

# Optimal Ridge Regularization for Out-of-Distribution Prediction

Pratik Patil<sup>1</sup>   Jin-Hong Du<sup>2</sup>   Ryan J. Tibshirani<sup>1</sup>

<sup>1</sup>University of California, Berkeley

<sup>2</sup>Carnegie Mellon University

JSM 2024

# Ridge regression in high dimensions

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

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

and its prediction risk:

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

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

$$R(\widehat{\beta}^\lambda) \rightarrow \mathcal{R}(\lambda, \phi)$$

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

## Optimal ridge regression under in-distribution

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

- (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)  $\mathcal{R}(\lambda^*, \phi)$  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  $\mathcal{R}(\lambda^*, \phi)$ ?*

# Summary of results

## Optimal regularization landscape in ridge regression.

| $\Sigma$            | $\beta$   | $\Sigma_0$ | $\beta_0$ | $\phi \leq 1$ | $\lambda_{\min}$ | Arb. Mod.    | Arb. SNR     | Arb. Spec.   | Arb. Geometry   | Additional Specific Data Conditions | $\lambda^*$ | Reference                  |
|---------------------|-----------|------------|-----------|---------------|------------------|--------------|--------------|--------------|---|-------------------------------------|-------------|----------------------------|
| In-distribution     |           |            |           |               |                  |              |              |              |   |                                     |             |                            |
| $\otimes$           | $\circ$   | $\Sigma$   | $\beta$   | all           | zero             | $\times$     | $\checkmark$ | $\times$     |   |                                     | +           | [DW, Thm. 2.1]             |
| $\circ$             | $\otimes$ | $\Sigma$   | $\beta$   | all           | zero             | $\times$     | $\checkmark$ | $\times$     |   |                                     | +           | [HMRT, Cor. 5]             |
|                     |           |            |           | under         | neg              | $\times$     | $\checkmark$ | $\times$     |   |                                     | +           | [WX, Prop. 6]              |
|                     |           |            |           | over          | neg              | $\times$     | $\times$     | $\times$     | Strict misalignment of $(\Sigma, \beta)$                  |                                     | +           | [WX, Thm. 4]               |
|                     |           |            |           | over          | neg              | $\times$     | $\times$     | $\times$     | Strict alignment of $(\Sigma, \beta)$                     |                                     | -           | [WX, Thm. 4, Prop. 7]      |
| $\otimes$           | $\otimes$ | $\Sigma$   | $\beta$   | over          | zero             | $\times$     | $\times$     | $\times$     | and/or special feature model                              |                                     | 0           | [RMR, Cor. 2]              |
|                     |           |            |           | under         | neg*             | $\checkmark$ | $\checkmark$ | $\checkmark$ |   |                                     | +           | <b>Theorem 2 (1)</b>       |
|                     |           |            |           | over          | neg*             | $\checkmark$ | $\checkmark$ | $\checkmark$ | General alignment of $(\Sigma, \beta, \sigma^2)$          |                                     | -           | <b>Theorem 2 (2)</b>       |
| Out-of-distribution |           |            |           |               |                  |              |              |              |   |                                     |             |                            |
| $\otimes$           | $\circ$   | $\Sigma_0$ | $\beta$   | all           | neg*             | $\checkmark$ | $\checkmark$ | $\checkmark$ |   |                                     | +           | <b>Proposition 3</b>       |
| $\otimes$           | $\otimes$ | $\Sigma_0$ | $\beta$   | under         | neg*             | $\checkmark$ | $\checkmark$ | $\checkmark$ |   |                                     | +           | <b>Theorem 4 (1)</b>       |
| $\otimes$           | $\otimes$ | $I$        | $\beta$   | over          | neg*             | $\checkmark$ | $\checkmark$ | $\checkmark$ |   |                                     | +           | <b>Theorem 4 (2)</b>       |
| $\circ$             | $\otimes$ | $\Sigma_0$ | $\beta$   | over          | neg*             | $\checkmark$ | $\checkmark$ | $\checkmark$ | General alignment of $(\Sigma_0, \beta, \sigma^2)$        |                                     | -           | <b>Theorem 4 (3)</b>       |
|                     |           |            |           | under         | neg*             | $\checkmark$ | $\checkmark$ | $\checkmark$ | General alignment of $(\Sigma, \beta, \beta_0)$           |                                     | -           | <b>Theorem 5 (1), (39)</b> |
| $\otimes$           | $\otimes$ | $\Sigma$   | $\beta_0$ | under         | neg*             | $\checkmark$ | $\checkmark$ | $\checkmark$ | General misalignment of $(\Sigma, \beta, \beta_0)$        |                                     | +           | <b>Theorem 5 (1), (39)</b> |
|                     |           |            |           | over          | neg*             | $\checkmark$ | $\checkmark$ | $\checkmark$ | General alignment of $(\Sigma, \beta, \beta_0, \sigma^2)$ |                                     | -           | <b>Theorem 5 (2)</b>       |

# Data assumptions and lower bound on negative regularization

Data assumptions:

- ▶ Covariate: Each feature vector  $\mathbf{x}_i$  for  $i \in [n]$  can be decomposed as  $\mathbf{x}_i = \Sigma^{1/2} \mathbf{z}_i$ , where  $\mathbf{z}_i \in \mathbb{R}^p$  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:

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

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

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

# Out-of-distribution risk characterization

The OOD risk asymptotics read that

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

where

$$\mathcal{B} = \mu^2 \cdot \boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \mu \mathbf{I})^{-1} (\tilde{\mathbf{v}} \boldsymbol{\Sigma} + \boldsymbol{\Sigma}_0) (\boldsymbol{\Sigma} + \mu \mathbf{I})^{-1} \boldsymbol{\beta},$$

$$\mathcal{V} = \sigma^2 \tilde{\mathbf{v}},$$

$$\mathcal{E} = 2\mu \cdot \boldsymbol{\beta}^\top (\boldsymbol{\Sigma} + \mu \mathbf{I})^{-1} \boldsymbol{\Sigma}_0 (\boldsymbol{\beta}_0 - \boldsymbol{\beta}),$$

$$\kappa^2 = (\boldsymbol{\beta}_0 - \boldsymbol{\beta})^\top \boldsymbol{\Sigma}_0 (\boldsymbol{\beta}_0 - \boldsymbol{\beta}) + \sigma_0^2.$$

The optimal regularization is defined as

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

# 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:

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

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

- ▶ **Alignment condition** (3) captures how well the signal  $\mathbf{B}$  is aligned with the feature covariance  $\Sigma$ .
- ▶  $\lambda^*$  could be **negative** in the overparameterized regime when  $p > n$ .



# Illustration (optimal IND regularization)

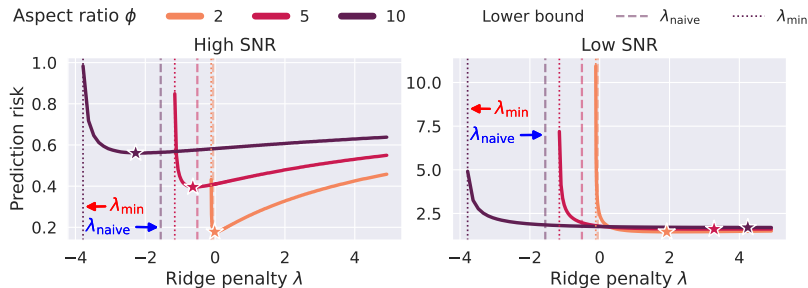


Figure: Illustration of negative or positive optimal regularization under general alignment.

- ▶  $\lambda^*$  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  $\Sigma_0 = \mathbf{I}$  (corresponding to the estimation risk), then we have  $\lambda^* \geq 0$ .
3. (*Overparameterized*) When  $\phi > 1$ , if  $\Sigma = \mathbf{I}$  and

$$\bar{\text{tr}}[\Sigma_0 \mathbf{B}] > \bar{\text{tr}}[\Sigma_0] \left( \bar{\text{tr}}[\mathbf{B}] + \frac{(1 + \mu(0, \phi))^3}{\mu(0, \phi)^3} \sigma^2 \right), \quad (4)$$

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

- ▶ The isotropic test covariance case ( $\Sigma_0 = \mathbf{I}$ ) is similar to underparameterized cases.
- ▶ Alignment condition (4) captures how well the signal  $\mathbf{B}$  is aligned with the covariance matrix of test features  $\Sigma_0$ .
- ▶  $\lambda^*$  can be **negative** even in the isotropic train covariance case ( $\Sigma = \mathbf{I}$ ).

## 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:

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

where  $\mathbf{B} = \beta \beta^\top$  and  $\mathbf{B}_0 = \beta_0 \beta_0^\top$ , 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)

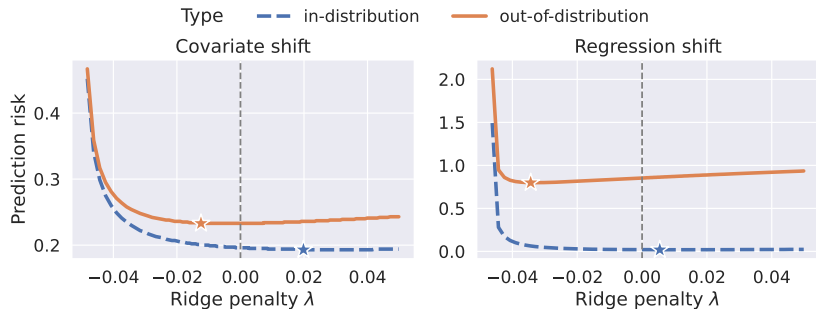


Figure: Covariate and regression shift can lead to negative optimal regularization in both underparameterized and overparameterized regimes.

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  $\phi$ .

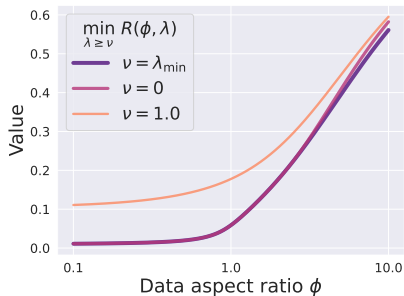
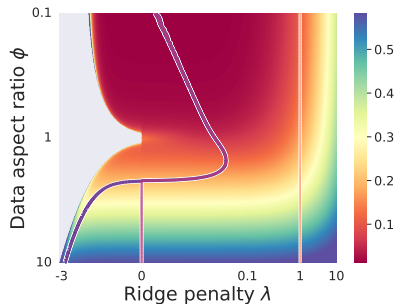


Figure: Ridge regression optimized over  $\lambda \geq \nu$  for different thresholds  $\nu$  has monotonic risk profile.

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