

PS5841

Data Science in Finance & Insurance

Regularization

Yubo Wang

Autumn 2022

Regularization

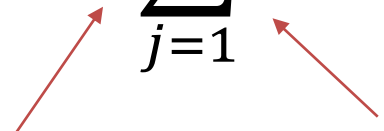
- Variance reduction
- Loss with penalty

Ridge Regression

- Model

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

- 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, $\hat{\beta}_\lambda^R$ minimizes

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

no intercept

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

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

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

Penalty Perspective

- Minimize

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

- Equivalent formulation

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

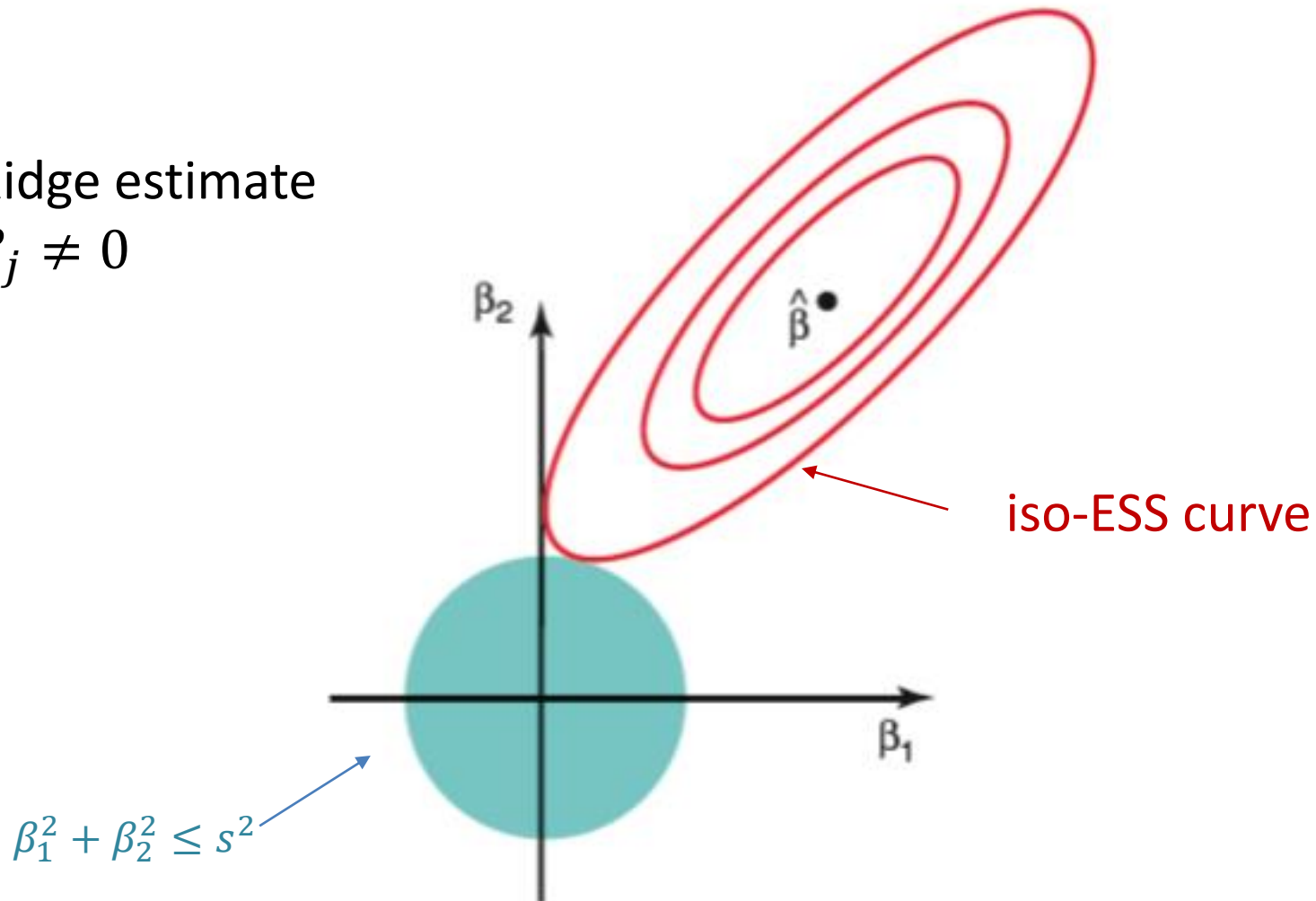
$$\text{Subject to } \sum_{j=1}^p \beta_j^2 \leq s^2$$

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

Geometry Perspective

Ridge estimate

$$\beta_j \neq 0$$



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

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

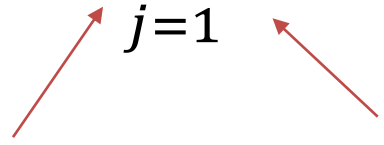
- Minimize $R(\lambda)$, $\lambda = \frac{\sigma^2}{\tau^2}$
- $\hat{\boldsymbol{\beta}}^R$ maximizes the posterior

LASSO

- Model

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

- Least Absolute Shrinkage and Selection Operator (LASSO) $\hat{\boldsymbol{\beta}}_\lambda^L$ minimizes

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


tuning parameter λ

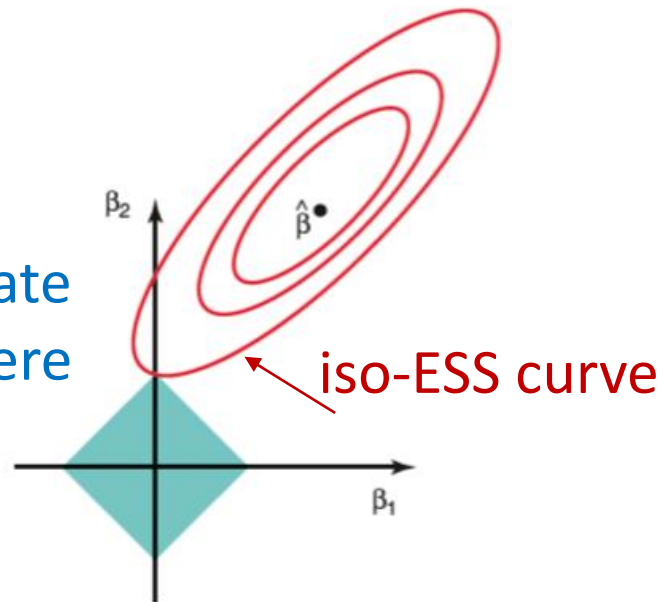
no intercept

Penalty & Geometry Perspective

- Minimize

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

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

Decision Tree Pruning

- A “big” tree risks overfitting data

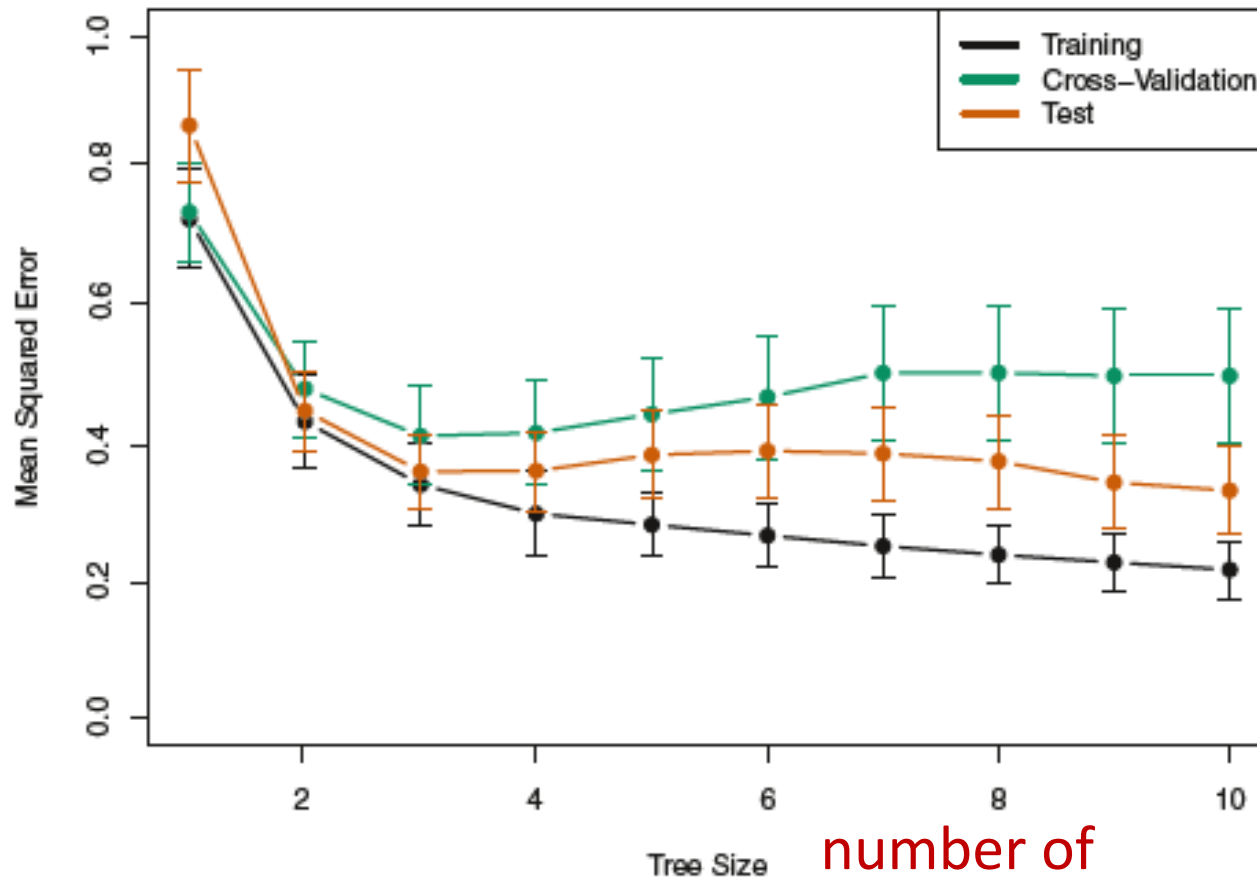
Cost Complexity Pruning

- Loss

$$\text{error_rate} + \alpha|T|$$
$$MSE + \alpha|T|$$

- tuning parameter ($\alpha \geq 0$)
- Resulting tree is a subtree $T \subset T_0$ which minimizes the loss

Example: Pruning



number of
terminal nodes

That was

