

1. Ridge regression solves the regularized least squares problem

$$\hat{\beta}_T = \arg \min_{\beta} (y - X\beta)^T (y - X\beta) + T\beta^T \beta$$

Here we can use the singular value decomposition of the design matrix X .

$$X = UDV^T$$

- ★ U is an orthogonal matrix of size $N \times N$
- ★ D is a diagonal matrix of size $N \times P$ with non-negative diagonal elements sorted in descending order.
- ★ V is an orthogonal matrix of size $P \times P$.

By substituting the SVD of X into the ridge regression problem,

$$\begin{aligned}\beta_T &= \arg \min_{\beta} (y - UDV^T \beta)^T (y - UDV^T \beta) + T\beta^T \beta \\ &= \arg \min_{\beta} (y - UD\hat{\beta})^T (y - UD\hat{\beta}) + T(\hat{\beta}^T \hat{\beta})\end{aligned}$$

where $\hat{\beta} = V^T \beta$

Here, we expand the above expression:

$$\begin{aligned}(y - UD\hat{\beta})^T (y - UD\hat{\beta}) + T\hat{\beta}^T \hat{\beta} &= y^T y - 2\hat{\beta}^T U^T y D + \hat{\beta}^T D^T \\ &= y^T y - 2\hat{\beta}^T U^T y D + \hat{\beta}^T (D^T D + T I) \hat{\beta}\end{aligned}$$

$$\begin{aligned}(y - UD\hat{\beta})^T (y - UD\hat{\beta}) + T\hat{\beta}^T \hat{\beta} &= y^T y - 2\hat{\beta}^T U^T y D + \hat{\beta}^T D^T \hat{\beta} + T\hat{\beta}^T \hat{\beta} \\ &= y^T y - 2\hat{\beta}^T U^T y D + \hat{\beta}^T (D^T D + T I) \hat{\beta}\end{aligned}$$

This upper expression is minimized when

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

$$\begin{aligned}
\beta_T &= V(D^T D + \tau I)^{-1} U^T y D \\
&= V(D^T D + \tau I)^{-1} V^T V U^T y D \\
&= V(D^T D + \tau I)^{-1} V^T X^T y \\
&= X^T y (V D (D^T D + \tau I)^{-1} V^T)^T \\
&= X^T y (V D^T D (D^T D + \tau I)^{-1} V^T) \\
&= X^T y (V D^T (D^T D + \tau I)^{-1} V^T V) \\
&= X^T y (V D^T (D^T D + \tau I)^{-1}) \\
&= X^T y (V D^{-1} (D^T D + \tau I)^{-1}) \\
&= X^T y (D^T D + \tau I)^{-1} D^{-T} V^T \\
&= X^T y S_T^{-1} S^{-1}
\end{aligned}$$

And here, $S = X^T X$ is the ordinary scatter matrix and $S_T = X^T X + \tau I$ is the regularized scatter matrix

Expectation of β_T is given by:

$$E[\beta_T] = X^T y S_T^{-1} S^{-1} = (S_T)^{-1} S \beta^*$$

$\therefore \beta^*$ is the true model coefficient

We substitute the expression for $\hat{\beta}$ back to find the covariance of β_T .

$$\begin{aligned}
(y - U D \hat{\beta})^T (y - U D \hat{\beta}) + \tau \hat{\beta}^T \hat{\beta} &= y^T y - 2 \hat{\beta}^T U^T y D + \hat{\beta}^T D^T D \hat{\beta} + \tau \hat{\beta}^T \hat{\beta} \\
&= y^T y - 2 y^T U D V^T \hat{\beta} + \hat{\beta}^T V D U^T U D V^T \hat{\beta} + \tau \hat{\beta}^T \hat{\beta} \\
&= y^T y - 2 y^T X \beta + \beta^T X^T X \beta + \tau \beta^T \beta \\
&= y^T y - 2 y^T X \beta + \beta^T S \beta + \tau \beta^T \beta \\
&= y^T y - 2 y^T X \beta + \beta^T (S + \tau I) \beta
\end{aligned}$$

$$S + \tau I = S_T$$

The covariance of β_T !

$$\begin{aligned}\text{Cov}[\beta_T] &= (S_T)^{-1} \sigma^2 \\ &= (S + \tau I)^{-1} \sigma^2 \\ &= (S_T^{-1})^T \sigma^2 \\ &= (S_T^{-1}) \sigma^2 \\ &= (S_T)^{-1} S (S_T)^{-1} \sigma^2\end{aligned}$$

The expectation and covariance matrix of the regularized solution, averaged over all possible training sets of size N , are !

$$E[\beta_T] = (S_T)^{-1} S \beta^*$$

&

$$\text{Cov}[\beta_T] = (S_T)^{-1} S (S_T)^{-1} \sigma^2$$

where $S = X^T X$ & $S_T = X^T X + \tau I$.

2) Expanding the square loss term

$$\frac{\partial}{\partial \beta} \left(- \sum_{i=1}^N (y_i^* - X_i \cdot \beta)^2 \right)$$

$$= \sum_{i=1}^N 2(y_i^* - X_i \cdot \beta)(-X_i)$$

$$\rightarrow -2 \sum_{i=1}^N X_i^T (y_i^* - X_i \cdot \beta)$$

$\rightarrow y_i^* \rightarrow y_i$ (true class label) & rewrite eqⁿ:

$$-2 \sum_{i=1}^N X_i^T (y_i - X_i \cdot \beta) = -2 \sum_{i=1}^N (X_i^T y_i - X_i^T X_i \cdot \beta)$$

Distribute the transpose:

$$-2 \sum_{i=1}^N (y_i^T X_i - \beta^T X_i^T X_i)$$

$\sum_{i=1}^N y_i^T X_i$ gets canceled out since classes are balanced

$$\& \sum_{i=1}^N y_i = 0$$

\rightarrow Simplify the equation:

$$-2 \sum_{i=1}^N (-\beta^T X_i^T X_i) = 2 \sum_{i=1}^N \beta^T X_i^T X_i$$

Re-arrange the terms:

$$2 \sum_{i=1}^N \beta^T X_i^T X_i = 2N \beta^T \hat{\Sigma} \quad \text{where } \hat{\Sigma} \text{ is covariance matrix}$$

$$E \sum_{i=1}^N X_i^T X_i = N \hat{\Sigma}$$

→ set the derivative to 0:

$$2N \beta^T \hat{\Sigma} + \frac{1}{4} (\mu_1 - \mu_{-1})^T (\mu_1 - \mu_{-1}) \beta = 0$$

→ Rearrange the eqⁿ:

$$\hat{\Sigma} \beta + \frac{1}{4N} (\mu_1 - \mu_{-1})^T (\mu_1 - \mu_{-1}) \beta = 0$$

→ Re write $(\mu_1 - \mu_{-1}) \beta$ as τ' for scalar τ'
 E ~~move~~ move second term to right hand side:

$$\sum \beta = - \frac{1}{4N} (\mu_1 - \mu_{-1})^T (\mu_1 - \mu_{-1}) \beta$$

$$\Rightarrow \hat{\Sigma} \beta = - \frac{1}{2N} (\mu_1 - \mu_{-1})^T (\mu_1 - \mu_{-1}) \beta$$

→ Rearrange the eqⁿ with $\tau = \frac{1}{2} - \frac{\tau'}{4}$:

$$\hat{\Sigma} \beta = \tau (\mu_1 - \mu_{-1})^T$$

Divide the eqⁿ by $\hat{\Sigma}$: we get

$$\hat{\beta}_{OLS} = \tau \hat{\Sigma}^{-1} (\mu_1 - \mu_{-1})^T \quad \text{with } \tau = \frac{1}{2} - \frac{\tau'}{4}$$