

# Penalized Regression

Vahid Partovi Nia

Advanced Machine Learning: Lecture 01

February 10, 2018



LS

MSE

Overfitting

Regularization

① LS

② MSE

③ Overfitting

④ Regularization



# Numerical Consideration

LS

MSE

Overfitting

Regularization

- Cholesky is faster than LU
- LU is faster than QR
- QR is faster than SVD



LS

MSE

Overfitting

Regularization

$$S(\beta) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta)$$

$$\min S(\beta) \Rightarrow (\mathbf{X}^\top \mathbf{X})\beta = \mathbf{X}^\top \mathbf{y}$$

Suppose  $\mathbf{Q}$  is an orthogonal (rotation) matrix. Then

$$S(\beta) = (\mathbf{Q}\mathbf{y} - \mathbf{Q}\mathbf{X}\beta)^\top (\mathbf{Q}\mathbf{y} - \mathbf{Q}\mathbf{X}\beta)$$


LS

MSE

Overfitting

Regularization

- 1 Decompose  $QX = \begin{pmatrix} R \\ 0 \end{pmatrix}$ , where  $R$  is upper triangular.
- 2 Partition  $Qy = \begin{pmatrix} q_1 \\ q_2 \end{pmatrix}$
- 3 Back-solve  $R\beta = q_1$



LS

MSE

Overfitting

Regularization

Suppose  $\mathbf{X}^\top \mathbf{X} = \mathbf{LU}$ ,  $\mathbf{L}$  is lower triangular and  $\mathbf{U}$  is upper triangular.

- ① Decompose  $\mathbf{X}^\top \mathbf{X} = \mathbf{LU}$
- ② Back-solve  $\mathbf{U}\mathbf{q}_1 = \mathbf{X}^\top \mathbf{y}$
- ③ Back-solve  $\mathbf{L}\boldsymbol{\beta}_2 = \mathbf{q}_1$



LS

MSE

Overfitting

Regularization

Suppose  $\mathbf{X}^\top \mathbf{X} = \mathbf{A}^\top \mathbf{A}$ , where  $\mathbf{A}$  is lower triangular.

- ① Decompose  $\mathbf{X}^\top \mathbf{X} = \mathbf{A}^\top \mathbf{A}$
- ② Back-solve  $\mathbf{A} \mathbf{q}_1 = \mathbf{X}^\top \mathbf{y}$
- ③ Back-solve  $\mathbf{A}^\top \boldsymbol{\beta} = \mathbf{q}_1$



LS

MSE

Overfitting

Regularization

Suppose  $\mathbf{X} = \mathbf{PDQ}$ , where  $\mathbf{D}$  is diagonal-zero and  $\mathbf{Q}$  is orthogonal. This means  $\mathbf{X}^\top \mathbf{X}$  is symmetric.

- ① Decompose  $\mathbf{X}^\top \mathbf{X} = \mathbf{Q}^\top \mathbf{D} \mathbf{Q}$ , where  $\mathbf{D}$  is diagonal, and  $\mathbf{Q}$  is orthogonal.
- ②  $\beta = \mathbf{Q}^\top \mathbf{D}^{-1} \mathbf{Q} \mathbf{X}^\top \mathbf{y}$





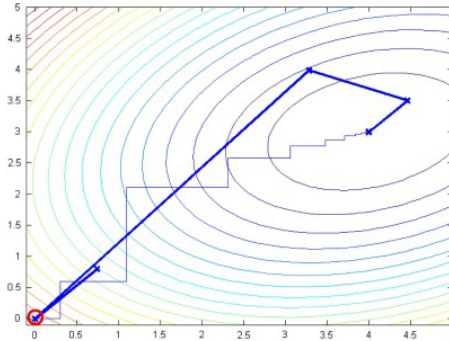
# coordinate vs conjugate

LS

MSE

Overfitting

Regularization



$$\hat{\beta}_j^{t+1} = \operatorname{argmin}_{\beta_j} S \left\{ (\hat{\beta}_{j-1}^t, \beta_j, \hat{\beta}_{j+1}^t) \right\}$$

$$\hat{\beta}^{t+1} = \hat{\beta}^t - \delta \left\{ \frac{\partial^2 S(\beta)}{\partial \beta \partial \beta^\top} \right\}^{-1} \Big|_{\beta=\hat{\beta}^t} \frac{\partial S(\beta)}{\partial \beta} \Big|_{\beta=\hat{\beta}^t}$$



LS

MSE

Overfitting

Regularization

Remember in  $y_i = \beta_1 x_{i1}$  the LS estimator is Suppose

$y_i = y_i - \bar{y}$ , so there is no need for  $\hat{\beta}_0$ .

$$\operatorname{argmin} \frac{1}{2} \sum_i (y_i - x_{i1}\beta_1 - \dots - x_{ij}\beta_j - \dots - x_{ip}\beta_p)^2$$



LS

MSE

Overfitting

Regularization

Remember in  $y_i = \beta_1 x_{i1}$  the LS estimator is Suppose

$y_i = y_i - \bar{y}$ , so there is no need for  $\hat{\beta}_0$ .

$$\operatorname{argmin} \frac{1}{2} \sum_i (y_i - x_{i1}\beta_1 - \dots - x_{ij}\beta_j - \dots - x_{ip}\beta_p)^2$$

$$\hat{\beta}_j = \frac{\sum_{i=1}^n x_{ij} r_i^t}{\sum_{i=1}^n x_{ij}^2}, \quad r_i^t = ?$$



What is the difference between a predictor  $\hat{\mathbf{y}}$  and an estimator  $\hat{\boldsymbol{\beta}}$ ?

- $\text{MSE}(\hat{\mathbf{y}}) = \mathbb{E}(\hat{\mathbf{y}} - \mathbf{y})^\top (\hat{\mathbf{y}} - \mathbf{y}) \approx \text{cross-validation.}$   
Overfitting!
- $\text{MSE}(\hat{\boldsymbol{\beta}}) = \mathbb{E}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \approx \text{bootstrap}$



LS

MSE

Overfitting

Regularization

Suppose a new  $\mathbf{x}_0$  is not observed in  $\mathbf{x}_i, i \in \{1, \dots, n\}$   
 $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i, \mathbb{E}(\varepsilon_i) = 0, \mathbb{V}(\varepsilon_i) = \sigma^2$

$$\text{MSE}\{\hat{y}(\mathbf{x}_0)\} = \mathbb{V}\{\hat{y}(\mathbf{x}_0)\} + \text{bias}^2\{\hat{y}(\mathbf{x}_0)\} + \sigma^2$$

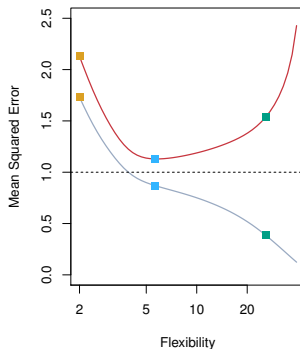
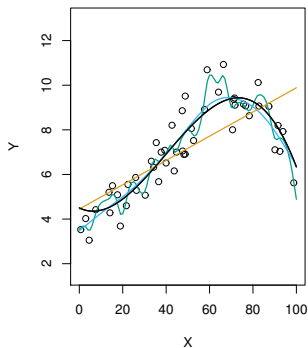


LS

MSE

Overfitting

Regularization



regplot code



Polynomial regression of  $y$  over one  $x$

$$y_i - \bar{y} = \beta_1 x_{i1} + \dots + \beta_k x_i^k + \varepsilon_i$$

- Why condition number of  $(\mathbf{X}^\top \mathbf{X})$  is important?
- What if the smallest eigenvalue of  $\mathbf{X}^\top \mathbf{X}$ , say  $\lambda_p \approx 0$ ?
- When does this happen in polynomial regression?
- How to find LS if the polynomial order  $k < n$ ?



## Polynomial regression of $y$ over one $x$

$$y_i - \bar{y} = \beta_1 x_{i1} + \dots + \beta_k x_i^k + \varepsilon_i$$

- Why condition number of  $(\mathbf{X}^\top \mathbf{X})$  is important?
- What if the smallest eigenvalue of  $\mathbf{X}^\top \mathbf{X}$ , say  $\lambda_p \approx 0$ ?
- When does this happen in polynomial regression?
- How to find LS if the polynomial order  $k < n$ ?
- How to find LS if the polynomial order  $k \geq n$ ?





LS

MSE

Overfitting

Regularization

Suppose the constant model  $y_i = \beta_0 + \varepsilon_i$

$$\hat{\beta}_0 = c\bar{y}$$

$$\hat{c} = \operatorname{argmin} \operatorname{MSE}(c)$$

$$\operatorname{MSE}(c) = \mathbb{E}(\hat{\beta}_0 - \beta_0)^2 \Rightarrow c = ?$$



LS

MSE

Overfitting

Regularization

Suppose multivariate mean estimation problem of dimension  $p \geq 3$  while

$$\mathbf{y}_{p \times 1} \sim \mathcal{N}(\boldsymbol{\beta}_{p \times 1}, \mathbf{I}_{p \times p}).$$

Stein showed

$$\left\{ 1 - \frac{(d-2)}{\sum_{j=1}^p y_j^2} \right\} \mathbf{y}$$

estimates  $\boldsymbol{\beta}$  better than  $\mathbf{y}$  in terms of MSE



LS

MSE

Overfitting

Regularization

For a known penalization constant  $\lambda$ , find

$$\text{Ridge : } \hat{\beta}_0(\lambda) = \operatorname{argmin} \frac{1}{2} \sum_{i=1}^n (y_i - \beta_0)^2 + \lambda \beta_0^2$$



LS

MSE

Overfitting

Regularization

For a known penalization constant  $\lambda$ , find

$$\text{Ridge : } \hat{\beta}_0(\lambda) = \operatorname{argmin} \frac{1}{2} \sum_{i=1}^n (y_i - \beta_0)^2 + \lambda \beta_0^2$$

$$\hat{\beta}_0(\lambda) = \frac{n}{n + \lambda} \bar{y}$$



For a known penalization constant  $\lambda$ , find

$$\text{Ridge : } \hat{\beta}_0(\lambda) = \operatorname{argmin} \frac{1}{2} \sum_{i=1}^n (y_i - \beta_0)^2 + \lambda \beta_0^2$$

$$\hat{\beta}_0(\lambda) = \frac{n}{n + \lambda} \bar{y}$$

$$\text{Lasso : } \hat{\beta}_0(\lambda) = \operatorname{argmin} \frac{1}{2} \sum_{i=1}^n (y_i - \beta_0)^2 + \lambda |\beta_0|$$

$$\hat{\beta}_0(\lambda) = \operatorname{sign}(\bar{y}) \left\{ \bar{y} - \frac{\lambda}{n} \right\}_+$$



# Plot $(\bar{y}, \hat{\beta}_0(\lambda))$

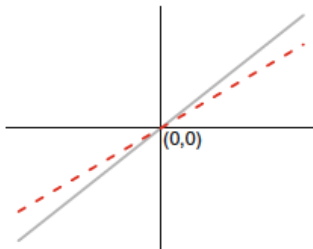
LS

MSE

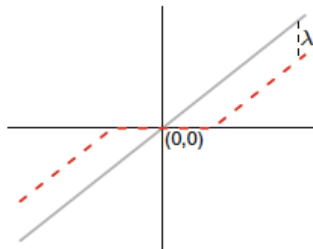
Overfitting

Regularization

Ridge



Lasso



LS

MSE

Overfitting

Regularization

$$S(\boldsymbol{\beta}) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}$$

Show that

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin} S(\boldsymbol{\beta}) = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

How do you compute?



LS

MSE

Overfitting

Regularization

$$S(\beta) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^\top \beta$$

Show that

$$\hat{\beta} = \operatorname{argmin} S(\beta) = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

How do you compute?

using Cholesky? LU? QR? SVD? coordinate? conjugate?





LS

MSE

Overfitting

Regularization

$$S(\boldsymbol{\beta}) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_{j=1}^p |\beta_j|$$

How do you compute?



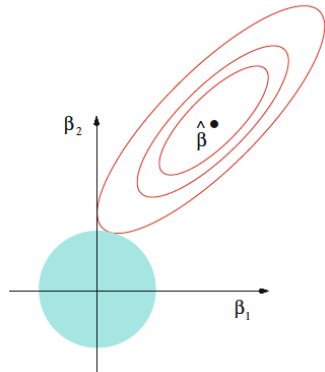
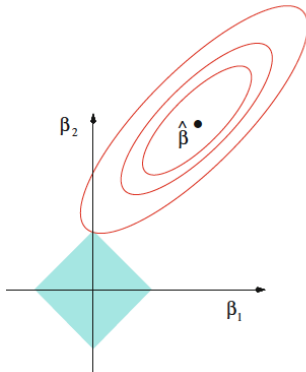
# Visual penalization

LS

MSE

Overfitting

Regularization



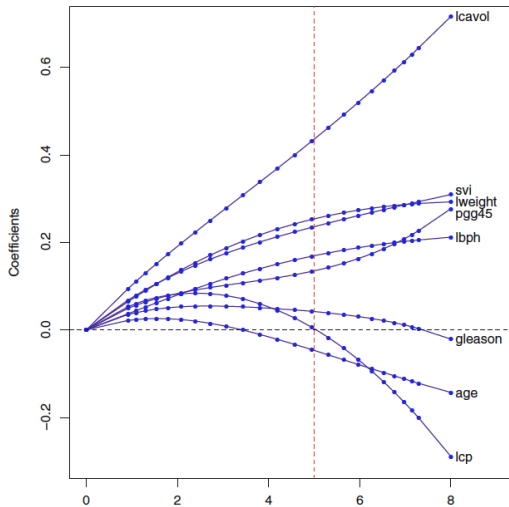
# Ridge coefficients

LS

MSE

Overfitting

Regularization



# Lasso coefficients

LS

MSE

Overfitting

Regularization

