

Introduction: Preconditioned Gradient Descent

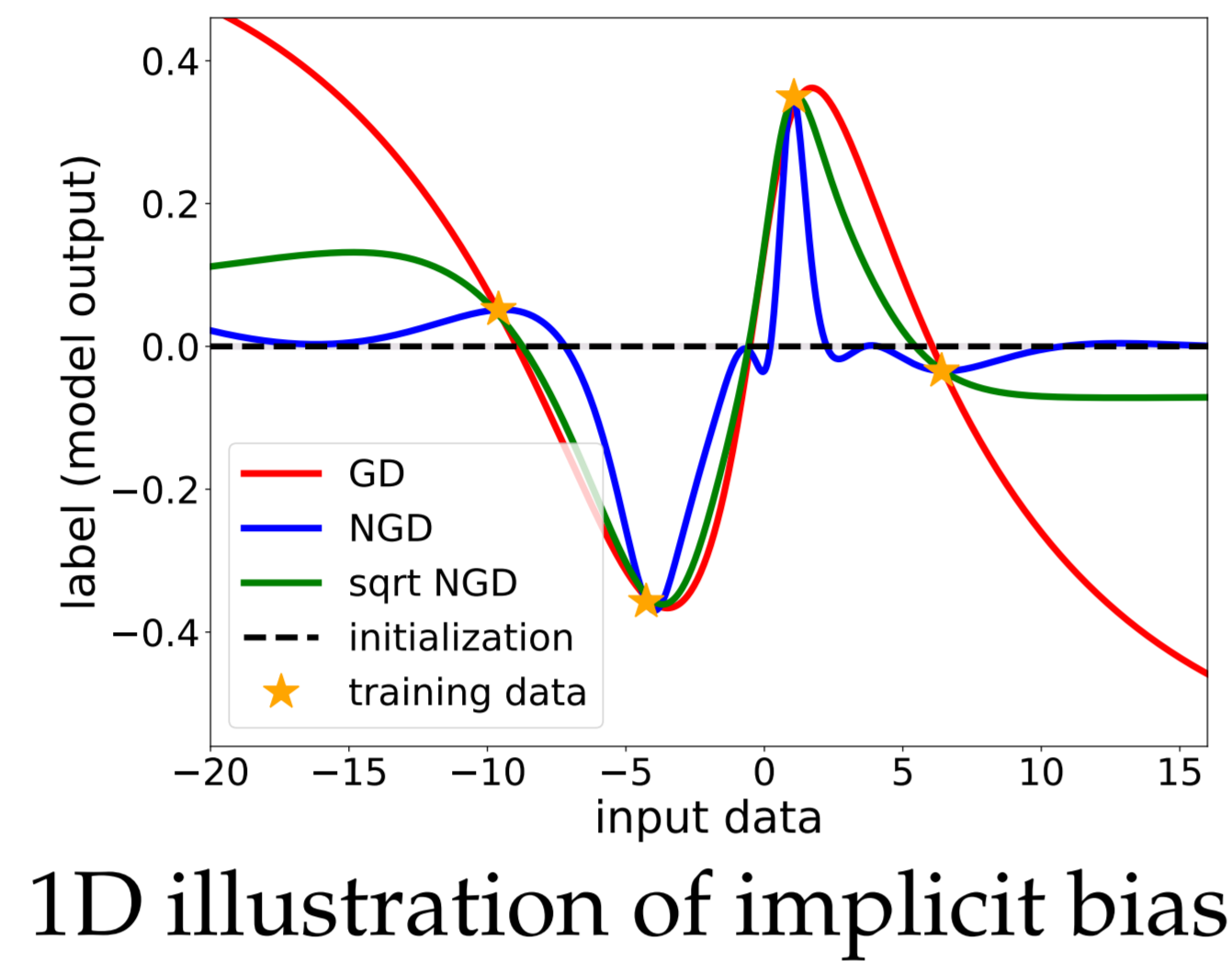
Update Rule: $\theta_{t+1} = \theta_t - \eta P(t) \nabla_{\theta_t} L(f_{\theta_t})$, $t = 0, 1, \dots$

Common choices of preconditioner P and corresponding algorithm:

- Inverse Fisher information matrix \Rightarrow *natural gradient descent* (NGD).
- Certain diagonal matrix \Rightarrow *adaptive gradient methods* (e.g. Adagrad, Adam).

Implicit Bias of Preconditioned Updates:

- Modern ML models (e.g. neural nets) are often **overparameterized**.
- Overparameterized models may **interpolate** training data *in different ways*.
- P alters properties of the interpolant.



Motivation of this work:

- How does *preconditioning affects generalization* under *interpolation*?
- Can we determine the *optimal preconditioner* for generalization?

Implicit Bias in Least Squares Regression

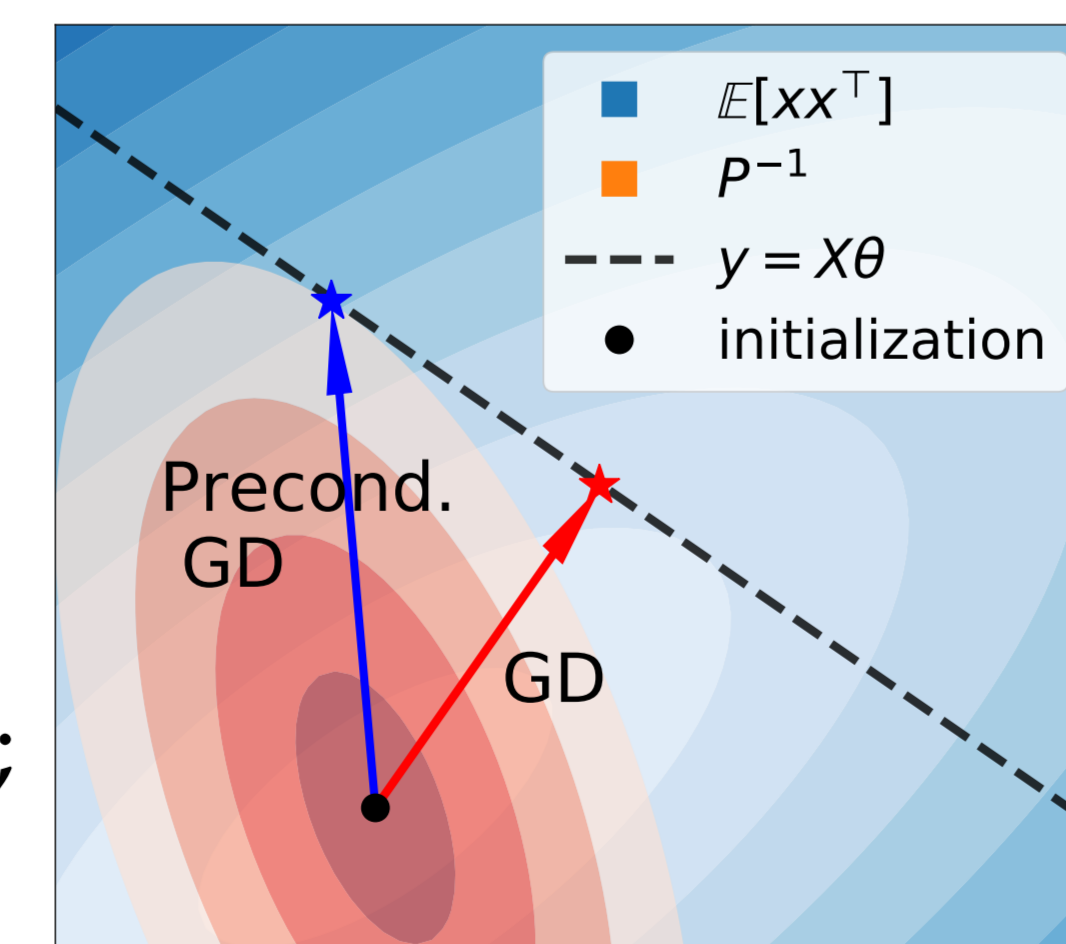
- **Student-teacher Setup.** $y_i = x_i^\top \theta_* + \varepsilon_i$, $1 \leq i \leq n$; $\mathbb{E}[xx^\top] = \Sigma_x \in \mathbb{R}^{d \times d}$.
- **Overparameterized Asymptotics.** $n, d \rightarrow \infty$, $d/n \rightarrow \gamma \in (1, \infty)$.
- **Update Rule.** Preconditioned gradient descent on *squared loss*:

$$d\theta(t) = P(t)X^\top(y - X\theta(t))dt, \quad \theta(0) = 0.$$

Stationary Solution ($t \rightarrow \infty$):

- **Gradient descent:** min ℓ_2 -norm interpolant.
- **Preconditioned GD:** for time-invariant and full-rank $P \Rightarrow$ minimum $\|\theta\|_{P^{-1}}$ interpolant.

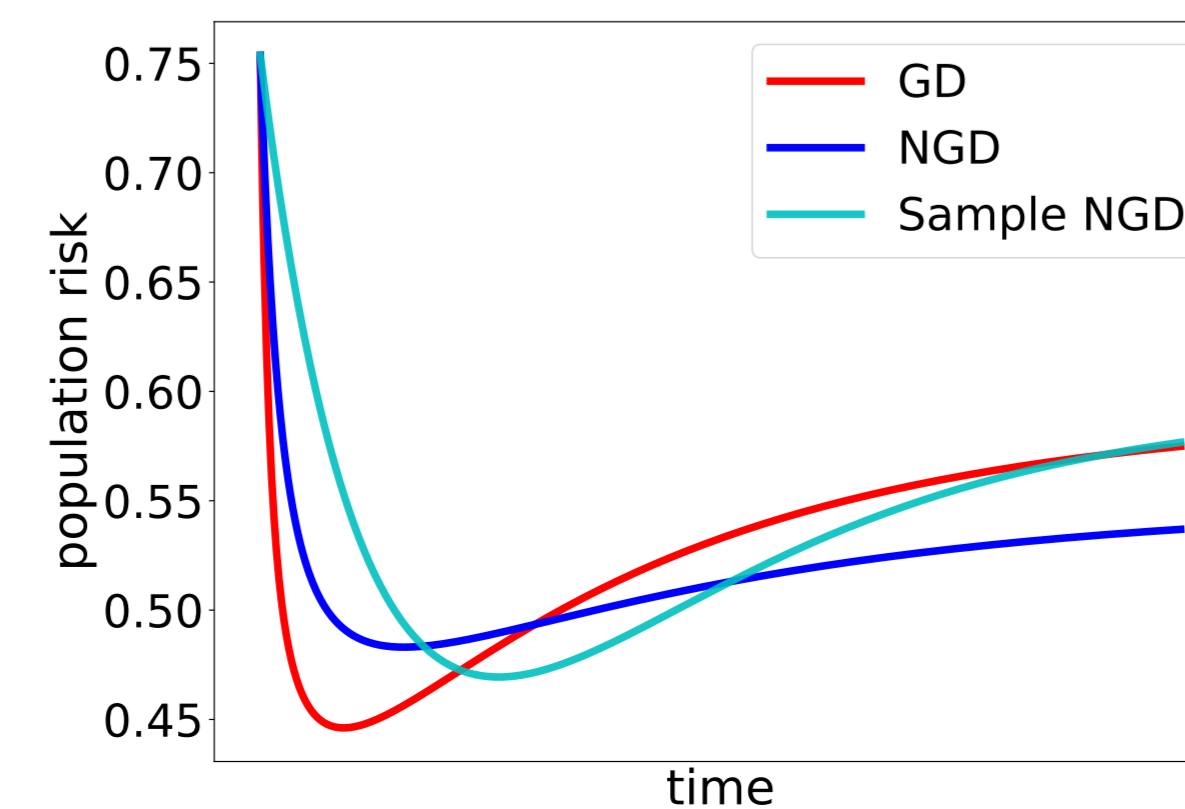
Common Argument: min ℓ_2 -norm solution generalizes well; therefore GD is better than preconditioned updates!



Question: Why is the ℓ_2 norm the best measure of generalization?

Noticeable Examples of Preconditioner:

- **Identity:** $P = I_d$ gives the min ℓ_2 norm interpolant (also true for momentum GD and SGD).
- **Population Fisher:** $P = F^{-1} = \Sigma_x^{-1}$ (NGD).
- **Variants of Sample Fisher:** $P = (X^\top X + \lambda I_d)^{-1}$ leads to the same solution as GD.



Bias-variance Decomposition of Generalization Error

Thm. (informal). Prediction risk of the min $\|\theta\|_{P^{-1}}$ solution is given as

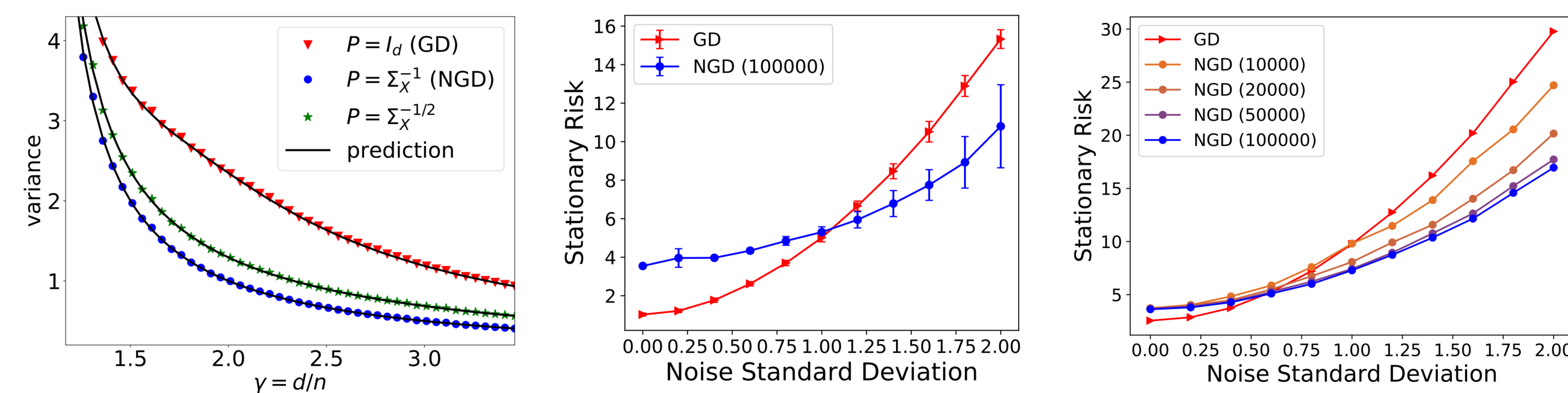
$$\mathbb{E}(\langle x, \theta_* \rangle - \langle x, \hat{\theta}_{P^{-1}} \rangle)^2 \xrightarrow{P} m'_0 m_0^{-2} \left(\underbrace{\gamma \mathbb{E}[v_x v_\theta (v_{xP} \cdot m_0 + 1)^{-2}]}_{\text{bias}} + \underbrace{\tilde{\sigma}^2}_{\text{variance}} \right),$$

where m_0 is the Stieltjes transform of $\frac{1}{n} X P X^\top$ evaluated at $\lambda \rightarrow 0_+$.

- **Bias term:** “Difficulty” in learning the *teacher model* θ_* .
- **Variance term:** “Stability” of learning under *label noise*.

Variance term: NGD is Optimal

Thm. (informal). NGD ($P = F^{-1}$) achieves lowest stationary variance.



(a) Linear regression. (b) 2-layer MLP (MNIST). (c) 2-layer MLP (CIFAR).

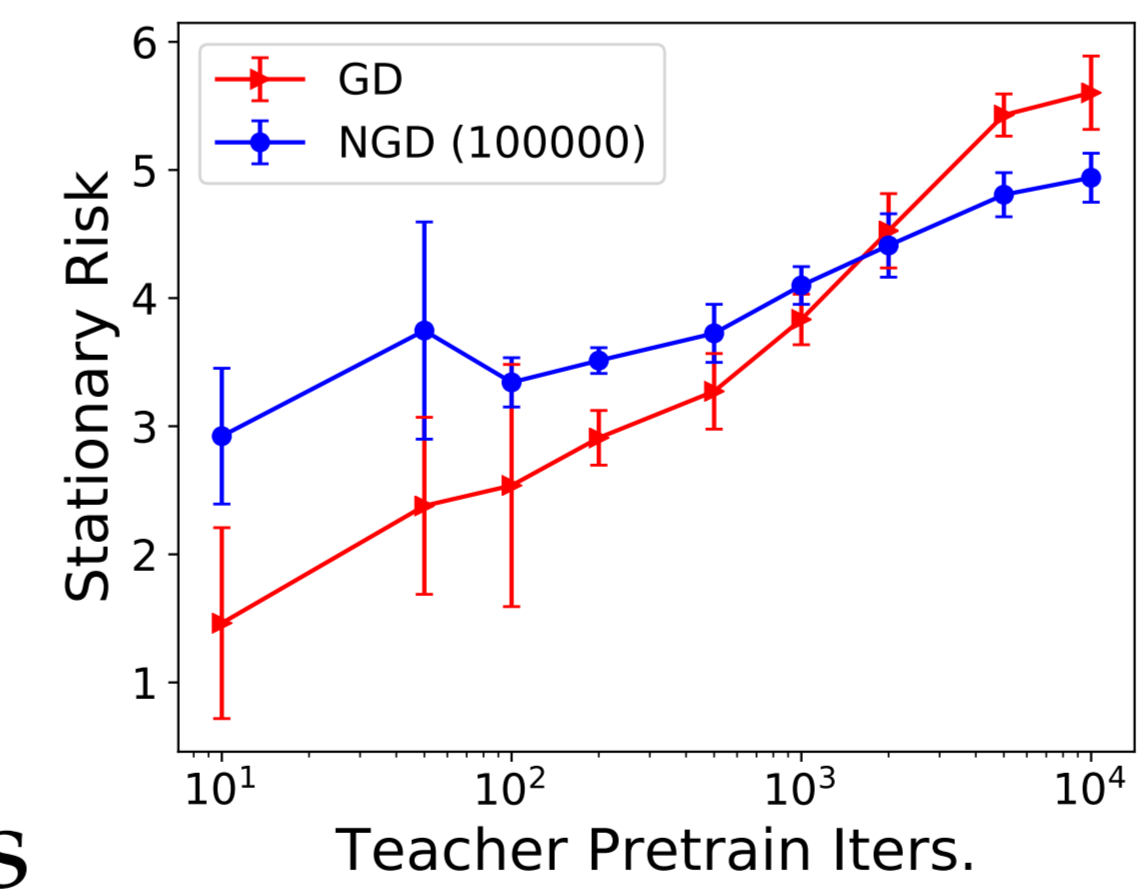
- (a)(b): labels are noisy (risk is *variance-dominated*) \Rightarrow NGD beneficial.
- (c): this advantage is present only for the **population Fisher**.

Misspecification \approx Label Noise:

Misspecified Model: $f_*(x) = x^\top \theta_* + f_*^c(x)$; residual f_*^c cannot be learned by student model.

Creating Misspecification in Neural Network:

- **Student:** small two-layer MLP.
- **Teacher:** ResNet-20 at varying training epochs



Bias term: No Free Lunch

General Prior: $\mathbb{E}[\theta_* \theta_*^\top] = d^{-1} \Sigma_\theta$, i.e. computing the *Bayes risk*.

Thm. (informal). Among all P codiagonalizable with Σ_x , bias is minimized by $P = U \text{diag}(U^\top \Sigma_\theta U) U^\top$, where U is the eigenvectors of Σ_x .

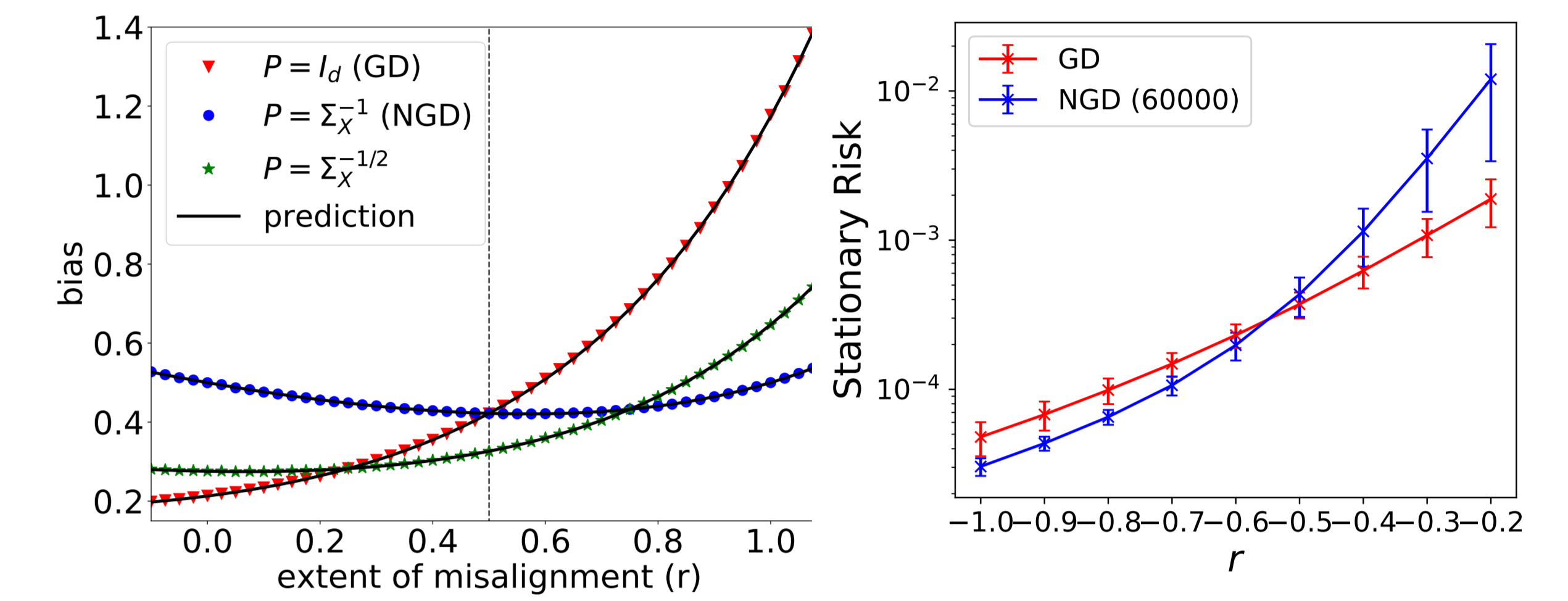
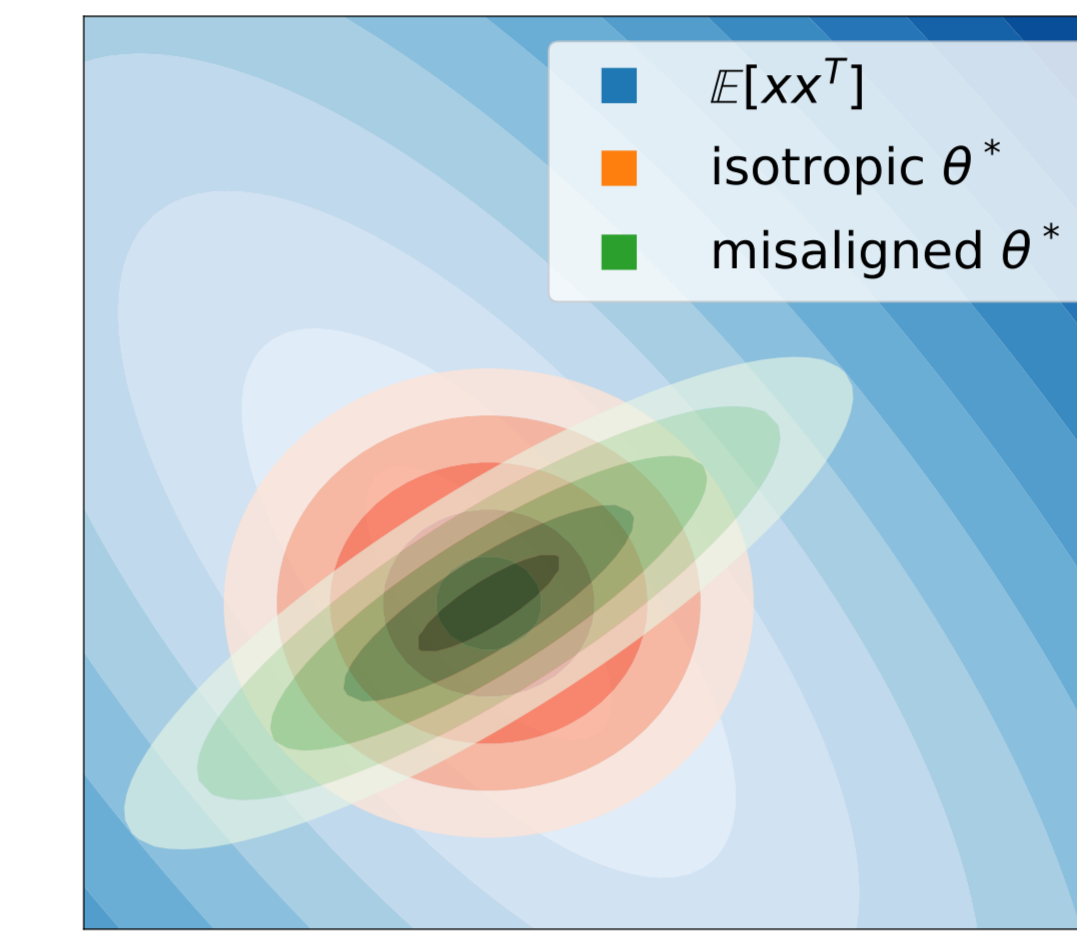
Remark: Setup extends previously assumed *isotropic prior* [Dobriban and Wager 18].

No-free-lunch: The optimal preconditioner depends on the “orientation” of teacher model θ_* , which is usually not known *a priori*.

Bias Term (continued): Alignment & Source Condition

- GD achieves optimal bias when teacher is **isotropic**: $\Sigma_\theta = I_d$.
- NGD is optimal under **misalignment**: $\Sigma_\theta = \Sigma_x^{-1}$ (“hard” problem).

Remark: We also show that this trend is roughly preserved under **early stopping**.



(a) Intuition of “alignment”. (b) Linear regression. (c) 2-layer MLP (MNIST).

Analogy to Source Condition ($\mathbb{E}\|\Sigma_x^{-r/2} \theta_*\|_2 < \infty$):

Prop. (informal). Consider $\Sigma_\theta = \Sigma_x^{-r}$. Then for some $r^* \in (0, 1)$, NGD achieves lower (higher) bias than GD if and only if $r > (<) r^*$.

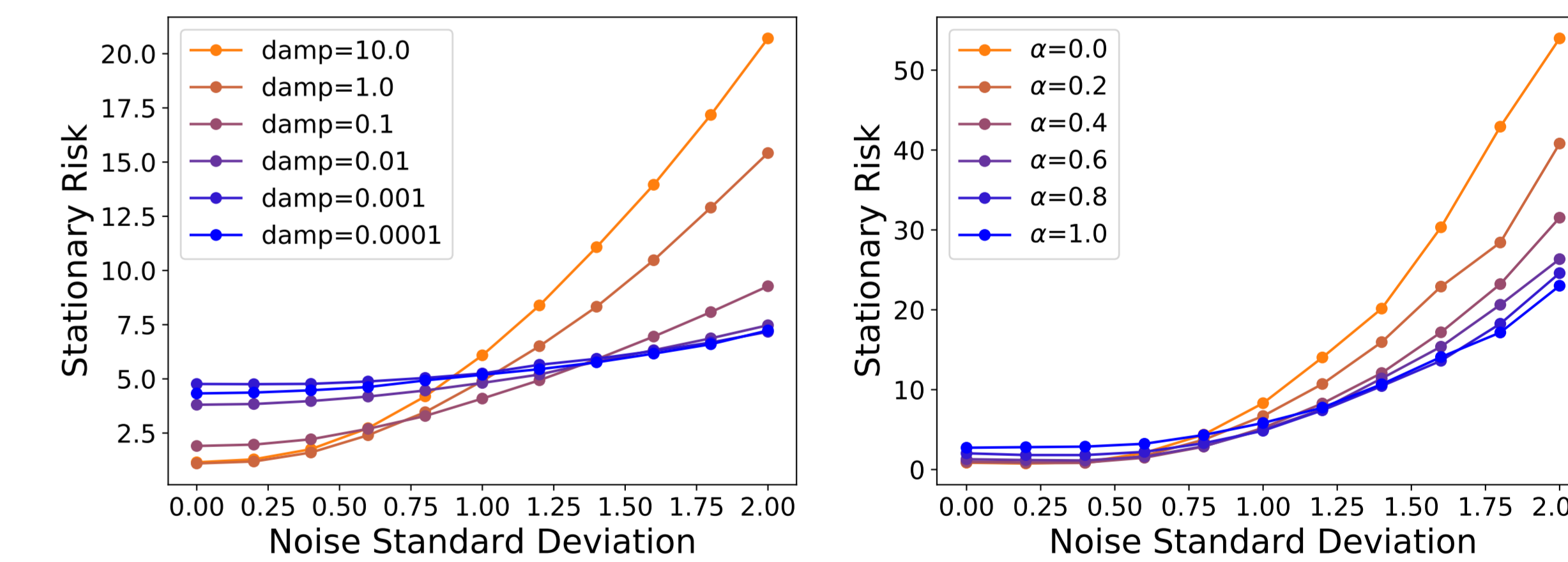
“Interpolating” Between GD and NGD

Question: Is it advantageous to “combine” GD and NGD?

Bias-variance Tradeoff:

- **Additive interp.:** $P_\alpha = (\Sigma_x + \alpha I_d)^{-1}$.
- **Geometric interp.:** $P_\alpha = \Sigma_x^{\alpha-1}$.

- Large $\alpha \Rightarrow$ GD-like update. Additive interp. (MLP).
- Small $\alpha \Rightarrow$ NGD-like update. Geometric interp. (MLP).



Message: At some SNR, *interpolating* between GD and NGD is beneficial.

Fast Decay in Population Risk:

Consider the following preconditioned update in the RKHS.

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

Remark: Update corresponds to *additive interpolation* between GD and NGD.

Thm. (informal). Preconditioned GD with *properly chosen* α achieves the minimax optimal rate $R(f_t) = \|S f_t - f^*\|_{L_2(P_X)}^2 = \tilde{O}\left(n^{-\frac{2rs}{2rs+1}}\right)$ in $t = \Theta(\log n)$ steps, whereas ordinary GD requires $t = \Theta\left(n^{\frac{2rs}{2rs+1}}\right)$ steps.

Message: Preconditioning can improve the *efficiency of learning*.