

PS5841

Data Science in Finance & Insurance

Shrinkage

Yubo Wang

Spring 2022

Ridge Regression

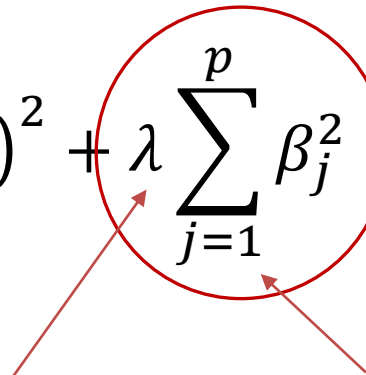
- Model

$$Y_i = \beta_0 + \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i$$

- Least squares: $\hat{\boldsymbol{\beta}}^{LS}$ minimizes

$$R_{LS} = ESS_{LS} = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \cdots - \beta_p x_{ip})^2$$

- Ridge regression: $\hat{\boldsymbol{\beta}}_{\lambda}^R$ minimizes

$$R(\lambda) = \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^T \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^p \beta_j^2$$


tuning parameter λ

no intercept

Solution

- Loss, where $\Lambda = \text{diag}(0, \lambda_1, \dots, \lambda_p)$

$$R(\lambda) = (\mathbf{y} - X\boldsymbol{\beta})^T (\mathbf{y} - X\boldsymbol{\beta}) + \boldsymbol{\beta}^T \Lambda \boldsymbol{\beta}$$

$$\frac{\partial R}{\partial \boldsymbol{\beta}} = -2X^T (\mathbf{y} - X\boldsymbol{\beta}) + 2\Lambda \boldsymbol{\beta}$$

$$\hat{\boldsymbol{\beta}}^R = (X^T X + \Lambda)^{-1} X^T \mathbf{y}$$

- Biased (when $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$, $\text{var}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}$)

$$E(\hat{\boldsymbol{\beta}}^R) = (X^T X + \Lambda)^{-1} X^T X \boldsymbol{\beta}$$

$$\text{Var}(\hat{\boldsymbol{\beta}}^R) = \sigma^2 (X^T X + \Lambda)^{-1} X^T X (X^T X + \Lambda)^{-1}$$

Equivalent Solution

- When features are centered, $\hat{\boldsymbol{\beta}}_{\lambda}^R$ minimizes

$$R(\lambda) = (\mathbf{y} - X\boldsymbol{\beta})^T (\mathbf{y} - X\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta}$$

$$\frac{\partial R}{\partial \boldsymbol{\beta}} = -2X^T (\mathbf{y} - X\boldsymbol{\beta}) + 2\lambda \boldsymbol{\beta}$$

$$\hat{\boldsymbol{\beta}}^R = (X^T X + \lambda \mathbf{I})^{-1} X^T \mathbf{y}$$

$$\hat{\beta}_0 = \sum_{i=1}^n y_i$$

no intercept

Recipe

- Scale features
 - Centered and normalized
 - Standardized
- Estimate the intercept with OLS
 - The intercept is a measure of the mean of the response when features are zero
 - When features are centered, we have $\hat{\beta}_0 = \bar{y}$
- Estimate the remaining coefficients by a Ridge regression without intercept using the scaled features

Penalty Perspective

- Minimize

$$R(\lambda) = (\mathbf{y} - X\boldsymbol{\beta})^T (\mathbf{y} - X\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta}$$

- Equivalent formulation

$$\text{minimize } R(\lambda = 0) = (\mathbf{y} - X\boldsymbol{\beta})^T (\mathbf{y} - X\boldsymbol{\beta})$$

$$\text{Subject to } \sum_{j=1}^p \beta_j^2 \leq s^2$$

- Tuning parameter λ
- Shrinks $\boldsymbol{\beta}$ and introduces bias

Bayesian Perspective

$$\mathbf{y} \sim \mathcal{N}(X\boldsymbol{\beta}, \sigma^2 \mathbf{I}), \quad \boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I}) \text{ prior}$$

- Posterior log-likelihood

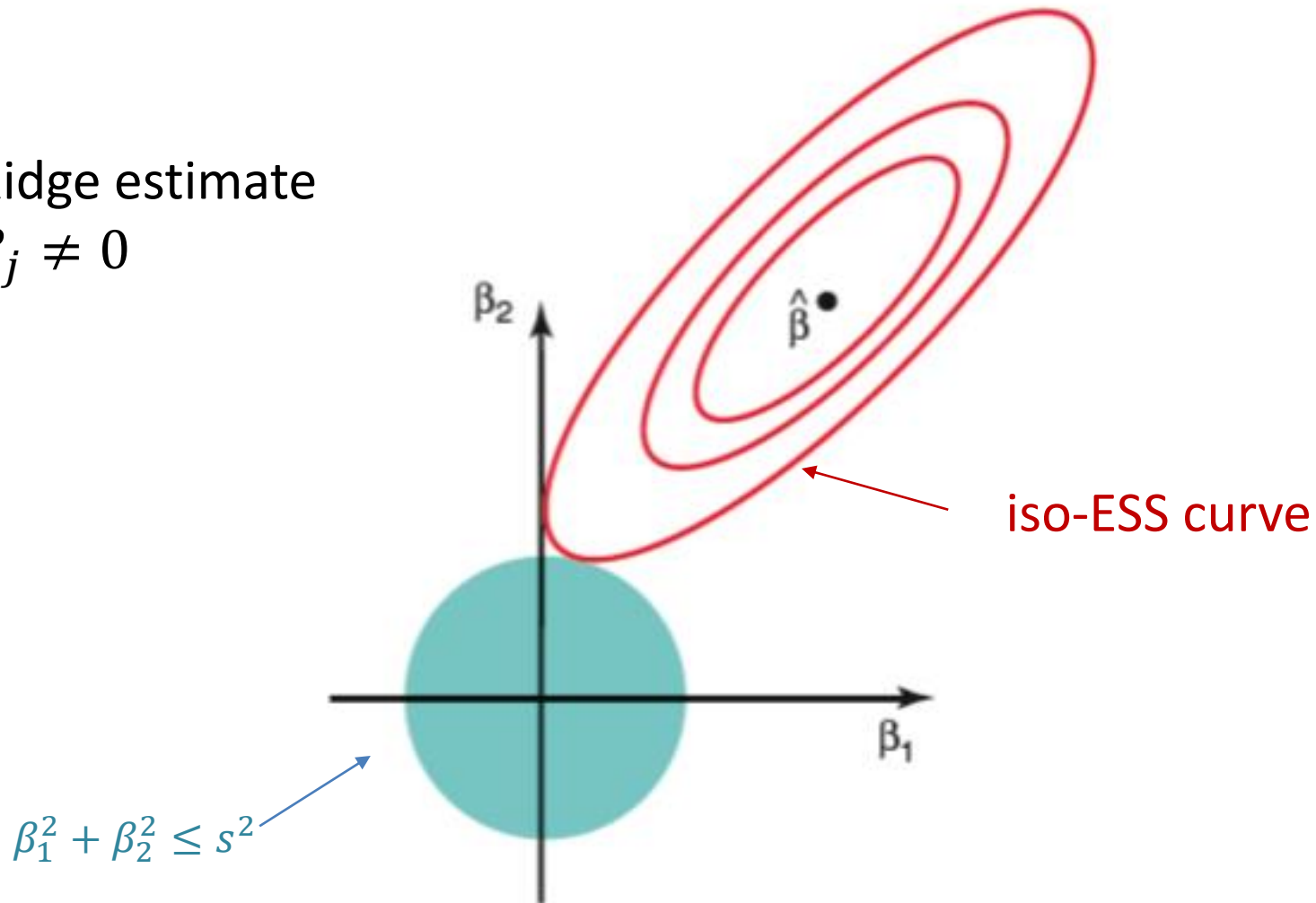
$$\begin{aligned} l(\boldsymbol{\beta}; \mathbf{y}) &\propto \text{posterior} \quad \text{likelihood} \quad \text{prior} \\ l(\boldsymbol{\beta}; \mathbf{y}) &\propto l(\mathbf{y}; \boldsymbol{\beta}) + l(\boldsymbol{\beta}) \\ &= -\frac{1}{2\sigma^2} [(\mathbf{y} - X\boldsymbol{\beta})^T (\mathbf{y} - X\boldsymbol{\beta}) + \frac{\sigma^2}{\tau^2} \boldsymbol{\beta}^T \boldsymbol{\beta}] \end{aligned}$$

- Minimize $R(\lambda)$, $\lambda = \frac{\sigma^2}{\tau^2}$
- $\hat{\boldsymbol{\beta}}^R$ maximizes the posterior

Geometry Perspective

Ridge estimate

$$\beta_j \neq 0$$



LASSO Regression

- Model

$$Y_i = \beta_0 + \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i$$

- Least squares: $\hat{\boldsymbol{\beta}}^{LS}$ minimizes

$$R_{LS} = ESS_{LS} = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \cdots - \beta_p x_{ip})^2$$

- Ridge regression: $\hat{\boldsymbol{\beta}}_{\lambda}^R$ minimizes

$$R(\lambda) = \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^T \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^p |\beta_j|$$

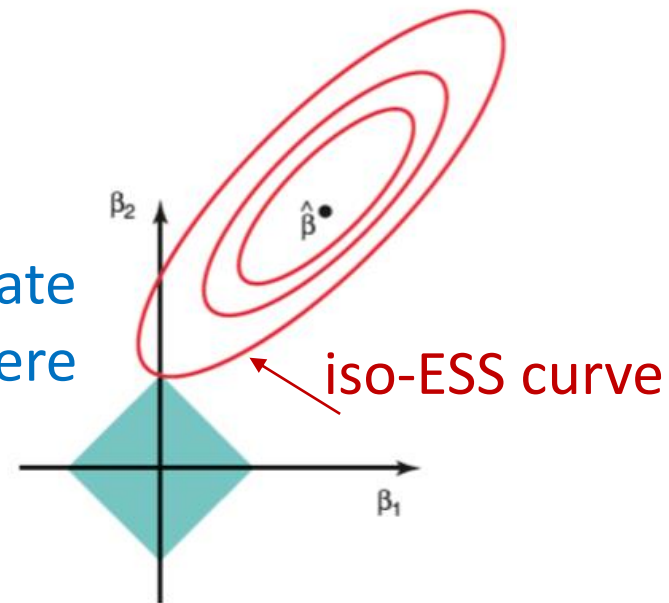
Generalized Penalty

- Minimize

$$ESS(\lambda) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_{j=1}^p |\beta_j|^q$$

- Ridge regression $q = 2$
- LASSO $q = 1$

LASSO estimate
 $\beta_1 = 0$!!! here



Bayesian Perspective

$$\mathbf{y} \sim \mathcal{N}(X\boldsymbol{\beta}, \sigma^2 \mathbf{I}), \quad \boldsymbol{\beta} \sim \text{Laplace}(\mathbf{0}, 2\tau^2 \mathbf{I})$$

prior

- Posterior log-likelihood

posterior

likelihood

prior

$$l(\boldsymbol{\beta}; \mathbf{y}) \propto l(\mathbf{y}; \boldsymbol{\beta}) + l(\boldsymbol{\beta})$$

$$= -\frac{1}{2\sigma^2} [(\mathbf{y} - X\boldsymbol{\beta})^T (\mathbf{y} - X\boldsymbol{\beta}) + \frac{\sigma^2}{\tau} \|\boldsymbol{\beta}\|]$$

- Minimize $R(\lambda)$, $\lambda = \frac{\sigma^2}{\tau}$
- $\hat{\boldsymbol{\beta}}^L$ maximizes the posterior

Observations

- $\hat{\beta}^{LS} = \hat{\beta}_{\lambda=0}^R = \hat{\beta}_{\lambda=0}^L$
- The LASSO can perform variable selection since it can yield **sparse models**
- Ridge does better when the response is a function of many predictors with coefficients of roughly equal size
- LASSO does better when only a small portion of predictors have substantial coefficients

That was

