambda^{\mu/2-1/2} K_{\mathbf{x}} K_{\mathbf{x}}^* \Sigma_\lambda^{-1/2+\mu/2} \Sigma_\lambda^{r-\mu/2} \right\| + \left\| \Sigma_\lambda^r \right\| \\
&\quad \leq C_\mu^2 \lambda^{-\mu/2} (2^{r-\mu/2} \vee \lambda^{r-\mu/2}) + (1+\lambda)^r \quad (\text{a.s.}), \\
&\text{(2nd order moment 1)} \quad \left\| \mathbb{E}_{\mathbf{x}} [\Sigma_\lambda^{-1/2} (K_{\mathbf{x}} K_{\mathbf{x}}^* - \Sigma) \Sigma_\lambda^{-1+2r} (K_{\mathbf{x}} K_{\mathbf{x}}^* - \Sigma) \Sigma_\lambda^{-1/2}] \right\| \\
&\quad \leq \left\| \Sigma_\lambda^{-1/2} \Sigma \Sigma_\lambda^{-1/2} \right\| \sup_{\mathbf{x} \in \text{supp}(P_X)} [K_{\mathbf{x}}^* \Sigma_\lambda^{-1/2+\mu/2} \Sigma_\lambda^{-\mu+2r} \Sigma_\lambda^{-1/2+\mu/2} K_{\mathbf{x}}] \\
&\quad \leq C_\mu^2 (2^{2r-\mu} \vee \lambda^{2r-\mu}), \\
&\text{(2nd order moment 2)} \quad \left\| \mathbb{E}_{\mathbf{x}} [\Sigma_\lambda^{-1/2+r} (K_{\mathbf{x}} K_{\mathbf{x}}^* - \Sigma) \Sigma_\lambda^{-1/2} \Sigma_\lambda^{-1/2} (K_{\mathbf{x}} K_{\mathbf{x}}^* - \Sigma) \Sigma_\lambda^{-1/2+r}] \right\| \\
&\quad \leq \left\| \mathbb{E}_{\mathbf{x}} [\Sigma_\lambda^{-1/2+r} K_{\mathbf{x}} K_{\mathbf{x}}^* \Sigma_\lambda^{-1} K_{\mathbf{x}} K_{\mathbf{x}}^* \Sigma_\lambda^{-1/2+r}] \right\| \\
&\quad \leq C_\mu^2 (2^{2r-\mu} \vee \lambda^{2r-\mu}) \mathbb{E}_{\mathbf{x}} [K_{\mathbf{x}}^* \Sigma_\lambda^{-1} K_{\mathbf{x}}] \\
&\quad = C_\mu^2 (2^{2r-\mu} \vee \lambda^{2r-\mu}) \text{tr}(\Sigma \Sigma_\lambda^{-1}) \\
&\quad = C_\mu^2 (2^{2r-\mu} \vee \lambda^{2r-\mu}) \mathcal{N}_\infty(\lambda).
\end{aligned}$$

We define this event as \mathcal{E}_3 . Therefore, the RHS of (D.23) can be further bounded by

$$\begin{aligned}
&\left[\left\| (I - \eta \Sigma_\lambda^{-1} \Sigma)^t \Sigma_\lambda^r \right\| + C t \eta \left\| \Sigma_\lambda^{-1/2} (\Sigma - \hat{\Sigma}) \Sigma_\lambda^{-1/2+r} \right\| \right] (1 + t\eta) \|h^*\|_{L_2(P_X)} \\
&\leq \left[\exp(-\eta t/2) \vee \left(\frac{1}{e} \frac{\lambda}{\eta t} \right)^r + t \eta \Xi'_n \right] (1 + t\eta) \|h^*\|_{L_2(P_X)}.
\end{aligned}$$

Finally, note that when $\lambda = \lambda^*$ and $2r \geq \mu$,

$$\Xi_n'^2 = \tilde{O} \left(\frac{\lambda^{*2r-\mu-1/s}}{n} + \frac{\lambda^{*2(r-\mu)}}{n^2} \right) \leq \tilde{O} \left(n^{-\frac{s(4r-\mu)}{2rs+1}} + n^{-\frac{s(4r-2\mu)+2}{2rs+1}} \right) \leq \tilde{O} \left(n^{-\frac{2rs}{2rs+1}} \right).$$

Step 3. Combining the calculations in Step 1 and 2 leads to the desired result. \square

D.8.3 AUXILIARY LEMMAS

Lemma 11. For $t \in \mathbb{N}$, $0 < \eta < 1$, $0 < \sigma \leq 1$ and $0 \leq \lambda$, it holds that

$$\left(1 - \eta \frac{\sigma}{\sigma + \lambda}\right)^t \sigma^r \leq \begin{cases} \exp(-\eta t/2) & (\sigma \geq \lambda) \\ \left(\frac{2r}{e} \frac{\lambda}{\eta t}\right)^r & (\sigma < \lambda) \end{cases}.$$

Proof. When $\sigma \geq \lambda$, we have

$$\left(1 - \eta \frac{\sigma}{\sigma + \lambda}\right)^t \sigma^r \leq \left(1 - \eta \frac{\sigma}{2\sigma}\right)^t \sigma^r = (1 - \eta/2)^t \sigma^r \leq \exp(-t\eta/2) \sigma^r \leq \exp(-t\eta/2)$$

due to $\sigma \leq 1$. On the other hand, note that

$$\begin{aligned} \left(1 - \eta \frac{\sigma}{\sigma + \lambda}\right)^t \sigma^r &\leq \exp\left(-\eta t \frac{\sigma}{\sigma + \lambda}\right) \times \left(\frac{\sigma \eta t}{\sigma + \lambda}\right)^r \left(\frac{\sigma + \lambda}{\eta t}\right)^r \\ &\leq \sup_{x>0} \exp(-x) x^r \left(\frac{\sigma + \lambda}{\eta t}\right)^r \leq \left(\frac{(\sigma + \lambda)r}{\eta t e}\right)^r, \end{aligned}$$

where we used $\sup_{x>0} \exp(-x) x^r = (r/e)^r$. \square

Lemma 12. For $t \in \mathbb{N}$, $0 < \eta$ and $0 \leq \sigma$ such that $\eta\sigma < 1$, it holds that $\eta \sum_{j=0}^{t-1} (1 - \eta\sigma)^j \sigma \leq 1$.

Proof. If $\sigma = 0$, then the statement is obvious. Assume that $\sigma > 0$, then

$$\sum_{j=0}^{t-1} (1 - \eta\sigma)^j \sigma = \frac{1 - (1 - \eta\sigma)^t}{1 - (1 - \eta\sigma)} \sigma = \frac{1}{\eta} [1 - (1 - \eta\sigma)^t] \leq \eta^{-1}.$$

This yields the desired claim. \square

Lemma 13. Under (A5-7), for any $0 < \lambda < 1$ and $q \leq r$, it holds that

$$\|\Sigma_\lambda^{-s} \bar{f}_t\|_{\mathcal{H}} \lesssim (1 + \lambda^{-(1/2+(q-r))_+} + \lambda t \eta \lambda^{-(3/2+(q-r))_+}) \|h^*\|_{L_2(P_X)}.$$

Proof. Recall that

$$\bar{f}_t = (I - \eta(\Sigma + \lambda I)^{-1} \Sigma) \bar{f}_{t-1} + \eta(\Sigma + \lambda I)^{-1} S^* f^* = \sum_{j=0}^{t-1} (I - \eta(\Sigma + \lambda I)^{-1} \Sigma)^j \eta(\Sigma + \lambda I)^{-1} S^* f^*.$$

Therefore, we obtain the following

$$\begin{aligned} \|\Sigma_\lambda^{-q} \bar{f}_t\|_{\mathcal{H}} &= \eta \left\| \sum_{j=0}^{t-1} (I - \eta \Sigma_\lambda^{-1} \Sigma)^j \Sigma_\lambda^{-1-q} S^* L^r h^* \right\|_{\mathcal{H}} \\ &= \eta \left\| \sum_{j=0}^{t-1} (I - \eta \Sigma_\lambda^{-1} \Sigma)^j \Sigma_\lambda^{-1} (\Sigma + \lambda I) \Sigma_\lambda^{-q-1} S^* L^r h^* \right\|_{\mathcal{H}} \\ &\leq \eta \left\| \sum_{j=0}^{t-1} (I - \eta \Sigma_\lambda^{-1} \Sigma)^j \Sigma_\lambda^{-1} \Sigma \Sigma_\lambda^{-q-1} S^* L^r h^* \right\|_{\mathcal{H}} + \lambda \eta \left\| \sum_{j=0}^{t-1} (I - \eta \Sigma_\lambda^{-1} \Sigma)^j \Sigma_\lambda^{-1} \Sigma_\lambda^{-q-1} S^* L^r h^* \right\|_{\mathcal{H}} \\ &\leq \eta \left\| \sum_{j=0}^{t-1} (I - \eta \Sigma_\lambda^{-1} \Sigma)^j \Sigma_\lambda^{-1} \Sigma \right\| \left\| \Sigma_\lambda^{-q-1} S^* L^r h^* \right\|_{\mathcal{H}} + \lambda \eta \left\| \sum_{j=0}^{t-1} (I - \eta \Sigma_\lambda^{-1} \Sigma)^j \Sigma_\lambda^{-1} \Sigma_\lambda^{-q-1} S^* L^r h^* \right\|_{\mathcal{H}} \\ &\leq \left\| \Sigma_\lambda^{-q-1} S^* L^r h^* \right\|_{\mathcal{H}} + \lambda t \eta \left\| \Sigma_\lambda^{-1} \Sigma_\lambda^{-q-1} S^* L^r h^* \right\|_{\mathcal{H}} \end{aligned}$$

$$\begin{aligned}
&\leq \|S^* L_\lambda^{-q-1+r} h^*\|_{\mathcal{H}} + \lambda t \eta \|S^* L_\lambda^{-q-2+r} h^*\|_{\mathcal{H}} \\
&\leq \sqrt{\langle h^*, L_\lambda^{-q-1+r} S S^* L_\lambda^{-q-1+r} h^* \rangle_{L_2(P_X)}} + \lambda t \eta \sqrt{\langle h^*, L_\lambda^{-q-2+r} S S^* L_\lambda^{-q-2+r} h^* \rangle_{L_2(P_X)}} \\
&= \sqrt{\langle h^*, L_\lambda^{-q-1+r} L L_\lambda^{-q-1+r} h^* \rangle_{L_2(P_X)}} + \lambda t \eta \sqrt{\langle h^*, L_\lambda^{-q-2+r} L L_\lambda^{-q-2+r} h^* \rangle_{L_2(P_X)}} \\
&\leq (\lambda^{-1/2-(q-r)} + \lambda t \eta \lambda^{-3/2-(q-r)}) \|h^*\|_{L_2(P_X)} \leq (1 + t \eta) \lambda^{-1/2-(q-r)} \|h^*\|_{L_2(P_X)}.
\end{aligned}$$

□

Lemma 14. Under (A5-7) and for $\lambda \in (0, 1)$, it holds that $\mathcal{N}_\infty(\lambda) \leq \text{tr}(\Sigma^{1/s}) \lambda^{-1/s}$, and $\mathcal{F}_\infty(\lambda) \leq 1/\lambda$.

Proof. For the first inequality, we have

$$\begin{aligned}
\mathcal{N}_\infty(\lambda) &= \text{tr}(\Sigma \Sigma_\lambda^{-1}) = \text{tr}\left(\Sigma^{1/s} \Sigma^{1-1/s} \Sigma_\lambda^{-(1-1/s)} \Sigma_\lambda^{-1/s}\right) \\
&\leq \text{tr}\left(\Sigma^{1/s} \Sigma^{1-1/s} \Sigma_\lambda^{-(1-1/s)}\right) \lambda^{-1/s} \leq \text{tr}(\Sigma^{1/s}) \lambda^{-1/s}.
\end{aligned}$$

As for the second inequality, note that

$$\mathcal{F}_\infty(\lambda) = \sup_{\mathbf{x}} \langle K_{\mathbf{x}}, \Sigma_\lambda^{-1} K_{\mathbf{x}} \rangle_{\mathcal{H}} \leq \sup_{\mathbf{x}} \lambda^{-1} \langle K_{\mathbf{x}}, K_{\mathbf{x}} \rangle_{\mathcal{H}} \leq \lambda^{-1} \sup_{\mathbf{x}} k(\mathbf{x}, \mathbf{x}) \leq \lambda^{-1}.$$

□

D.9 PROOF OF PROPOSITION 8

Proof. Part (a) is a simple combination of Bai & Yin (2008, Theorem 2) and assumption (A3), which implies $\|\Sigma_X\|_2$ and $\|\Sigma_X^{-1}\|_2$ are both finite. For part (b), the substitution error for the variance term (ignoring the scalar σ^2) can be bounded as

$$\begin{aligned}
|V^* - \hat{V}| &= \left| \text{tr}\left(\mathbf{F}^{-1} \mathbf{X}^\top (\mathbf{X} \mathbf{F}^{-1} \mathbf{X}^\top)^{-2} \mathbf{X} \mathbf{F}^{-1} \Sigma_X\right) - \text{tr}\left(\hat{\mathbf{F}}^{-1} \mathbf{X}^\top (\mathbf{X} \hat{\mathbf{F}}^{-1} \mathbf{X}^\top)^{-2} \mathbf{X} \hat{\mathbf{F}}^{-1} \Sigma_X\right) \right| \\
&\stackrel{(i)}{\leq} O(1) \left\| \mathbf{F}^{-1} - \hat{\mathbf{F}}^{-1} \right\|_2 \left(\text{tr}\left(\mathbf{X}^\top \mathbf{S}^{-2} \mathbf{X} \mathbf{F} \Sigma_X\right) + \sqrt{d} \left\| \mathbf{X}^\top \mathbf{S}^{-2} \mathbf{X} \right\|_2 \left\| \Sigma_X \hat{\mathbf{F}}^{-1} \right\|_F \right) \\
&\quad + \text{tr}\left(\mathbf{X} \hat{\mathbf{F}}^{-1} \Sigma_X \mathbf{F}^{-1} \mathbf{X}^\top\right) \left\| \mathbf{S}^{-2} - \hat{\mathbf{S}}^{-2} \right\|_2 \stackrel{(ii)}{=} O(\epsilon).
\end{aligned}$$

where we defined $\mathbf{S} = \mathbf{X} \mathbf{F}^{-1} \mathbf{X}^\top$ and $\hat{\mathbf{S}} = \mathbf{X} \hat{\mathbf{F}}^{-1} \mathbf{X}^\top$ in (i) and applied $\text{tr}(\mathbf{A} \mathbf{B}) \leq \lambda_{\max}(\mathbf{A} + \mathbf{A}^\top) \text{tr}(\mathbf{B})$ for positive semi-definite \mathbf{B} , as well as $\text{tr}(\mathbf{A} \mathbf{B}) \leq \sqrt{d} \|\mathbf{A}\|_2 \|\mathbf{B}\|_F$, and (ii) is due to (A3), $\psi > 1$, Wedin (1973, Theorem 4.1) and the following estimate,

$$\begin{aligned}
n_u^2 \left\| \mathbf{S}^{-2} - \hat{\mathbf{S}}^{-2} \right\|_2 &\leq \left\| n_u \mathbf{S}^{-1} - n_u \hat{\mathbf{S}}^{-1} \right\|_2 \left(n_u \left\| \mathbf{S}^{-1} \right\|_2 + n_u \left\| \hat{\mathbf{S}}^{-1} \right\|_2 \right) \\
&\stackrel{(i)}{=} O(1) \left\| n_u \mathbf{S}^{-1} \right\|_2 \left\| n_u \hat{\mathbf{S}}^{-1} \right\|_2 \left\| \mathbf{S}/n_u - \hat{\mathbf{S}}/n_u \right\|_2 \stackrel{(ii)}{=} O(\epsilon),
\end{aligned}$$

where (i) again follows from Wedin (1973, Theorem 4.1), and (ii) is due to (A1)(A3) and $\psi > 1$ (which implies $\|n_u \mathbf{S}^{-1}\|_2$ and $\|n_u \hat{\mathbf{S}}^{-1}\|_2$ are bounded a.s.). Finally, from part (a) we know that $\psi = \Theta(\epsilon^{-2})$ suffices to achieve ϵ -accurate approximation of \mathbf{F} in spectral norm. The substitution error for the bias term can be derived from similar calculation, the details of which we omit. □

D.10 PROOF OF PROPOSITION 9

Proof. Since $\Sigma_\theta = \Sigma_X^{-r}$, we can simplify the expressions by defining $v_x \triangleq h$ and thus $v_\theta = h^{-r}$. From Theorem 2 we have the following derivation of the GD bias under (A1)(A3),

$$B(\hat{\theta}_I) \rightarrow \frac{m'_1}{m_1^2} \mathbb{E} \frac{h \cdot h^{-r}}{(1 + h \cdot m_1)^2} = \frac{\mathbb{E} \frac{h^{1-r}}{(1 + h \cdot m_1)^2}}{1 - \gamma \mathbb{E} \frac{(h \cdot m_1)^2}{(1 + h \cdot m_1)^2}}, \quad (\text{D.24})$$

where $m_1 = \lim_{\lambda \rightarrow 0^+} m(-\lambda)$, and m satisfies

$$\frac{1}{m(-\lambda)} = \lambda + \gamma \mathbb{E} \left[\frac{h}{1 + h \cdot m(-\lambda)} \right].$$

Similarly, for NGD ($\mathbf{P} = \Sigma_X^{-1}$) we have

$$B(\hat{\boldsymbol{\theta}}_{\mathbf{F}^{-1}}) \rightarrow \frac{m'_2}{m_2^2} \mathbb{E} \frac{h \cdot h^{-r}}{(1 + m_2)^2} = \frac{\mathbb{E} \frac{h^{1-r}}{(1+m_2)^2}}{1 - \gamma \mathbb{E} \frac{m_2^2}{(1+m_2)^2}} = \frac{\mathbb{E} h^{1-r}}{(1 + m_2)^2 - \gamma m_2^2}, \quad (\text{D.25})$$

where standard calculation yields $m_2 = (\gamma - 1)^{-1}$, and thus $B(\hat{\boldsymbol{\theta}}_{\mathbf{F}^{-1}}) \rightarrow (1 - \gamma^{-1}) \mathbb{E} h^{1-r}$.

To compare the magnitude of (D.24) and (D.25), observe the following equivalence.

$$\begin{aligned} B(\hat{\boldsymbol{\theta}}_{\mathbf{I}}) &\leq B(\hat{\boldsymbol{\theta}}_{\mathbf{F}^{-1}}) \\ \Leftrightarrow \mathbb{E} \frac{h^{1-r}}{(1 + h \cdot m_1)^2} \cdot \frac{\gamma}{\gamma - 1} &\leq \left(1 - \gamma \mathbb{E} \frac{(h \cdot m_1)^2}{(1 + h \cdot m_1)^2} \right) \mathbb{E} h^{1-r}. \\ &\stackrel{(i)}{\Leftrightarrow} \mathbb{E} \frac{\zeta^{1-r}}{(1 + \zeta)^2} \mathbb{E} \frac{\zeta}{1 + \zeta} \leq \mathbb{E} \frac{\zeta}{(1 + \zeta)^2} \mathbb{E} \zeta^{1-r} \mathbb{E} \frac{1}{1 + \zeta}. \end{aligned} \quad (\text{D.26})$$

where (i) follows from the definition of m_1 and we defined $\zeta \triangleq h \cdot m_1$. Note that when $r \geq 1$ and h is not a point mass, we have

$$\mathbb{E} \frac{\zeta^{1-r}}{(1 + \zeta)^2} \mathbb{E} \frac{\zeta}{1 + \zeta} > \mathbb{E} \frac{\zeta}{(1 + \zeta)^2} \mathbb{E} \frac{\zeta^{1-r}}{1 + \zeta} > \mathbb{E} \frac{\zeta}{(1 + \zeta)^2} \mathbb{E} \zeta^{1-r} \mathbb{E} \frac{1}{1 + \zeta}.$$

On the other hand, when $r \leq 0$, following the exact same procedure we get

$$\mathbb{E} \frac{\zeta^{1-r}}{(1 + \zeta)^2} \mathbb{E} \frac{\zeta}{1 + \zeta} < \mathbb{E} \frac{\zeta}{(1 + \zeta)^2} \mathbb{E} \zeta^{1-r} \mathbb{E} \frac{1}{1 + \zeta}.$$

Combining the two cases completes the proof. \square

D.11 PROOF OF COROLLARY 10

Proof. Note that in this setting v_x takes value of $\frac{2}{1+\kappa}$ and $\frac{2\kappa}{1+\kappa}$ with probability 1/2 each. From (D.25) one can easily verify that for NGD,

$$B(\hat{\boldsymbol{\theta}}_{\mathbf{F}^{-1}}) \rightarrow \frac{2^{-r}(1 + \kappa^{1-r})}{(1 + \kappa)^{1-r}} \left(1 - \frac{1}{\gamma} \right).$$

For GD, the bias formula (D.24) can be simplified as

$$B(\hat{\boldsymbol{\theta}}_{\mathbf{I}}) \rightarrow \frac{1}{\gamma} \cdot \left(\frac{\left(\frac{2}{1+\kappa} \right)^{-r}}{(1 + \kappa + 2m_1)^2} + \frac{\kappa \left(\frac{2\kappa}{1+\kappa} \right)^{-r}}{(1 + \kappa + 2\kappa m_1)^2} \right) \cdot \left(\frac{m_1}{(1 + \kappa + 2m_1)^2} + \frac{\kappa m_1}{(1 + \kappa + 2\kappa m_1)^2} \right)^{-1}.$$

On the other hand, from standard numerical calculation one can show that when $\gamma > 1, \kappa \geq 1$,

$$m_1 = \frac{(\kappa + 1) \sqrt{\gamma^2(\kappa + 1)^2 + 4(1 - \gamma)(\kappa - 1)^2} + (2 - \gamma)(\kappa + 1)^2}{8(\gamma - 1)\kappa}.$$

Setting $B(\hat{\boldsymbol{\theta}}_{\mathbf{I}}) = B(\hat{\boldsymbol{\theta}}_{\mathbf{F}^{-1}})$ and solve for r , we have

$$r^* = \ln \left(\frac{c_4 - c_2}{c_1 - c_3} \right) / \ln \kappa, \quad (\text{D.27})$$

where

$$\begin{aligned} c_1 &= \left(1 - \frac{1}{\gamma}\right) \frac{1}{\kappa + 1}, \\ c_2 &= \left(1 - \frac{1}{\gamma}\right) \frac{\kappa}{\kappa + 1}, \\ c_3 &= \frac{1}{\gamma} \cdot \frac{(1 + \kappa + 2\kappa m_1)^2}{m_1(1 + \kappa + 2\kappa m_1)^2 + \kappa m_1(1 + \kappa + 2m_1)^2}, \\ c_4 &= \frac{1}{\gamma} \cdot \frac{\kappa(1 + \kappa + 2m_1)^2}{m_1(1 + \kappa + 2\kappa m_1)^2 + \kappa m_1(1 + \kappa + 2m_1)^2}. \end{aligned}$$

Hence, Proposition 9 (from which we know $r^* \in (0, 1)$) and the uniqueness of (D.27) implies that when $r \geq r^*$, $B(\hat{\theta}_I) \geq B(\hat{\theta}_{F-1})$, and vice versa. Finally, observe that in the special case of $\gamma = 2$, $m_1 = \frac{\kappa+1}{2\sqrt{\kappa}}$. Therefore, one can check that constants in (D.27) simplify to

$$c_1 - c_3 = \frac{1 - \sqrt{\kappa}}{2(\kappa + 1)}, \quad c_2 - c_4 = \frac{\sqrt{\kappa}(\sqrt{\kappa} - 1)}{2(\kappa + 1)},$$

which implies that $r^* = 1/2$. □

E EXPERIMENT SETUP

E.1 PROCESSING THE DATASETS

To obtain extra unlabeled data to estimate the Fisher, we zero pad pixels on the borders of each image before randomly cropping; a random horizontal flip is also applied for CIFAR10 images (Krizhevsky et al., 2009). We preprocess all images by dividing pixel values by 255 before centering them to be located within $[-0.5, 0.5]$ with the subtraction by $1/2$. For experiments on CIFAR10, we downsample the original images using a max pooling layer with kernel size 2 and stride 2.

E.2 SETUP AND IMPLEMENTATION FOR OPTIMIZERS

In all settings, GD uses a learning rate of 0.01 that is exponentially decayed every 1k updates with the parameter value 0.999. For NGD, we use a fixed learning rate of 0.03. Since inverting a parameter-by-parameter-sized Fisher estimate per iteration would be costly, we adopt the Hessian free approach (Martens, 2010) which computes approximate matrix-inverse-vector products using the conjugate gradient (CG) method (Nocedal & Wright, 2006; Boyd et al., 2004). For each approximate inversion, we run CG for 200 iterations starting from the solution returned by the previous CG run. The precise number of CG iterations and the initialization heuristic roughly follow Martens & Sutskever (2012). For the first run of CG, we initialize the vector from a standard Gaussian, and run CG for 5k iterations. To ensure invertibility, we apply a very small amount of damping (0.00001) in most scenarios. For geometric interpolation experiments between GD and NGD, we use the singular value decomposition to compute the minus α power of the Fisher, as CG is not applicable in this scenario.

E.3 OTHER DETAILS

For experiments in the label noise and misspecification sections, we pretrain the teacher using the Adam optimizer (Kingma & Ba, 2014) with its default hyperparameters and a learning rate of 0.001.

For experiments in the misalignment section, we downsample all images twice using max pooling with kernel size 2 and stride 2. Moreover, only for experiments in this section, we implement natural gradient descent by exactly computing the Fisher on a large batch of unlabeled data and inverting the matrix by calling PyTorch’s `torch.inverse` before right multiplying the gradient.