

Ridge Regression, a.k.a $\hat{\beta}^{Ridge} = \arg \min \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$

Andersen Ang

First created: 2014. Last update : 2017-Feb

1 Summary

Consider the regression problem

$$y_{n \times 1} = X_{n \times p} \beta_{p \times 1} + \varepsilon_{n \times 1} \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I)$$

Now consider $n = p$, our goal is to estimate β .

The Ordinary Least Square Estimator $\hat{\beta}^{OLS} = \arg \min \|y - X\beta\|_2^2$ has the analytical closed form $\hat{\beta}^{OLS} = (X^T X)^{-1} X^T y$ and distribution $\hat{\beta}^{OLS} \sim \mathcal{N}(\beta_0, \sigma^2 (X^T X)^{-1})$ for β_0 is the true value.

The OLS estimator had the problem of *overfitting*. So a penalized OLS, called the *Ridge Regression* is invented. It has the following form

$$\hat{\beta} = \arg \min \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$$

2 The Ridge Regression derivation

Denote the true β be β_0 and our estimate be $\hat{\beta}$. Our $\hat{\beta}$ can be found by solving the optimization problem

$$\hat{\beta} = \arg \min \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$$

Thus cost function J is a sum of the error and the size of the β

$$J = \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$$

Since for a vector v , $\|v\|_2^2 = v^T v$, so

$$\begin{aligned} \|y - X\beta\|_2^2 &= (y - X\hat{\beta})^T (y - X\hat{\beta}) = y^T y - y^T X\hat{\beta} - \hat{\beta}^T X^T y + \hat{\beta}^T X^T X\hat{\beta} \\ \|\beta\|_2^2 &= \beta^T \beta \end{aligned}$$

Thus

$$J = y^T y - y^T X\hat{\beta} - \hat{\beta}^T X^T y + \hat{\beta}^T X^T X\hat{\beta} + \lambda \hat{\beta}^T \hat{\beta}$$

Now, as we want to minimize J by setting $\frac{\partial J}{\partial \hat{\beta}}$

$$\frac{\partial J}{\partial \hat{\beta}} = \frac{\partial}{\partial \hat{\beta}} [y^T y - y^T X\hat{\beta} - \hat{\beta}^T X^T y + \hat{\beta}^T X^T X\hat{\beta} + \lambda \hat{\beta}^T \hat{\beta}]$$

Formula for Matrix calculus

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} = \mathbf{0} \quad \frac{\partial \mathbf{x}}{\partial \mathbf{x}} = \mathbf{I} \quad \frac{\partial \mathbf{A}\mathbf{x}}{\partial \mathbf{x}} = \mathbf{A} \quad \frac{\partial \mathbf{x}^T \mathbf{A}}{\partial \mathbf{x}} = \mathbf{A}^T \quad \frac{\partial \mathbf{x}^T \mathbf{A}\mathbf{x}}{\partial \mathbf{x}} = \begin{cases} (\mathbf{A} + \mathbf{A}^T) \mathbf{x} & \text{General } \mathbf{A} \\ 2\mathbf{A}\mathbf{x} & \text{Symmetric } \mathbf{A} \end{cases}$$

Since $X^T X$ is symmetric matrix, and therefore

$$\begin{aligned}\frac{\partial J}{\partial \hat{\beta}} &= \left[-y^T X - y^T X + 2X^T X \hat{\beta} + 2\lambda \hat{\beta} \right] \\ \frac{\partial \|y - X \hat{\beta}\|_2^2}{\partial \hat{\beta}} &= -2X^T y + 2(X^T X + \lambda I) \hat{\beta}\end{aligned}$$

Now set this derivative equal to zero, we get

$$-2X^T y + 2(X^T X + \lambda I) \hat{\beta} = 0$$

$$X^T y = (X^T X + \lambda I) \hat{\beta}$$

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

Therefore, now we have the analytical solution of $\hat{\beta}$ that minimize J .

3 The derivation of $\hat{\beta}^{Ridge} \sim \mathcal{N} \left(\left[\left(I + \lambda (X^T X)^{-1} \right)^{-1} - I \right] \beta_0, \sigma^2 (X^T X + \lambda I)^{-1} X^T \right)$

It can be shown that $\hat{\beta} = \left(I + \lambda (X^T X)^{-1} \right)^{-1} \beta_0 + (X^T X + \lambda I)^{-1} X^T \varepsilon$

First, the original expression of $\hat{\beta}$ is

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

Take out the $X^T X$

$$\hat{\beta} = \left(I + \lambda (X^T X)^{-1} \right)^{-1} (X^T X)^{-1} X^T y$$

Since $y = X\beta + \varepsilon$, where β is the “true β ”, denoted as β_0 , then

$$\hat{\beta} = \left(I + \lambda (X^T X)^{-1} \right)^{-1} (X^T X)^{-1} X^T [X\beta_0 + \varepsilon]$$

Expand

$$\hat{\beta} = \left(I + \lambda (X^T X)^{-1} \right)^{-1} (X^T X)^{-1} X^T X \beta_0 + \left(I + \lambda (X^T X)^{-1} \right)^{-1} (X^T X)^{-1} X^T \varepsilon$$

$$\hat{\beta} = \left(I + \lambda (X^T X)^{-1} \right)^{-1} \beta_0 + \left(I + \lambda (X^T X)^{-1} \right)^{-1} (X^T X)^{-1} X^T \varepsilon$$

For the second term, move the $X^T X$ back

$$\hat{\beta} = \left(I + \lambda (X^T X)^{-1} \right)^{-1} \beta_0 + (X^T X + \lambda I)^{-1} X^T \varepsilon$$

Now we can show that $\hat{\beta}^{Ridge} \sim \mathcal{N} \left(\left[\left(I + \lambda (X^T X)^{-1} \right)^{-1} - I \right] \beta_0, \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} \right)$

First, consider the bias (the mean of $\hat{\beta} - \beta_0$)

$$\mathbb{E}(\hat{\beta} - \beta_0) = \mathbb{E} \left[\left(I + \lambda (X^T X)^{-1} \right)^{-1} \beta_0 + (X^T X + \lambda I)^{-1} X^T \varepsilon - \beta_0 \right]$$

Recall that for $y = X\beta + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_{p \times p})$. That means ε is having zero mean and ε_i are independent from each other ($\mathbb{E}[\varepsilon \varepsilon^T] = \sigma^2 I$)

$$\mathbb{E}(\hat{\beta} - \beta_0) = \mathbb{E} \left[\left((I + \lambda (X^T X)^{-1})^{-1} - I \right) \beta_0 \right] + (X^T X + \lambda I)^{-1} X^T \underbrace{\mathbb{E} \varepsilon}_0$$

So

$$Bias(\hat{\beta}^{Ridge}) = \mathbb{E}(\hat{\beta} - \beta_0) = \mathbb{E} \left(\left((I + \lambda (X^T X)^{-1})^{-1} - I \right) \beta_0 \right)$$

Therefore, $\hat{\beta}^{OLS}$ is a *biased* estimator of β

Now check for variance

$$Var[\hat{\beta}] = \mathbb{E} \left[(\hat{\beta} - \beta_0) (\hat{\beta} - \beta_0)^T \right]$$

Since only the noise part ε have nonzero variance, so

$$Var[\hat{\beta}] = \mathbb{E} \left[\left((X^T X + \lambda I)^{-1} X^T \varepsilon \right) \left((X^T X + \lambda I)^{-1} X^T \varepsilon \right)^T \right]$$

Since $(PQ)^T = Q^T P^T$ and $(ABC)^T = C^T B^T A^T$

$$Var[\hat{\beta}] = \mathbb{E} \left[(X^T X + \lambda I)^{-1} X^T \varepsilon \varepsilon^T X (X^T X + \lambda I)^{-1} \right]$$

$$Var[\hat{\beta}] = (X^T X + \lambda I)^{-1} X^T \mathbb{E}[\varepsilon \varepsilon^T] X (X^T X + \lambda I)^{-1}$$

Recall $\mathbb{E}[\varepsilon \varepsilon^T] = \sigma^2 I$

$$Var[\hat{\beta}] = (X^T X + \lambda I)^{-1} X^T \sigma^2 X (X^T X + \lambda I)^{-1}$$

$$Var[\hat{\beta}] = \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}$$

Therefore

$$\hat{\beta}^{Ridge} \sim \mathcal{N} \left(\left[(I + \lambda (X^T X)^{-1})^{-1} - I \right] \beta_0, \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} \right)$$