

AI2611 Machine Learning (Spring 2023)

Assignment 1

Xiangyuan Xue (521030910387)

Ridge regression is a generalized linear regression which introduces an L-2 regularization term. The object function can be specified as

$$L(\beta) = (\mathbf{Z}\beta - \mathbf{y})^T (\mathbf{Z}\beta - \mathbf{y}) + \lambda \cdot \beta^T \beta$$

where $\mathbf{Z} \in \mathbb{R}^{m \times n}$, $\mathbf{y} \in \mathbb{R}^m$, $\beta \in \mathbb{R}^n$ and $\lambda \in \mathbb{R}$. Let the derivative of be zero, we have

$$\nabla_{\beta} L(\beta) = 2 (\mathbf{Z}^T \mathbf{Z} \beta - \mathbf{Z}^T \mathbf{y} + \lambda \beta) = 0$$

which yields

$$\hat{\beta}_{\lambda}^{\text{ridge}} = (\mathbf{Z}^T \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}^T \mathbf{y}$$

namely the closed-form solution of ridge regression.

However, numerical computation of $(\mathbf{Z}^T \mathbf{Z} + \lambda \mathbf{I})^{-1}$ is expensive. Instead, we could use singular value decomposition (SVD) to lower the computation cost. Suppose

$$\mathbf{Z} = \mathbf{U} \mathbf{D} \mathbf{V}^T$$

where $\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p)$ is an $n \times p$ orthogonal matrix, $\mathbf{D} = \text{diag}(d_1, d_2, \dots, d_p)$ is a $p \times p$ diagonal matrix consisting singular values $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$ and $\mathbf{V} = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_p)$ is a $p \times p$ orthogonal matrix. Note that $\mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V} = \mathbf{I}_p$ and $\mathbf{D}^T \mathbf{D} = \mathbf{D}^2$, thus it holds that

$$\begin{aligned} \mathbf{Z}^T \mathbf{Z} &= (\mathbf{U} \mathbf{D} \mathbf{V}^T)^T (\mathbf{U} \mathbf{D} \mathbf{V}^T) \\ &= \mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{U} \mathbf{D} \mathbf{V}^T \\ &= \mathbf{V} \mathbf{D}^T \mathbf{D} \mathbf{V}^T \\ &= \mathbf{V} \mathbf{D}^2 \mathbf{V}^T \end{aligned}$$

Then the closed-form solution can be rewritten as

$$\begin{aligned} \hat{\beta}_{\lambda}^{\text{ridge}} &= (\mathbf{V} \mathbf{D}^2 \mathbf{V}^T + \lambda \mathbf{I}_p)^{-1} \mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{y} \\ &= [\mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I}_p) \mathbf{V}^T]^{-1} \mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{y} \\ &= (\mathbf{V}^T)^{-1} (\mathbf{D}^2 + \lambda \mathbf{I}_p)^{-1} \mathbf{V}^{-1} \mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I}_p)^{-1} \mathbf{D}^T \mathbf{U}^T \mathbf{y} \end{aligned}$$

Since $\mathbf{D}^2 + \lambda \mathbf{I}_p = \text{diag}_j(d_j^2 + \lambda)$, we have $(\mathbf{D}^2 + \lambda \mathbf{I}_p)^{-1} = \text{diag}_j\left(\frac{1}{d_j^2 + \lambda}\right)$, then

$$\begin{aligned}\hat{\boldsymbol{\beta}}_{\lambda}^{\text{ridge}} &= \mathbf{V} \text{diag}_j\left(\frac{1}{d_j^2 + \lambda}\right) \mathbf{D} \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V} \text{diag}_j\left(\frac{1}{d_j^2 + \lambda}\right) \text{diag}_j(d_j) \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V} \text{diag}_j\left(\frac{d_j}{d_j^2 + \lambda}\right) \mathbf{U}^T \mathbf{y}\end{aligned}$$

which is just the required form. Since SVD can be computed efficiently, the computation of ridge regression can be largely accelerated.