



High-dimensional Analysis of Knowledge Distillation: Weak-to-Strong Generalization and Scaling Laws

Arxiv 2024

M. Emrullah Ildiz, Halil Alperen Gozeten, Ege Onur Taga, Marco Mondelli, Samet Oymak

Problem Setup

Stage 1: Surrogate Model.

- Data distribution $(\tilde{\mathbf{x}}, \tilde{y}) \sim \mathcal{D}_s$
 - $\tilde{y} = \tilde{\mathbf{x}}^\top \boldsymbol{\beta}_\star + \tilde{z}$
 - $\boldsymbol{\beta}_\star \in \mathbb{R}^p$, p : Feature dimension
 - $\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_s)$, $\boldsymbol{\Sigma}_s$: Covariance matrix of distribution \mathcal{D}_s
 - $\tilde{z} \sim \mathcal{N}(0, \sigma_s^2)$: independent of $\tilde{\mathbf{x}}$
- Surrogate dataset $\{(\tilde{\mathbf{x}}_i, \tilde{y}_i)_{i=1}^m\}$: i.i.d. from \mathcal{D}_s
- Estimator $\boldsymbol{\beta}^s$: both under- and over-parameterized settings
 - $\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1^\top, \dots, \tilde{\mathbf{x}}_m^\top]^\top \in \mathbb{R}^{m \times p}$, $\tilde{\mathbf{y}} = [\tilde{y}_1, \dots, \tilde{y}_m]^\top \in \mathbb{R}^m$
 - Under-parametrized ($m \geq p$): Quadratic loss
 - Over-parametrized ($m < p$): Minimum norm interpolator

$$\boldsymbol{\beta}^s = \text{Est}(\tilde{\mathbf{X}}, \tilde{\mathbf{y}}) := \begin{cases} \arg \min_{\boldsymbol{\beta}} \|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}\|_2^2, & \text{if } m \geq p, \\ \arg \min_{\boldsymbol{\beta}} \{\|\boldsymbol{\beta}\|_2^2 : \tilde{\mathbf{X}}\boldsymbol{\beta} = \tilde{\mathbf{y}}\} & \text{if } m < p, \end{cases}$$

Problem Setup

Stage 2: Target Model.

- Data distribution $(\mathbf{x}, y^s) \sim \mathcal{D}_t(\boldsymbol{\beta}^s)$
 - $y^s = \mathbf{x}^\top \boldsymbol{\beta}^s + z$
 - $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_t)$, $\boldsymbol{\Sigma}_t$: Covariance matrix of distribution $\mathcal{D}_t(\boldsymbol{\beta}^s)$
 - $z \sim \mathcal{N}(0, \sigma_t^2)$: independent of \mathbf{x}
- Target Dataset $\{(\mathbf{x}_i, y_i^s)_{i=1}^n\}$: i.i.d. from $\mathcal{D}_t(\boldsymbol{\beta}^s)$.
- Estimator $\boldsymbol{\beta}^{s2t}$

$$\boldsymbol{\beta}^{s2t} = \text{Est}(\mathbf{X}, \mathbf{y}^s),$$

$$\text{where } \mathbf{X} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top]^\top \in \mathbb{R}^{n \times p}, \mathbf{y}^s = [y_1^s, \dots, y_n^s]^\top \in \mathbb{R}^n$$

Excess (population) Risk for any estimator $\hat{\boldsymbol{\beta}}$

$$\mathcal{R}(\hat{\boldsymbol{\beta}}) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_t(\boldsymbol{\beta}_*)} [(y - \mathbf{x}^\top \hat{\boldsymbol{\beta}})^2] - \sigma_t^2 = \|\boldsymbol{\Sigma}_t^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_*)\|_2^2.$$

\Rightarrow How close the estimator $\hat{\boldsymbol{\beta}}$ is to $\boldsymbol{\beta}_*$?

Two Reference Models

Reference 1: Standard Target Model. (No Distillation)

- Access to the ground-truth parameter β_\star through labeling
- Target Dataset $\{(\mathbf{x}_i, y_i)_{i=1}^n\}$: i.i.d. from $\mathcal{D}_t(\beta_\star)$.
- Estimator β^t :

$$\beta^t := \text{Est}(\mathbf{X}, \mathbf{y}),$$

where $\mathbf{X} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top]^\top \in \mathbb{R}^{n \times p}$, $\mathbf{y} = [y_1, \dots, y_n]^\top \in \mathbb{R}^n$.

Reference 2: Covariance Shift model. [1, 2]

- Data distribution $(\mathbf{x}, y) \sim \mathcal{D}_s^{cs}$
 - $y = \mathbf{x}^\top \beta_\star + z$
 - $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma_s)$, $z \sim \mathcal{N}(0, \sigma_t^2)$
- Target Dataset $\{(\mathbf{x}_i, y_i)_{i=1}^n\}$: i.i.d. from \mathcal{D}_s^{cs}
- Estimator $\hat{\beta}^{cs}$:

$$\hat{\beta}^{cs} := \text{Est}(\mathbf{X}, \mathbf{y}),$$

where $\mathbf{X} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top]^\top \in \mathbb{R}^{n \times p}$, $\mathbf{y} = [y_1, \dots, y_n]^\top \in \mathbb{R}^n$.

Distribution of Ridge(less) Estimator from [3]

Setup: Linear Regression Model

- $y = \mathbf{x}^\top \boldsymbol{\beta}_0 + z$, $\mathbf{x}, \boldsymbol{\beta} \in \mathbb{R}^p$, $y, z \in \mathbb{R}$
- $\mathbb{E}[\mathbf{x}] = 0$, $Cov(\mathbf{x}) = \boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$, $\mathbb{E}[z] = 0$, $Var(z) = \sigma^2$.
- $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times p}$, $\mathbf{y} = [y_1, \dots, y_n]^\top \in \mathbb{R}^n$.

Estimator: Ridge estimator $\hat{\boldsymbol{\beta}}$ with regularization $\eta > 0$ is defined as

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \left\{ \frac{1}{2p} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \frac{\eta}{2} \|\boldsymbol{\beta}\|^2 \right\} = \frac{1}{n} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} + \eta \mathbf{I}_n \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

and the Ridgeless estimator (also known as minimum-norm interpolator) is defined as

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \{ \|\boldsymbol{\beta}\|^2 : \mathbf{y} = \mathbf{X}\boldsymbol{\beta} \} = (\mathbf{X}^\top \mathbf{X})^- \mathbf{X}^\top \mathbf{y}, \quad \mathbf{A}^- : \text{pseudo-inverse of } \mathbf{A}$$

Asymptotic Regime: Fixed $\kappa_t = p/n > 1$ with $n, p \rightarrow \infty$.

Gaussian Sequence Model: For a given pair of $(\boldsymbol{\Sigma}, \boldsymbol{\beta}_0)$ and noise level $\gamma > 0$,

$$y_{(\boldsymbol{\Sigma}, \boldsymbol{\beta}_0)}^{seq}(\gamma) := \boldsymbol{\Sigma}^{1/2} \boldsymbol{\beta}_0 + \frac{\gamma \mathbf{g}}{\sqrt{p}}, \quad \mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p).$$

Distribution of Ridge(less) Estimator from [3]

Ridge estimator with regularization $\tau \geq 0$ in the Gaussian sequence model:

$$\begin{aligned}\hat{\beta}_{(\Sigma, \beta_0)}^{seq}(\gamma; \tau) &:= \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \|\Sigma^{1/2} \beta - y_{(\Sigma, \beta_0)}^{seq}(\gamma)\|^2 + \frac{\tau}{2} \|\beta\|^2 \right\} \\ &= (\Sigma + \tau I_p)^{-1} \Sigma^{1/2} \left(\Sigma^{1/2} \beta_0 + \frac{\gamma g}{\sqrt{p}} \right)\end{aligned}$$

Distributional Characterization (Informal).

For any $\eta \geq 0$, there exists a unique pair $(\gamma_{\eta,*}, \tau_{\eta,*}) \in (0, \infty)^2$ determined via a fixed point equation, such that the distribution $\hat{\beta}$ is about the same as that of $\hat{\beta}_{(\Sigma, \beta)}^{seq}(\gamma_{\eta,*}; \tau_{\eta,*})$.

Formally, for any 1-Lipschitz function $g : \mathbb{R}^n \rightarrow \mathbb{R}$ and any $K > 0$, with high probability,

$$\sup_{\eta \in [0, K]} \left| g(\hat{\beta}) - \mathbb{E} \left[g(\hat{\beta}_{(\Sigma, \beta)}^{seq}(\gamma_{\eta,*}; \tau_{\eta,*})) \right] \right| \approx 0.$$

Main Idea from [3]: Analyze the behavior of Ridge(less) estimator $\hat{\beta}$ using $\hat{\beta}_{(\Sigma, \beta)}^{seq}(\gamma_{\eta,*}; \tau_{\eta,*})$!

Fixed Point Equation from [3]

For $\gamma, \tau > 0$, define the error $\text{err}_{(\Sigma, \beta_0)}(\gamma; \tau)$ and the degrees-of-freedom $\text{dof}_{(\Sigma, \beta_0)}(\gamma; \tau)$ as

$$\begin{aligned}\text{err}_{(\Sigma, \beta_0)}(\gamma; \tau) &:= \left\| \Sigma^{1/2} \left(\hat{\beta}_{(\Sigma, \beta)}^{seq}(\gamma; \tau) - \beta_0 \right) \right\|^2 \\ \text{dof}_{(\Sigma, \beta_0)}(\gamma; \tau) &:= \left\langle \frac{\gamma \mathbf{g}}{\sqrt{p}}, \left(\hat{\beta}_{(\Sigma, \beta)}^{seq}(\gamma; \tau) - \beta_0 \right) \right\rangle\end{aligned}$$

Fixed Point Equation: For $\eta \geq 0$,

$$\begin{cases} \kappa_t^{-1} \gamma^2 = \sigma^2 + \mathbb{E}[\text{err}_{(\Sigma, \beta_0)}(\gamma; \tau)], \\ \kappa_t^{-1} - \frac{\eta}{\tau} = \frac{1}{p} \text{tr}((\Sigma + \tau \mathbf{I}_p)^{-1} \Sigma) = \frac{1}{\gamma^2} \mathbb{E}[\text{dof}_{(\Sigma, \beta_0)}(\gamma; \tau)]. \end{cases}$$

For any $\eta \geq 0$, the fixed equation has a unique solution $(\gamma_{\eta,*}, \tau_{\eta,*}) \in (0, \infty)^2$.

Returning to the Original Target Model

Recall the setup of the target model:

- Data distribution $(\mathbf{x}, y^s) \sim \mathcal{D}_t(\boldsymbol{\beta}^s)$
 - $y^s = \mathbf{x}^\top \boldsymbol{\beta}^s + z$
 - $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_t)$, $\boldsymbol{\Sigma}_t$: Covariance matrix of distribution $\mathcal{D}_t(\boldsymbol{\beta}^s)$
 - $z \sim \mathcal{N}(0, \sigma_t^2)$: independent of \mathbf{x}
- Target Dataset $\{(\mathbf{x}_i, y_i^s)_{i=1}^n\}$: i.i.d. from $\mathcal{D}_t(\boldsymbol{\beta}^s)$.
- Estimator $\boldsymbol{\beta}^{s2t}$

$$\boldsymbol{\beta}^{s2t} = \text{Est}(\mathbf{X}, \mathbf{y}^s),$$

$$\text{where } \mathbf{X} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top]^\top \in \mathbb{R}^{n \times p}, \mathbf{y}^s = [y_1^s, \dots, y_n^s]^\top \in \mathbb{R}^n$$

Main Idea: $\boldsymbol{\beta}^{s2t} \approx \hat{\boldsymbol{\beta}}_{(\boldsymbol{\Sigma}_t, \boldsymbol{\beta}^s)}^{seq}(\gamma_t; \tau_t)$, where γ_t, τ_t satisfies the fixed point equation

$$\begin{cases} \gamma_t^2 = \kappa_t (\sigma_t^2 + \mathbb{E}[\text{err}_{(\boldsymbol{\Sigma}_t, \boldsymbol{\beta}^s)}(\gamma_t; \tau_t)]) , \\ \kappa_t^{-1} = \frac{1}{p} \text{tr}((\boldsymbol{\Sigma}_t + \tau_t \mathbf{I}_p)^{-1} \boldsymbol{\Sigma}_t) \end{cases}$$

and

$$\hat{\boldsymbol{\beta}}_{(\boldsymbol{\Sigma}_t, \boldsymbol{\beta}^s)}^{seq}(\gamma_t; \tau_t) = (\boldsymbol{\Sigma}_t + \tau_t \mathbf{I}_p)^{-1} \boldsymbol{\Sigma}_t^{1/2} \left(\boldsymbol{\Sigma}_t^{1/2} \boldsymbol{\beta}^s + \frac{\gamma_t \mathbf{g}_t}{\sqrt{p}} \right), \quad \mathbf{g}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p).$$

Asymptotic Risk Estimate

Excess Risk $\mathcal{R}(\hat{\beta})$ for any estimator $\hat{\beta} \in \mathbb{R}^p$:

$$\mathcal{R}(\hat{\beta}) := \mathbb{E}_{(x,y) \sim \mathcal{D}_t(\beta_\star)} [(y - x^\top \hat{\beta})^2] - \sigma_t^2 = \|\Sigma_t^{1/2}(\hat{\beta} - \beta_\star)\|_2^2.$$

Risk of the surrogate-to-target problem: $\mathcal{R}(\beta^{s2t})$:

$$\mathcal{R}(\beta^{s2t}) = \|\Sigma_t^{1/2}(\beta^{s2t} - \beta_\star)\|_2^2$$

Asymptotic Risk Estimate $\bar{\mathcal{R}}_{\kappa_t, \sigma_t}^{s2t}(\Sigma_t, \beta_\star, \beta^s)$:

$$\bar{\mathcal{R}}_{\kappa_t, \sigma_t}^{s2t}(\Sigma_t, \beta_\star, \beta^s) = \mathbb{E}_{g_t} [\mathcal{R}(\hat{\beta}_{(\Sigma_t, \beta^s)}^{seq}(\gamma_t; \tau_t))] = \mathbb{E}_{g_t} \left[\left\| \Sigma_t^{1/2} \left(\hat{\beta}_{(\Sigma_t, \beta^s)}^{seq}(\gamma_t; \tau_t) - \beta_\star \right) \right\|_2^2 \right]$$

Recall the definition of $\hat{\beta}_{(\Sigma_t, \beta^s)}^{seq}(\gamma_t; \tau_t)$:

$$\begin{aligned} \hat{\beta}_{(\Sigma_t, \beta^s)}^{seq}(\gamma_t; \tau_t) &= (\Sigma_t + \tau_t I_p)^{-1} \Sigma_t^{1/2} \left(\Sigma_t^{1/2} \beta^s + \frac{\gamma_t g_t}{\sqrt{p}} \right) \\ &= \underbrace{(\Sigma_t + \tau_t I_p)^{-1} \Sigma_t}_{:= \theta_1} \beta^s + \underbrace{\gamma_t (\Sigma_t + \tau_t I_p)^{-1} \Sigma_t^{1/2} \frac{g_t}{\sqrt{p}}}_{:= \theta_2} \\ &= \theta_1 \beta^s + \gamma_t \theta_2 \end{aligned}$$

Note that $\mathbb{E}_{\mathbf{g}_t}[\boldsymbol{\theta}_2] = 0$ and $\boldsymbol{\theta}_1$ is not a function of \mathbf{g}_t .

$$\begin{aligned}\bar{\mathcal{R}}_{\kappa_t, \sigma_t}^{s2t}(\boldsymbol{\Sigma}_t, \boldsymbol{\beta}_\star, \boldsymbol{\beta}^s) &= \mathbb{E}_{\mathbf{g}_t} \left[(\boldsymbol{\theta}_1 \boldsymbol{\beta}^s + \gamma_t \boldsymbol{\theta}_2 - \boldsymbol{\beta}_\star)^\top \boldsymbol{\Sigma}_t (\boldsymbol{\theta}_1 \boldsymbol{\beta}^s + \gamma_t \boldsymbol{\theta}_2 - \boldsymbol{\beta}_\star) \right] \\ &= (\boldsymbol{\theta}_1 \boldsymbol{\beta}^s - \boldsymbol{\beta}_\star)^\top \boldsymbol{\Sigma}_t (\boldsymbol{\theta}_1 \boldsymbol{\beta}^s - \boldsymbol{\beta}_\star) + \gamma_t^2 \mathbb{E}_{\mathbf{g}_t} [\boldsymbol{\theta}_2^\top \boldsymbol{\Sigma}_t \boldsymbol{\theta}_2]\end{aligned}$$

The former term can be expressed as

$$\begin{aligned}(\boldsymbol{\theta}_1 (\boldsymbol{\beta}^s - \boldsymbol{\beta}_\star) - (\mathbf{I} - \boldsymbol{\theta}_1) \boldsymbol{\beta}_\star)^\top \boldsymbol{\Sigma}_t (\boldsymbol{\theta}_1 (\boldsymbol{\beta}^s - \boldsymbol{\beta}_\star) - (\mathbf{I} - \boldsymbol{\theta}_1) \boldsymbol{\beta}_\star) \\ = (\boldsymbol{\beta}^s - \boldsymbol{\beta}_\star)^\top \boldsymbol{\theta}_1^\top \boldsymbol{\Sigma}_t \boldsymbol{\theta}_1 (\boldsymbol{\beta}^s - \boldsymbol{\beta}_\star) + \boldsymbol{\beta}_\star^\top (\mathbf{I} - \boldsymbol{\theta}_1)^\top \boldsymbol{\Sigma}_t (\mathbf{I} - \boldsymbol{\theta}_1) \boldsymbol{\beta}_\star \\ - 2 \boldsymbol{\beta}_\star^\top (\mathbf{I} - \boldsymbol{\theta}_1)^\top \boldsymbol{\Sigma}_t (\boldsymbol{\beta}^s - \boldsymbol{\beta}_\star)\end{aligned}$$

Hence, we can formulate Asymptotic risk estimate as

$$\begin{aligned}\bar{\mathcal{R}}_{\kappa_t, \sigma_t}^{s2t}(\boldsymbol{\Sigma}_t, \boldsymbol{\beta}_\star, \boldsymbol{\beta}^s) &:= (\boldsymbol{\beta}^s - \boldsymbol{\beta}_\star)^\top \boldsymbol{\theta}_1^\top \boldsymbol{\Sigma}_t \boldsymbol{\theta}_1 (\boldsymbol{\beta}^s - \boldsymbol{\beta}_\star) + \gamma_t^2 \mathbb{E}_{\mathbf{g}_t} [\boldsymbol{\theta}_2^\top \boldsymbol{\Sigma}_t \boldsymbol{\theta}_2] \\ &\quad + \boldsymbol{\beta}_\star^\top (\mathbf{I} - \boldsymbol{\theta}_1)^\top \boldsymbol{\Sigma}_t (\mathbf{I} - \boldsymbol{\theta}_1) \boldsymbol{\beta}_\star - 2 \boldsymbol{\beta}_\star^\top (\mathbf{I} - \boldsymbol{\theta}_1)^\top \boldsymbol{\Sigma}_t (\boldsymbol{\beta}^s - \boldsymbol{\beta}_\star)\end{aligned}$$

Recall the fixed point equation:

$$\begin{cases} \gamma_t^2 = \kappa_t (\sigma_t^2 + \mathbb{E}[\text{err}_{(\boldsymbol{\Sigma}_t, \boldsymbol{\beta}^s)}(\gamma_t; \tau_t)]) = \kappa_t (\sigma_t^2 + \bar{\mathcal{R}}_{\kappa_t, \sigma_t}^{s2t}(\boldsymbol{\Sigma}_t, \boldsymbol{\beta}^s, \boldsymbol{\beta}^s)) \\ \kappa_t^{-1} = \frac{1}{p} \text{tr}((\boldsymbol{\Sigma}_t + \tau_t \mathbf{I}_p)^{-1} \boldsymbol{\Sigma}_t) \end{cases}$$

Non-Asymptotic Characterization of the Risk

Theorem 1 (Theorem 2.3 in [3])

Suppose that some constant $M_t > 1$, we have $1/M_t \leq \kappa_t, \sigma_t^2 \leq M_t$ and $\|\Sigma_t\|_{op}, \|\Sigma_t^{-1}\|_{op} \leq M_t$. Let $B_p(R) := \{\mathbf{x} \in \mathbb{R}^p : \|\mathbf{x}\|_2 \leq R\}$. Then, there exists a constant $C = C(M_t)$ such that, for any $\epsilon \in (0, 1/2]$, the following holds with $R + 1 < M_t$:

$$\sup_{\beta_*, \beta^s \in B_p(R)} P(|\mathcal{R}(\beta^{s2t}) - \bar{\mathcal{R}}_{\kappa_t, \sigma_t}^{s2t}(\Sigma_t, \beta_*, \beta^s)| \geq \epsilon) \leq Cpe^{-p\epsilon^4/C}.$$

Meaning of Theorem 1: In the asymptotic regime, $\mathcal{R}(\beta^{s2t}) \xrightarrow{p} \bar{\mathcal{R}}_{\kappa_t, \sigma_t}^{s2t}(\Sigma_t, \beta_*, \beta^s)$

Proposition 1

Let $\Omega = \frac{\text{tr}(\Sigma_t^2(\Sigma_t + \tau_t I_p)^{-2})}{n}$. The optimal surrogate β^s minimizing the asymptotic risk $\bar{\mathcal{R}}_{\kappa_t, \sigma_t}^{s2t}(\Sigma_t, \beta_*, \beta^s)$ is

$$\beta^{s*} = \left((\Sigma_t + \tau_t I_p)^{-1} \Sigma_t + \frac{\Omega \tau_t^2}{1 - \Omega} \Sigma_t^{-1} (\Sigma_t + \tau_t I_p)^{-1} \right) \beta_*.$$

Proof: Differentiate $\bar{\mathcal{R}}_{\kappa_t, \sigma_t}^{s2t}(\Sigma_t, \beta_*, \beta^s)$ with respect to β^s .

Optimal Surrogate Parameter is not β_\star

Corollary 1

Without loss of generality, suppose that Σ_t is diagonal. Let $(\lambda_i)_{i=1}^p$ be the eigenvalues of Σ_t in non-increasing order and let $\xi_i = \frac{\tau_t}{\lambda_i + \tau_t}$ for $i \in [p]$. Then, the following results hold:

1. $(\beta^{s*})_i = (\beta_\star)_i \left((1 - \xi_i) + \xi_i \frac{\Omega}{1 - \Omega} \frac{\xi_i}{1 - \xi_i} \right)^{-1}$ for every $i \in [p]$.
2. $|(\beta^{s*})_i| > |(\beta_\star)_i|$ if and only if $1 - \xi_i > \Omega = \frac{\sum_{j=1}^p (1 - \xi_j)^2}{\sum_{j=1}^p (1 - \xi_j)}$ for every $i \in [p]$.
3. $\beta^{s*} = \beta_\star$ if and only if the covariance matrix $\Sigma_t = cI_p$ for some $c \in \mathbb{R}$.

Meaning of Corollary 1

1. Optimal surrogate parameter β^{s*} only depends on the covariance spectrum λ_i .
2. For regions where λ_i is small, $1 - \xi_i$ is also small, leading to further amplification through the ground truth parameter β_\star . Conversely, for regions with large λ_i , β_\star induces shrinkage.
3. The threshold Ω corresponds to the ratio of the sample second moment to the sample first moment, arising from the trade-off between bias and variance terms of $\bar{\mathcal{R}}_{\kappa_t, \sigma_t}^{s2t}(\Sigma_t, \beta_\star, \beta^s)$.
4. Unless the eigenvalues of the Σ_t are constant, there is potential for improvement by using the surrogate parameter β^s rather than using β_\star .

Weak-to-Strong Generalization

How to design *weak model*?

- Surrogate model uses fewer features, $p_s < p$, by a mask operation $\mathcal{M}(\mathbf{x})$.

- **Masked Surrogate Model:** $(\mathcal{M}(\tilde{\mathbf{x}}), \tilde{y}) \sim \mathcal{D}_s^{p_s}$: Masked distribution

- $\tilde{y} = \mathcal{M}(\tilde{\mathbf{x}})^\top \mathcal{M}(\boldsymbol{\beta}_*) + \tilde{z}$

- $\boldsymbol{\beta}_* \in \mathbb{R}^p$, $\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_s)$, $\tilde{z} \sim \mathcal{N}(0, \sigma_s^2)$.

\Rightarrow Estimator $\boldsymbol{\beta}^s = \text{Est}(\tilde{\mathbf{X}}, \tilde{\mathbf{y}})$

- **Masked Target Model:** $(\mathbf{x}, y^s) \sim \mathcal{D}_t^{p_s}(\boldsymbol{\beta}^s)$

- $y^s = \mathcal{M}(\mathbf{x})^\top \boldsymbol{\beta}^s + z$

- $\boldsymbol{\beta}^s \in \mathbb{R}^{p_s}$, $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_t)$, $z \sim \mathcal{N}(0, \sigma_t^2)$.

\Rightarrow Estimator $\boldsymbol{\beta}^{s2t} = \text{Est}(\mathbf{X}, \mathbf{y})$

- **Standard Target Model:** $(\mathbf{x}, y^s) \sim \mathcal{D}_t^p(\boldsymbol{\beta}_*)$

- $y = \mathbf{x}^\top \boldsymbol{\beta}_* + z$

- $\boldsymbol{\beta}_* \in \mathbb{R}^p$, $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_t)$, $z \sim \mathcal{N}(0, \sigma_t^2)$.

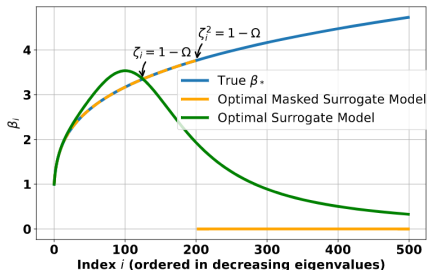
\Rightarrow Estimator $\boldsymbol{\beta}^t = \text{Est}(\mathbf{X}, \mathbf{y})$

Leveraging Information from a Weak Teacher

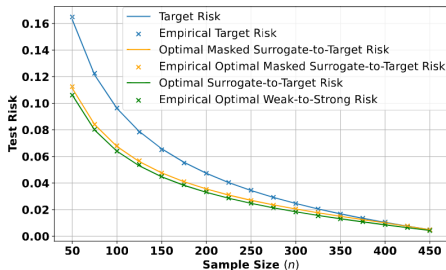
Proposition 2

Assume that $\Sigma_s = \Sigma_t$ and Σ_t is diagonal. In the absence of model shift ($\mathcal{M}(\beta_*) = \beta^s$), the following results hold:

1. If the mask operation \mathcal{M} select all the features that satisfy $1 - \xi_i^2 > \Omega$, then the surrogate-to-target model outperforms the standard target model in the asymptotic risk.
2. Let \mathcal{M} represent the all possible \mathcal{M} , where $|\mathcal{M}| = 2^p$. The optimal \mathcal{M}^* for the asymptotic risk within \mathcal{M} is the one that selects all features satisfying $1 - \xi^2 > \Omega$.



(a) Ground-truth and surrogate model weights



(b) Test risks as a function of sample size

Fundamental Limits and Scaling Laws

Observation 2. For any covariance matrix $\Sigma \in \mathbb{R}^{p \times p}$, there exists an orthonormal matrix $U \in \mathbb{R}^{p \times p}$ such that the transformation of $\mathbf{x} \rightarrow U^\top \mathbf{x}$ and $\beta \rightarrow U^\top \beta$ does not affect the labels \mathbf{y} but ensures that the covariance matrix is diagonal. \Rightarrow Consider only the diagonal covariance!

Omniscient Test Risk Estimate

Fix $p > n \geq 1$. Let $\Sigma = U \text{diag}(\lambda) U^\top$ and $\bar{\beta} = U^\top \beta_\star$. Then, the excess test risk estimate is the following:

$$\mathcal{R}(\hat{\beta}) \approx \mathbb{E}_{\hat{\beta} \sim D(\beta_\star)} \left[(y - \mathbf{x}^\top \hat{\beta})^2 \right] - \sigma^2 = \frac{\sigma^2 \Omega + \mathcal{B}(\bar{\beta})}{1 - \Omega},$$

$$\text{where } n = \sum_{i=1}^p \frac{1}{\lambda_i + \tau}, \quad \xi_i = \frac{\tau}{\lambda_i + \tau}, \quad \Omega = \frac{1}{n} \sum_{i=1}^p (1 - \xi_i)^2, \quad \mathcal{B}(\bar{\beta}) = \sum_{i=1}^p \lambda_i \xi_i^2 \bar{\beta}_i^2.$$

Asymptotic behavior of omniscient risk: As $n, p \rightarrow \infty$ with a fixed ratio $p/n = \kappa$, the approximation becomes an equality. (Same as Theorem 1)

What is $D(\beta_\star)$?

$$\hat{\beta} = (\Sigma + \tau \mathbf{I}_p)^{-1} \Sigma^{1/2} \left(\Sigma^{1/2} \beta_\star + \frac{\gamma \mathbf{g}}{\sqrt{p}} \right), \quad (\gamma, \tau) : \text{sol. of fixed point eq., } \mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p).$$

Framework for Degree Analysis

Setup

- $p, n \rightarrow \infty$ with a fixed ratio $\kappa = p/n > 1$. (Over-parametrized regime)
- Covariance Σ : diagonal, $\Sigma_{i,i} = \lambda_i = i^{-\alpha}$ for $\alpha > 1$.

Equations for the τ_t and Ω :

$$\sum_{i=1}^{\infty} \frac{\tau_t}{\lambda_i + \tau_t} = n, \quad n\Omega = \sum_{i=1}^{\infty} \left(\frac{i^{-\alpha}}{i^{-\alpha} + \tau_t} \right)^2$$

Proposition 3. When $\lambda_i = i^{-\alpha}$,

$$\tau_t = cn^{-\alpha}(1 + O(n^{-1})), \Omega = \frac{\alpha - 1}{\alpha} - O(n^{-1}) \text{ for } c = \left(\frac{\pi}{\alpha \sin(\pi/\alpha)} \right)^\alpha.$$

Proposition 4. When $C_1 := \frac{\alpha \sin(\pi/\alpha)}{n(\alpha-1)^{1/\alpha}}$ and $C_2 := \frac{\alpha \sin(\pi/\alpha)}{n(\sqrt{\alpha}-1)^{1/\alpha}}$. Then, the indices i for $\xi_i < 1 - \Omega$ are $i < nC_1 + O(1)$; while the indices for $\xi_i^2 < 1 - \Omega$ are $i < nC_2 + O(1)$.

Meaning of Proposition 4: As sample size n increases, the cut-off indices of β^{s*} and the optimal mask \mathcal{M}^* increase linearly in n .

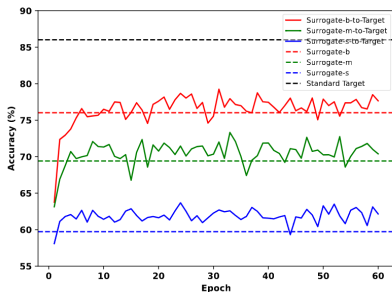
Scaling Law for S2T Model

Scaling Law

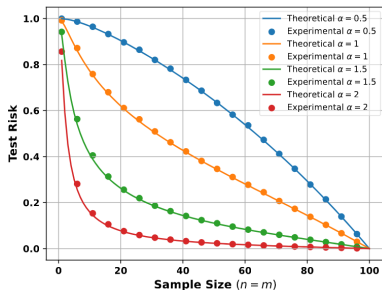
Assume $\lambda_i = i^{-\alpha}$ and $\lambda_i \beta_i^2 = i^{-\beta}$ for $\alpha, \beta > 1$. In the asymptotic regime ($p \rightarrow \infty$), the excess risk of the surrogate-to-target model with an *optimal* surrogate parameter scales the same as the standard target model. Specifically, we have

$$\mathcal{R}^*(\beta^{s^{2t}}) = \Theta(n^{-\min(2\alpha, \beta-1)}) = \mathcal{R}(\beta^t),$$

where the minimum surrogate-to-target risk $\mathcal{R}^*(\beta^{s^{2t}})$ is attained by β^{s^*} .



(a) Weak-to-strong on CIFAR-10



(b) Comparison of theoretical and experimental risks

Asymptotic Risk for the Two-Stage Model

Definition 3. Recall the definition of τ_t and γ_t in Theorem 1. Let $\kappa_s = p/m > 1$ and define $\tau_s \in \mathbb{R}$ similarly to τ_t . We define the random variable $X_{(\kappa_s, \sigma_s^2)}^s$ based on $\mathbf{g}_s \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and the function $\gamma_s : \mathbb{R}^p \rightarrow \mathbb{R}$ as follows:

$$X_{(\kappa_s, \sigma_s^2)}^s(\Sigma_s, \beta_\star, \mathbf{g}_s) := (\Sigma_s + \tau_s \mathbf{I}_p)^{-1} \Sigma_s^{1/2} \left(\Sigma_s^{1/2} \beta_\star + \frac{\gamma_s(\beta_\star) \mathbf{g}_s}{\sqrt{p}} \right)$$

$$\gamma_s^2(\beta_\star) := \kappa_s \left(\sigma_s^2 + \mathbb{E}_{\mathbf{g}_s} \left[\Sigma_s^{1/2} \left(X_{(\kappa_s, \sigma_s^2)}^s(\Sigma_s, \beta_\star, \mathbf{g}_s) - \beta_\star \right)^2 \right] \right)$$

Let $\dot{\kappa} = (\kappa_s, \kappa_t)$, $\dot{\Sigma} = (\Sigma_s, \Sigma_t)$ and $\dot{\sigma} = (\sigma_s, \sigma_t)$. Then, the asymptotic risk estimate is

$$\begin{aligned} \bar{\mathcal{R}}_{\dot{\kappa}, \dot{\sigma}}(\dot{\Sigma}, \beta_\star) &= \|\sigma_t^{1/2} (\mathbf{I} - (\Sigma_t + \tau_t \mathbf{I})^{-1} \Sigma_t (\Sigma_s + \tau_s \mathbf{I}) \Sigma_s) \beta_\star\|_2^2 \\ &+ \frac{\mathbb{E}_{\beta^s \sim X^s} [\gamma_t^2(\beta^s)]}{p} \text{tr}(\Sigma_t^2 (\Sigma_t + \tau_t \mathbf{I})^{-2}) \\ &+ \frac{\gamma_s^2(\beta_\star)}{p} \text{tr} \left(\Sigma_s^{1/2} (\Sigma_s + \tau_s \mathbf{I}) \Sigma_t (\Sigma_t + \tau_t \mathbf{I})^{-1} \Sigma_t (\Sigma_t + \tau_t \mathbf{I})^{-1} \Sigma_t (\Sigma_s + \tau_s \mathbf{I})^{-1} \Sigma_s^{1/2} \right). \end{aligned}$$

Non-Asymptotic Risk of Two-Stage Model

Theorem 2

Suppose that some constant $M_t > 1$, we have $1/M_t \leq \kappa_s, \sigma_s^2, \kappa_t, \sigma_t^2 \leq M_t$ and $\|\Sigma_s\|_{op}, \|\Sigma_s^{-1}\|_{op}, \|\Sigma_t\|_{op}, \|\Sigma_t^{-1}\|_{op} \leq M_t$. Let $B_p(R) := \{\mathbf{x} \in \mathbb{R}^p : \|\mathbf{x}\|_2 \leq R\}$. Then, there exists a constant $C = C(M_t)$ such that, for any $\epsilon \in (0, 1/2]$, the following holds with $R + 1 < M_t$:

$$\sup_{\beta_\star \in B_p(R)} P\left(|\mathcal{R}(\beta^{s^{2t}}) - \bar{\mathcal{R}}_{\hat{\kappa}, \hat{\sigma}}(\hat{\Sigma}, \beta_\star)| \geq \epsilon\right) \leq Cpe^{-p\epsilon^4/C}.$$

Future research direction

- Extend the two-stage process to multiple stages.
- Apply the two-stage learning to data pruning, using the surrogate model to decide whether keep or discard each data during the training of the target model.

Summary

Problem Setup

- Surrogate model: $\tilde{y} = \tilde{\mathbf{x}}^\top \boldsymbol{\beta}_\star + \tilde{z}$, $\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_s)$, $\tilde{z} \sim \mathcal{N}(0, \sigma_s^2) \rightarrow$ estimate: $\boldsymbol{\beta}^s$
- Target model: $y^s = \mathbf{x}^\top \boldsymbol{\beta}^s + z$, $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_t)$, $z \sim \mathcal{N}(0, \sigma_t^2) \rightarrow$ estimate: $\boldsymbol{\beta}^{s2t}$
- Standard target model: $y = \mathbf{x}^\top \boldsymbol{\beta}_\star + z$, $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_t)$, $z \sim \mathcal{N}(0, \sigma_t^2) \rightarrow$ estimate: $\boldsymbol{\beta}^t$
- Excess Risk: $\mathcal{R}(\hat{\boldsymbol{\beta}}) = \|\boldsymbol{\Sigma}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_\star)\|^2$

Asymptotic Risk[3]

- Asymptotic regime: $n, p \rightarrow \infty$, $\kappa = p/n > 1$
- Theorem 1, Theorem 2: non-asymptotic risk $\rightarrow 0$, excess risk \rightarrow asymptotic risk
- Only consider the asymptotic risk for further analysis

Optimal Surrogate Parameter

- $\boldsymbol{\beta}_\star$ is not a optimal surrogate parameter! \rightarrow room for improvement from the surrogate model
- Optimal surrogate parameter only depends on the eigenvalues of the covariance matrix.
- Threshold: $1 - \xi > \Omega \rightarrow$ originates from the trade-off between bias and variance.

Summary

Weak-to-Strong Generalization

- Masked surrogate model: Only uses fewer features by a mask operation \mathcal{M}
- Select all features with $1 - \xi^2 > \Omega$, the target model outperforms the standard target model.
 \Rightarrow Weak-to-strong model outperforms strong model!

Scaling law

- Assumption: Σ : diagonal, $\Sigma_{i,i} = \lambda_i = i^{-\alpha}$, $\lambda_i \beta_i^2 = i^{-\beta}$
- $\mathcal{R}^*(\beta^{s2t}) = \Theta(n^{-\min(2\alpha, \beta-1)}) = \mathcal{R}(\beta^t)$

Experiment

- CIFAR-10: Surrogate model < Target model < Standard target model
- Why Target model < Standard target model? \rightarrow Due to the lack of a feature selection process

Future direction

- Multi-stage process (Multi-round distillation)
- Data pruning \rightarrow Utilize the surrogate model for data selection

References

- [1] N. Mallinar, A. Zane, S. Frei, and B. Yu, “Minimum-norm interpolation under covariate shift,” *arXiv preprint arXiv:2404.00522*, 2024.
- [2] P. Patil, J.-H. Du, and R. J. Tibshirani, “Optimal ridge regularization for out-of-distribution prediction,” *arXiv preprint arXiv:2404.01233*, 2024.
- [3] Q. Han and X. Xu, “The distribution of ridgeless least squares interpolators,” *arXiv preprint arXiv:2307.02044*, 2023.