

# 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>Department of Statistics, UCB <sup>2</sup>Department of Statistics, CMU

## Ridge regression in high dimensions

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

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

and its prediction risk:

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

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

$$R(\hat{\beta}^\lambda) \rightarrow \mathcal{R}(\lambda, \phi), \quad p/n \rightarrow \phi \in (0, \infty)$$

where  $p$  is feature size,  $n$  is sample size, and  $\phi$  is the *aspect ratio*.

**Distribution shifts.** We consider two types of distribution shifts:

(1) *Covariate shift*: where  $P_{\mathbf{x}_0} \neq P_{\mathbf{x}}$  but  $P_{y_0|\mathbf{x}_0} = P_{y|\mathbf{x}}$ .

(2) *Regression shift*: where  $P_{y_0|\mathbf{x}_0} \neq P_{y|\mathbf{x}}$  but  $P_{\mathbf{x}_0} = P_{\mathbf{x}}$ .

**Questions of interest.** We answer two out-of-distribution problems:

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

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

**Data assumptions.** Feature distribution: 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  $4^+$  moments for some  $\mu > 0$ . Response distribution: Each response variable  $y_i$  for  $i \in [n]$  has mean 0, and bounded  $4^+$  moments.

**Lower bound on ridge regularization.** 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:  $\lambda_{\min}(\phi) = \mu_{\min} - \phi \bar{\text{tr}}[\Sigma(\Sigma + \mu_{\min} \mathbf{I})^{-1}]$ .

## Summary of results

$\Sigma$	$\beta$	$\Sigma_0$	$\beta_0$	$\phi \leq 1$	$\lambda_{\min}$	Arb. Mod.	Arb. SNR	Arb. Spec.	Additional Geometry	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]
				over	zero	$\times$	$\times$	$\times$	and/or special feature model		0	[RMR, Cor. 2]
				under	neg*	$\checkmark$	$\checkmark$	$\checkmark$			+	Theorem 2 (1)
				over	neg*	$\checkmark$	$\checkmark$	$\checkmark$	General alignment of $(\Sigma, \beta, \sigma^2)$		-	Theorem 2 (2)
Out-of-distribution												
$\otimes$	$\circ$	$\Sigma_0$	$\beta$	all	neg*	$\checkmark$	$\checkmark$	$\checkmark$			+	Proposition 3
$\otimes$	$\otimes$	$\Sigma_0$	$\beta$	under	neg*	$\checkmark$	$\checkmark$	$\checkmark$			+	Theorem 4 (1)
$\otimes$	$\otimes$	$\mathbf{I}$	$\beta$	over	neg*	$\checkmark$	$\checkmark$	$\checkmark$			+	Theorem 4 (2)
$\circ$	$\otimes$	$\Sigma_0$	$\beta$	over	neg*	$\checkmark$	$\checkmark$	$\checkmark$	General alignment of $(\Sigma_0, \beta, \sigma^2)$		-	Theorem 4 (3)
				under	neg*	$\checkmark$	$\checkmark$	$\checkmark$	General alignment of $(\Sigma, \beta, \beta_0)$		-	Theorem 5 (1), (39)
				under	neg*	$\checkmark$	$\checkmark$	$\checkmark$	General misalignment of $(\Sigma, \beta, \beta_0)$		+	Theorem 5 (1), (39)
				over	neg*	$\checkmark$	$\checkmark$	$\checkmark$	General alignment of $(\Sigma, \beta, \beta_0, \sigma^2)$		-	Theorem 5 (2)

## Out-of-distribution risk characterization

**Proposition 1 (Deterministic equivalents for OOD risk).** The asymptotic OOD risk decomposes into:

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

where

$$\mathcal{B} = \mu^2 \cdot \beta^\top (\Sigma + \mu \mathbf{I})^{-1} (\tilde{v} \Sigma + \Sigma_0) (\Sigma + \mu \mathbf{I})^{-1} \beta,$$

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

$$\mathcal{E} = 2\mu \cdot \beta^\top (\Sigma + \mu \mathbf{I})^{-1} \Sigma_0 (\beta_0 - \beta),$$

$$\kappa^2 = (\beta_0 - \beta)^\top \Sigma_0 (\beta_0 - \beta) + \sigma_0^2.$$

The optimal regularization is defined as  $\lambda^* \in \arg \min_{\lambda \geq \lambda_{\min}(\phi)} \mathcal{R}(\lambda, \phi)$ .

## Optimal regularization sign characterization (IND)

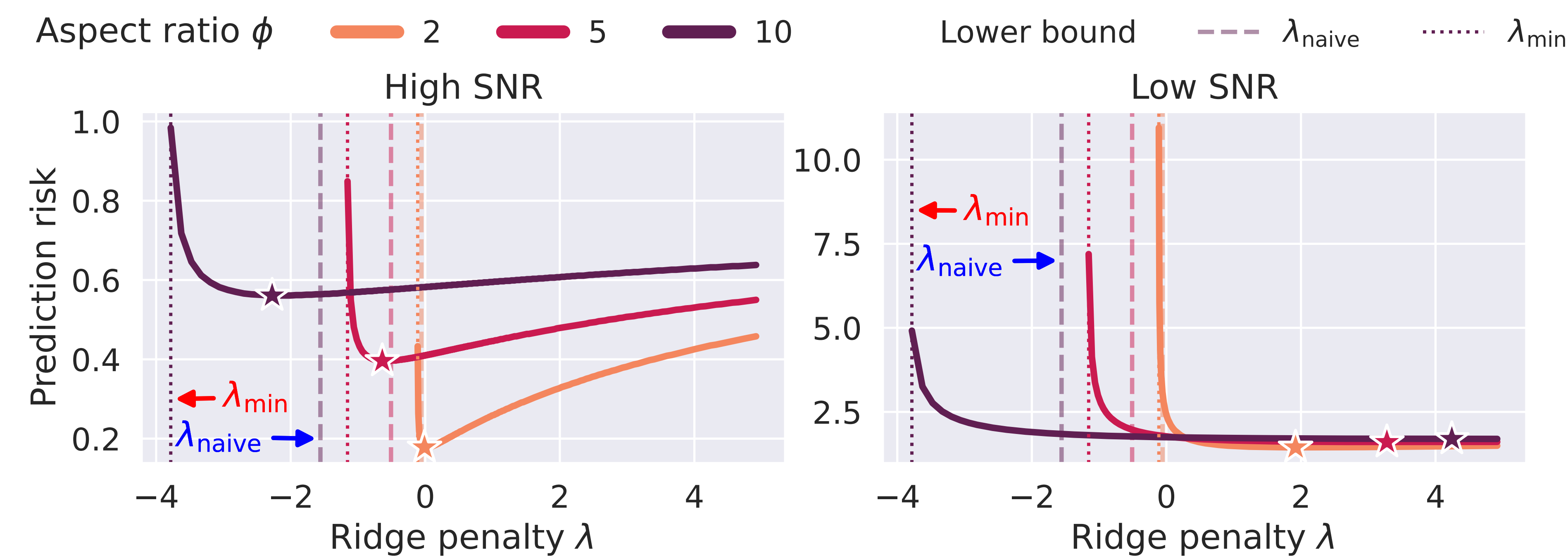


Illustration of negative or positive optimal regularization under general alignment.

**Theorem 2 (Optimal regularization sign, no shift)** Assume  $\Sigma_0 = \Sigma$  and  $\beta_0 = \beta$ .

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:

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

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

### References:

[HMRT] Trevor Hastie et al. "Surprises in high-dimensional ridgeless least squares interpolation". In: *The Annals of Statistics* 50.2 (2022), pp. 949–986

[DW] Edgar Dobriban and Stefan Wager. "High-dimensional asymptotics of prediction: Ridge regression and classification". In: *The Annals of Statistics* 46.1 (2018), pp. 247–279

[RMR] Dominic Richards, Jaouad Mourtada, and Lorenzo Rosasco. "Asymptotics of ridge (less) regression under general source condition". In: *International Conference on Artificial Intelligence and Statistics*. 2021

[WX] Denny Wu and Ji Xu. "On the Optimal Weighted  $\ell_2$  Regularization in Overparameterized Linear Regression". In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 10112–10123

## Optimal regularization sign characterization (OOD)

**Theorem 4 (Optimal regularization, covariate shift).** Assume  $\Sigma_0 \neq \Sigma$  and  $\beta_0 = \beta$ .

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

2. (*Overparameterized*) When  $\phi > 1$ , if  $\Sigma_0 = \mathbf{I}$  (estimation risk), we have  $\lambda^* \geq 0$ .

3. (*Overparameterized*) When  $\phi > 1$ , if  $\Sigma = \mathbf{I}$  and

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

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

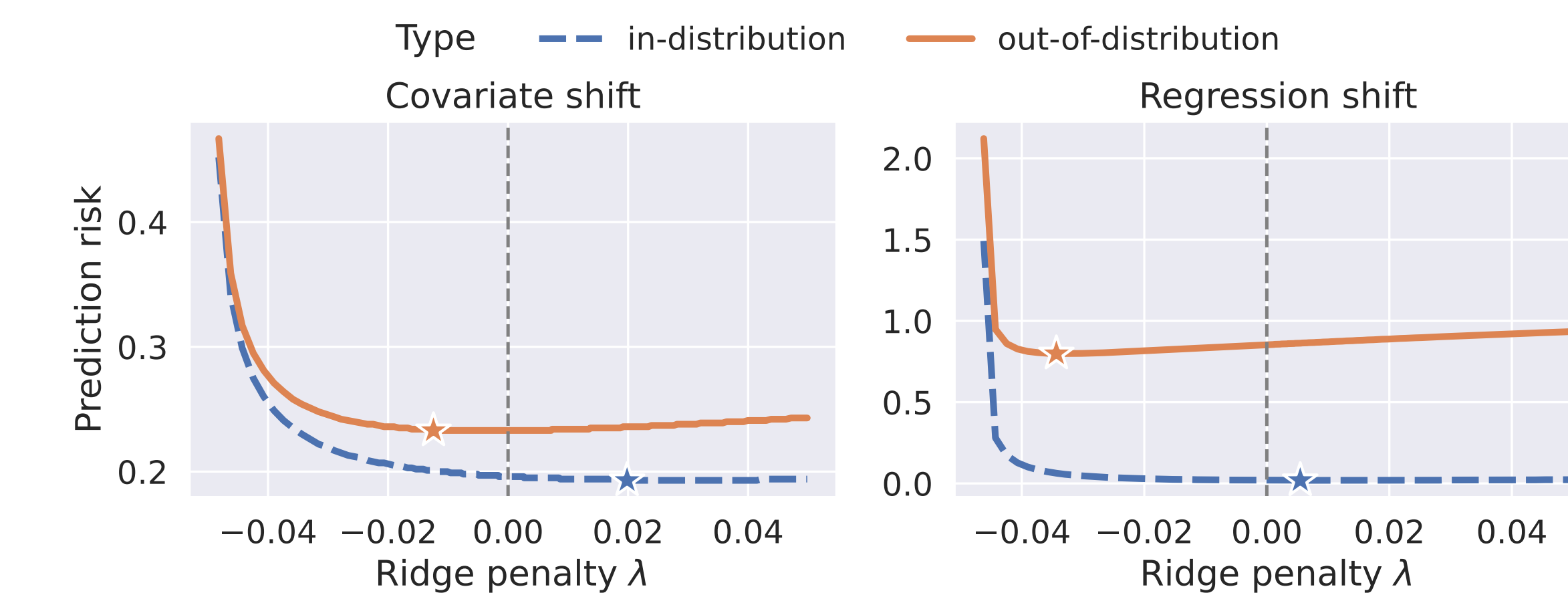
**Theorem 5 (Optimal regularization, regression shift).** Assume  $\Sigma_0 = \Sigma$  and  $\beta_0 \neq \beta$ .

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

$$\bar{\text{tr}}[\mathbf{B}_0 \Sigma^2 (\Sigma + \mu \mathbf{I})^{-2}] > \bar{\text{tr}}[\mathbf{B} \Sigma^2 (\Sigma + \mu \mathbf{I})^{-2}], \quad (4)$$

where  $\mathbf{B} = \beta\beta^\top$  and  $\mathbf{B}_0 = \beta_0\beta_0^\top$ , we have  $\lambda^* < 0$ .

2. (*Overparameterized*) When  $\phi > 1$ , if conditions (2) and (4) hold, we have  $\lambda^* < 0$ .



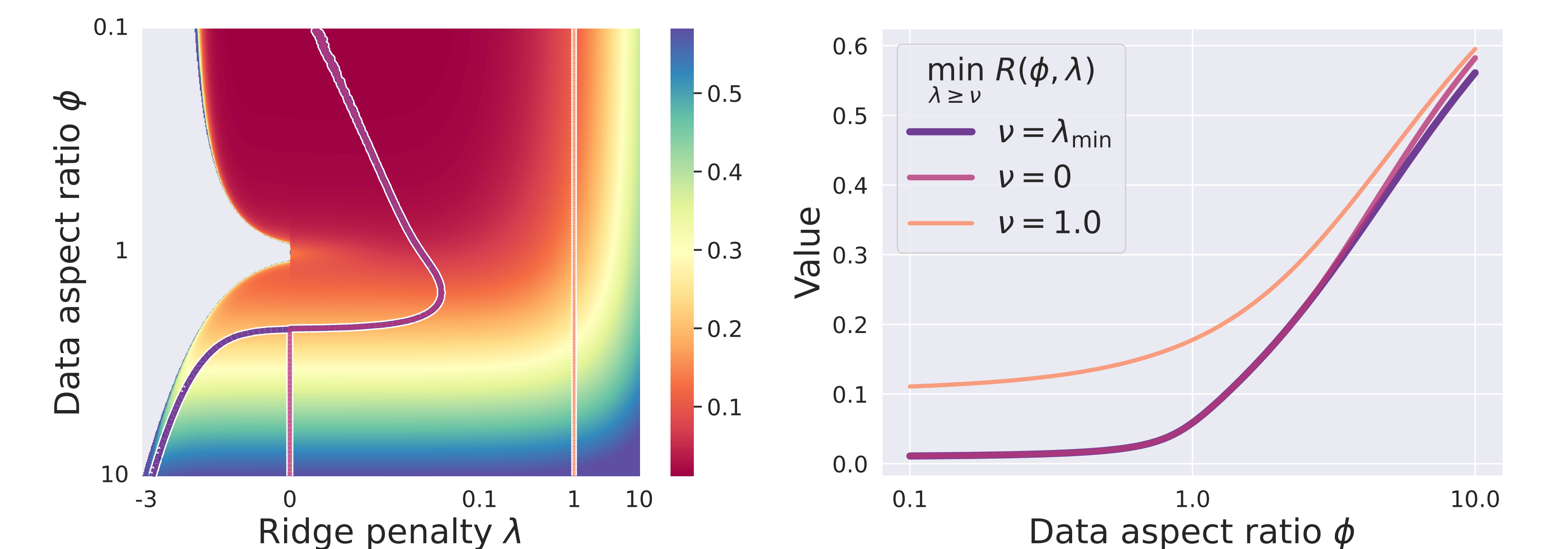
Covariate and regression shift can lead to negative optimal regularization in both overparameterized and underparameterized regimes.

## Optimal risk monotonicity (both IND and OOD)

**Theorem 7 (Optimal regularization, regression shift).** For  $\lambda \geq \lambda_{\min}(\phi)$ , for all  $\varepsilon > 0$  small enough, the risk of optimal ridge predictor satisfies:

$$\min_{\lambda \geq \lambda_{\min}(\phi) + \varepsilon} R(\hat{\beta}^\lambda) \simeq \min_{\lambda \geq \lambda_{\min}(\phi)} \mathcal{R}(\lambda, \phi), \quad (5)$$

and right side of (5) is monotonically increasing in  $\phi$  if SNR and  $\sigma_0^2$  are fixed. In addition, when  $\beta = \beta_0$  it is monotonically increasing in SNR if  $\phi$ ,  $\sigma^2$ , and  $\sigma_0^2$  are fixed.



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