# Optimal Ridge Regularization for Out-of-Distribution Prediction

Pratik Patil[1]    Jin-Hong Du[2]    Ryan J. Tibshirani[1]

[1]University of California, Berkeley
[2]Carnegie Mellon University

# Ridge regression in high dimensions

Recent interests in high-dimensional ridge regression concern the ridge estimator:

$$\widehat{\boldsymbol{\beta}}^\lambda = (\boldsymbol{X}^\top \boldsymbol{X}/n + \lambda \boldsymbol{I}_p)^\dagger \boldsymbol{X}^\top \boldsymbol{y}/n,$$

and its prediction risk:

$$R(\widehat{\boldsymbol{\beta}}^\lambda) = \mathbb{E}_{\boldsymbol{x}_0, y_0}[(y_0 - \boldsymbol{x}_0^\top \widehat{\boldsymbol{\beta}}^\lambda)^2 \mid \boldsymbol{X}, \boldsymbol{y}].$$

The goal is to study the behavior of its asymptotic prediction risk:

$$R(\widehat{\boldsymbol{\beta}}^\lambda) \to \mathscr{R}(\lambda, \phi)$$

as the feature size $p$ and the sample size $n$ diverge proportionally to an *aspect ratio* $p/n \to \phi \in (0, \infty)$.

# Optimal ridge regression under in-distribution

For high-dimensional ridge regression, two questions for the optimal in-distribution asymptotic risk $\min_{\lambda \geq \lambda_{\min}} \mathscr{R}(\lambda, \phi)$:

(Q1) What is the behavior of the *optimal ridge penalty*, as a function of parameters such as signal-to-noise ratio, data aspect ratio, feature correlations, and signal structure?

(Q2) What is the behavior of the *optimally tuned ridge risk*, as a function of these same problem parameters?

Known results provide partial answers:

(A1) $\lambda^* = \phi/\text{SNR} > 0$ in the isotropic cases when $\lambda_{\min} = 0$, while $\lambda^* < 0$ in some anisotropic cases (both signal and features) and overparameterized regimes.

(A2) $\mathscr{R}(\lambda^*, \phi)$ is monotonically increasing in $\phi$.

# Ridge regression under distribution shifts (motivation)

We consider two types of distribution shifts:

(i) *Covariate shift*: where $P_{x_0} \neq P_x$ but $P_{y_0|x_0} = P_{y|x}$.

(ii) *Regression shift*: where $P_{y_0|x_0} \neq P_{y|x}$ but $P_{x_0} = P_x$.

and answer two out-of-distribution problems:

(Q1') *How does distribution shift alter optimal regularization $\lambda^*$?*

(Q2') *How does distribution shift alter optimal risk behavior $\mathscr{R}(\lambda^*, \phi)$?*

## Optimal regularization landscape in ridge regression.

| Σ | β | Σ₀ | β₀ | $\phi \lesssim 1$ | $\lambda_{\min}$ | Arb. Mod. | Arb. SNR | Arb. Spec. | Additional Specific Data Geometry Conditions | $\lambda^\star$ | Reference |
|---|---|---|---|---|---|---|---|---|---|---|---|
| *In-distribution* | | | | | | | | | | | |
| ⊗ | ○ | Σ | β | all | zero | ✗ | ✓ | ✗ | | + | [DW, Thm. 2.1] |
| ○ | ⊗ | Σ | β | all | zero | ✗ | ✓ | ✗ | | + | [HMRT, Cor. 5] |
| | | | | under | neg | ✗ | ✓ | ✗ | | + | [WX, Prop. 6] |
| | | | | over | neg | ✗ | ✗ | ✗ | Strict misalignment of $(\Sigma, \beta)$ | + | [WX, Thm. 4] |
| | | | | over | neg | ✗ | ✗ | ✗ | Strict alignment of $(\Sigma, \beta)$ | − | [WX, Thm. 4, Prop. 7] |
| ⊗ | ⊗ | Σ | β | over | zero | ✗ | ✗ | ✗ | and/or special feature model | 0 | [RMR, Cor. 2] |
| | | | | under | neg$^\star$ | ✓ | ✓ | ✓ | | + | **Theorem 2** (1) |
| | | | | over | neg$^\star$ | ✓ | ✓ | ✓ | General alignment of $(\Sigma, \beta, \sigma^2)$ | − | **Theorem 2** (2) |
| *Out-of-distribution* | | | | | | | | | | | |
| ⊗ | ○ | Σ₀ | β | all | neg$^\star$ | ✓ | ✓ | ✓ | | + | **Proposition 3** |
| ⊗ | ⊗ | Σ₀ | β | under | neg$^\star$ | ✓ | ✓ | ✓ | | + | **Theorem 4** (1) |
| ⊗ | ⊗ | I | β | over | neg$^\star$ | ✓ | ✓ | ✓ | | + | **Theorem 4** (2) |
| ○ | ⊗ | Σ₀ | β | over | neg$^\star$ | ✓ | ✓ | ✓ | General alignment of $(\Sigma_0, \beta, \sigma^2)$ | − | **Theorem 4** (3) |
| | | | | under | neg$^\star$ | ✓ | ✓ | ✓ | General alignment of $(\Sigma, \beta, \beta_0)$ | − | **Theorem 5** (1), (39) |
| ⊗ | ⊗ | Σ | β₀ | under | neg$^\star$ | ✓ | ✓ | ✓ | General misalignment of $(\Sigma, \beta, \beta_0)$ | + | **Theorem 5** (1), (39) |
| | | | | over | neg$^\star$ | ✓ | ✓ | ✓ | General alignment of $(\Sigma, \beta, \beta_0, \sigma^2)$ | − | **Theorem 5** (2) |

# Data assumptions and lower bound on negative regularization

Data assumptions:

► Covariate: Each feature vector $x_i$ for $i \in [n]$ can be decomposed as $x_i = \Sigma^{1/2} z_i$, where $z_i \in \mathbb{R}^p$ contains i.i.d. entries $z_{ij}$ for $j \in [p]$ with mean $0$, variance $1$, and bounded $4 + \mu$ moments for some $\mu > 0$. (RMT structure and bounded moment)

► Response: Each response variable $y_i$ for $i \in [n]$ has mean $0$, and bounded $4 + \mu$ moments. (model-free)

Lower bound on $\lambda$: Let $\mu_{\min} \in \mathbb{R}$ be the unique solution, satisfying $\mu_{\min} > -r_{\min}$, to the equation:

$$1 = \phi \, \bar{\mathrm{tr}}[\Sigma^2 (\Sigma + \mu_{\min} I)^{-2}],$$

and let $\lambda_{\min}(\phi)$ be given by:

$$\lambda_{\min}(\phi) = \mu_{\min} - \phi \, \bar{\mathrm{tr}}[\Sigma (\Sigma + \mu_{\min} I)^{-1}].$$

# Out-of-distribution risk characterization

The OOD risk asymptotics read that

$$\mathscr{R}(\lambda, \phi) := \underbrace{\mathscr{B}(\lambda, \phi)}_{\text{bias}} + \underbrace{\mathscr{V}(\lambda, \phi)}_{\text{variance}} + \underbrace{\mathscr{E}(\lambda, \phi)}_{\text{extra bias}} + \underbrace{\kappa^2}_{\text{irreducible error}}, \quad (1)$$

where

$$\mathscr{B} = \mu^2 \cdot \boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \mu \boldsymbol{I})^{-1} (\widetilde{\nu} \boldsymbol{\Sigma} + \boldsymbol{\Sigma}_0) (\boldsymbol{\Sigma} + \mu \boldsymbol{I})^{-1} \boldsymbol{\beta},$$
$$\mathscr{V} = \sigma^2 \widetilde{\nu},$$
$$\mathscr{E} = 2\mu \cdot \boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \mu \boldsymbol{I})^{-1} \boldsymbol{\Sigma}_0 (\boldsymbol{\beta}_0 - \boldsymbol{\beta}),$$
$$\kappa^2 = (\boldsymbol{\beta}_0 - \boldsymbol{\beta})^\top \boldsymbol{\Sigma}_0 (\boldsymbol{\beta}_0 - \boldsymbol{\beta}) + \sigma_0^2.$$

The optimal regularization is defined as

$$\lambda^* \in \operatorname*{argmin}_{\lambda \geq \lambda_{\min}(\phi)} \mathscr{R}(\lambda, \phi). \quad (2)$$

# Optimal regularization sign characterization (IND)

### Theorem (Optimal regularization sign for IND risk)

1. *(Underparameterized) When $\phi < 1$, we have $\lambda^* \geq 0$.*

2. *(Overparameterized) When $\phi > 1$, if for all $v < 1/\mu(0, \phi)$, the following general alignment holds:*

$$\frac{\bar{\text{tr}}[\boldsymbol{B}\boldsymbol{\Sigma}(v\boldsymbol{\Sigma} + \boldsymbol{I})^{-2}] + \sigma^2}{\bar{\text{tr}}[\boldsymbol{B}\boldsymbol{\Sigma}(v\boldsymbol{\Sigma} + \boldsymbol{I})^{-3}] + \sigma^2} > \frac{\bar{\text{tr}}[\boldsymbol{\Sigma}(v\boldsymbol{\Sigma} + \boldsymbol{I})^{-2}]}{\bar{\text{tr}}[\boldsymbol{\Sigma}(v\boldsymbol{\Sigma} + \boldsymbol{I})^{-3}]}, \tag{3}$$

*where $\boldsymbol{B} = \boldsymbol{\beta}\boldsymbol{\beta}^\top$, then we have $\lambda^* < 0$.*

▶ Alignment condition (3) captures how well the signal $\boldsymbol{B}$ is aligned with the feature covariance $\boldsymbol{\Sigma}$.

▶ $\lambda^*$ could be negative in the overparameterized regime when $p > n$.
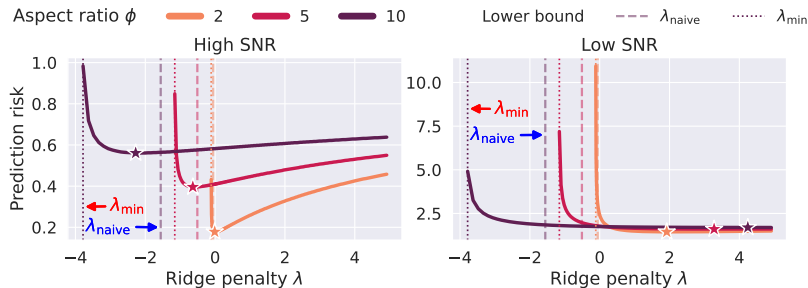
# Illustration (optimal IND regularization)



Figure: Illustration of negative or positive optimal regularization under general alignment.

▶ $\lambda^*$ can be smaller than the previous bound.

▶ The more the alignment (seen as a function of SNR), the lower $\lambda^*$; the more the misalignment, the higher $\lambda^*$ (seen as a function of SNR).

# Optimal regularization sign characterization (OOD, covariate shift)

1. *(Underparameterized)* When $\phi < 1$, we have $\lambda^* \geq 0$.
2. *(Overparameterized)* When $\phi > 1$, if $\boldsymbol{\Sigma}_0 = \boldsymbol{I}$ (corresponding to the estimation risk), then we have $\lambda^* \geq 0$.
3. *(Overparameterized)* When $\phi > 1$, if $\boldsymbol{\Sigma} = \boldsymbol{I}$ and

$$\bar{\mathrm{tr}}[\boldsymbol{\Sigma}_0 \boldsymbol{B}] > \bar{\mathrm{tr}}[\boldsymbol{\Sigma}_0] \left( \bar{\mathrm{tr}}[\boldsymbol{B}] + \frac{(1 + \mu(0, \phi))^3}{\mu(0, \phi)^3} \sigma^2 \right), \tag{4}$$

where $\boldsymbol{B} = \boldsymbol{\beta}\boldsymbol{\beta}^\top$, then we have $\lambda^* < 0$.

▶ The isotropic test covariance case ($\boldsymbol{\Sigma}_0 = \boldsymbol{I}$) is similar to underparameterized cases.

▶ Alignment condition (4) captures how the well the signal $\boldsymbol{B}$ aligned with covariance matrix of test features $\boldsymbol{\Sigma}_0$.

▶ $\lambda^*$ can be negative even in the isotropic train covariance case ($\boldsymbol{\Sigma} = \boldsymbol{I}$).

# Optimal regularization sign characterization (OOD, label shift)

1. *(Underparameterized)* When $\phi < 1$, if $\sigma^2 = o(1)$ and for all $\mu \geq 0$, the following general alignment holds:

$$\bar{\text{tr}}[\boldsymbol{B}_0 \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \mu \boldsymbol{I})^{-2}] > \bar{\text{tr}}[\boldsymbol{B} \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \mu \boldsymbol{I})^{-2}], \tag{5}$$

   where $\boldsymbol{B} = \boldsymbol{\beta}\boldsymbol{\beta}^\top$ and $\boldsymbol{B}_0 = \boldsymbol{\beta}_0\boldsymbol{\beta}^\top$, then we have $\lambda^* < 0$.

2. *(Overparameterized)* When $\phi > 1$, if the general alignment conditions (3) and (5) hold, then we have $\lambda^* < 0$.

▶ $\lambda^*$ can be negative even if the design is underparameterized!

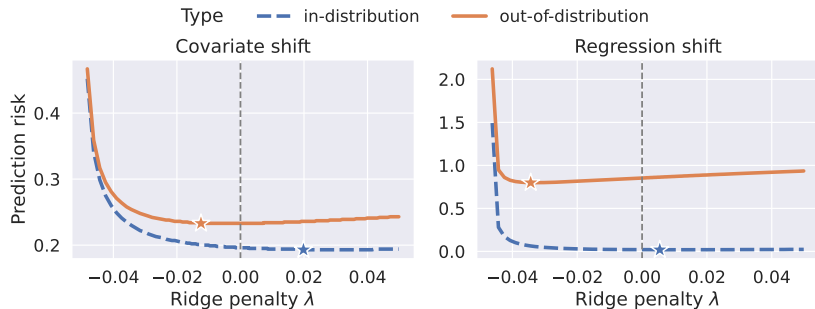# Illustration (optimal OOD regularization)



Figure: Covariate and regression shift can lead to negative optimal regularization in both underparameterized and overparameterized regimes.

The design is isotropic on the left.
The design is underparameterized on the right.

# Optimal risk monotonicity

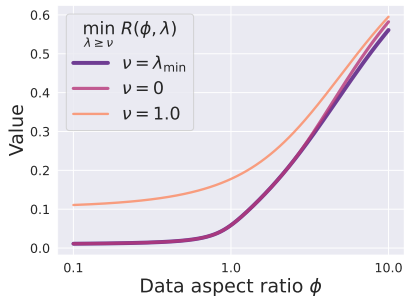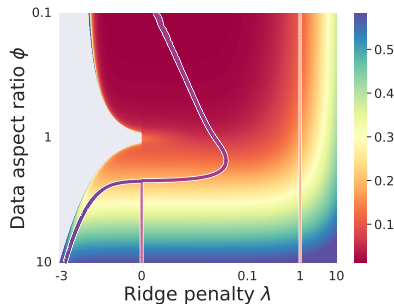The map $\phi \mapsto \min_{\lambda \geq \lambda_{\min}(\phi)} \mathcal{R}(\lambda, \phi)$ is monotonically increasing in $\phi$.



Figure: Ridge regression optimized over $\lambda \geq \nu$ for different thresholds $\nu$ has monotonic risk profile.

Previous result holds for positive $\lambda$ and IND risks.
Current result holds for negative $\lambda$ and OOD risks.