Ridge Regression and PCA

Archer Yang

McGill University

September 25, 2024

Reading assignment: Chapter 3.2-3.4 The Elements of Statistical Learning

EFFICIENT COMPUTATION OF RIDGE REGRESSION

Ridge regression

 \blacksquare The ridge penalized least squares estimator of $\hat{\boldsymbol{\beta}}^{(\lambda)}$ is defined by

$$\hat{\boldsymbol{\beta}}^{(\lambda)} = \underset{\boldsymbol{\beta}}{\arg\min} \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_2^2, \tag{1}$$

where $\lambda \ge 0$

Solve
$$\partial f(\beta^{(\lambda)}) = 0$$
 for $\beta^{(\lambda)}$, which is

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

$$(\mathbf{X}^{\mathsf{T}}\mathbf{X} + \lambda\mathbf{I}_{p})\beta^{(\lambda)} = \mathbf{X}^{\mathsf{T}}\mathbf{y} \qquad (2)$$

$$\beta^{(\lambda)} = (\mathbf{X}^{\mathsf{T}}\mathbf{X} + \lambda\mathbf{I}_{p})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}, \qquad (3)$$

- (2) can be solved with a linear system solver
- Or use the less efficient closed-form solution in (3)
- Both methods have to be recomputed for a different value of λ, computationally expensive when p is large and inefficient if we wish to compute β^(λ) for multiple values of λ.

■ Using the full SVD of $X = UDV^{\top}$, where $U = \mathbb{R}^{n \times p}$, $V \in \mathbb{R}^{p \times p}$ and $D = \mathbb{R}^{p \times p}$

$$\boldsymbol{X}^{\top}\boldsymbol{X} = \boldsymbol{V}\boldsymbol{D}^{\top}\boldsymbol{U}^{\top}\boldsymbol{U}\boldsymbol{D}\boldsymbol{V}^{\top} = \boldsymbol{V}\boldsymbol{D}^{\top}\boldsymbol{D}\boldsymbol{V}^{\top}$$

So write (2) as

$$(\boldsymbol{V}\boldsymbol{D}^{\top}\boldsymbol{D}\boldsymbol{V}^{\top} + \lambda \boldsymbol{I}_{p})\boldsymbol{\beta}^{(\lambda)} = \boldsymbol{V}\boldsymbol{D}^{\top}\boldsymbol{U}^{\top}\boldsymbol{y}$$

and replacing \boldsymbol{I}_p with $\boldsymbol{V}\boldsymbol{V}^{ op}$ gives

$$\boldsymbol{V}(\boldsymbol{D}^{\top}\boldsymbol{D} + \lambda \boldsymbol{I}_{p})\boldsymbol{V}^{\top}\boldsymbol{\beta}^{(\lambda)} = \boldsymbol{V}\boldsymbol{D}^{\top}\boldsymbol{U}^{\top}\boldsymbol{y}$$

• $V(D^{\top}D + \lambda I_p)V^{\top}$ is the eigen-decomposition of $X^{\top}X + \lambda I_p$, which is positive semi-definite.

• Assuming that $\lambda > 0$,

$$\beta^{(\lambda)} = V(D^{\top}D + \lambda I_p)^{-1}V^{\top}VD^{\top}U^{\top}y$$
$$= V(D^{\top}D + \lambda I_p)^{-1}D^{\top}U^{\top}y$$
$$= VMU^{\top}y$$

where $\boldsymbol{M} = (\boldsymbol{D}^{\top}\boldsymbol{D} + \lambda \boldsymbol{I}_p)^{-1}\boldsymbol{D}^{\top} = \text{diag}(m_1, \cdots, m_p) \in \mathbb{R}^{p \times p}$ is diagonal where

$$m_j = \begin{cases} d_j / (d_j^2 + \lambda) & j = 1, \cdots, r \\ 0 & j = r, \cdots, p \end{cases}$$

• We can use the reduced SVD on X, for rank $(X) = r \le \min(n, p)$,

$$\boldsymbol{X} = \boldsymbol{U}_r \boldsymbol{D}_r \boldsymbol{V}_r^{\mathsf{T}}$$

where \boldsymbol{U}_r = $\boldsymbol{U}_{[,1,\cdots,r]}$ and \boldsymbol{V}_r = $\boldsymbol{V}_{[,1,\cdots,r]}$, therefore

$$\boldsymbol{\beta}^{(\lambda)} = \boldsymbol{V}_r \boldsymbol{M}_r \boldsymbol{U}_r^\top \boldsymbol{y}$$

where $M_r = \text{diag}(m_1, \dots, m_r) \in \mathbb{R}^{r \times r}$ with diagonal values $m_j = d_j/(d_j^2 + \lambda).$

CONNECTION BETWEEN RIDGE REGRESSION AND PCA

OLS solution

Using the SVD we can write the OLS fitted vector as

$$\hat{\boldsymbol{y}}^{\text{OLS}} = \boldsymbol{X} \hat{\boldsymbol{\beta}}^{\text{OLS}} = \boldsymbol{X} (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top} \boldsymbol{y}$$
$$= \boldsymbol{U} \boldsymbol{D} \boldsymbol{V}^{\top} (\boldsymbol{V} \boldsymbol{D} \boldsymbol{U}^{\top} \boldsymbol{U} \boldsymbol{D} \boldsymbol{V}^{\top})^{-1} \boldsymbol{V} \boldsymbol{D} \boldsymbol{U}^{\top} \boldsymbol{y}$$
$$= \boldsymbol{U} \boldsymbol{U}^{\top} \boldsymbol{y}.$$

U^Ty are the coordinates of y with respect to the basis U.
We see that XV = X(v₁, ..., v_r) = (u₁, ..., u_r)D = UD.
u_j can be viewed as scaled version of the coordinates z_j, because

$$\boldsymbol{z}_j = \boldsymbol{X} \boldsymbol{v}_j = \boldsymbol{u}_j d_j$$
 $\operatorname{Var}(\boldsymbol{z}_j) = \operatorname{Var}(\boldsymbol{X} \boldsymbol{v}_j) = \frac{d_1^2}{N}$

Can view \hat{y}^{OLS} is a reconstruced version of y

Ridge solutions

On the other hand, the ridge solutions are

$$\boldsymbol{X} \hat{\boldsymbol{\beta}}^{\text{ridge}} = \boldsymbol{X} (\boldsymbol{X}^{\top} \boldsymbol{X} + \lambda \boldsymbol{I}_{p})^{-1} \boldsymbol{X}^{\top} \boldsymbol{y}$$
$$= \boldsymbol{U} \boldsymbol{D} (\boldsymbol{D}^{2} + \lambda \boldsymbol{I}_{p})^{-1} \boldsymbol{D} \boldsymbol{U}^{\top} \boldsymbol{y}$$
$$= \sum_{j=1}^{p} \boldsymbol{u}_{j} \frac{d_{j}^{2}}{d_{j}^{2} + \lambda} \boldsymbol{u}_{j}^{\top} \boldsymbol{y}$$
$$= \boldsymbol{U} \boldsymbol{H} \boldsymbol{U}^{\top} \boldsymbol{y}$$

where

$$\boldsymbol{H} = \boldsymbol{D}(\boldsymbol{D}^{\top}\boldsymbol{D} + \lambda \boldsymbol{I}_p)^{-1}\boldsymbol{D}$$

is diagonal with

$$h_j = d_j^2 / (d_j^2 + \lambda)$$
 $j = 1, \dots, p$

Ridge solutions

- Like $\hat{\boldsymbol{\beta}}^{\text{OLS}}$, $\hat{\boldsymbol{\beta}}^{\text{ridge}}$ projects \mathbf{y} to $\mathbf{U}^T \mathbf{y}$ with respect to the orthonormal basis \mathbf{U} .
- It then shrinks these projected y by the factors $d_j^2/(d_j^2 + \lambda)$.

• And then reconstructs \mathbf{y} using $\mathbf{UHU}^T \mathbf{y} = \sum_{j=1}^p u_j \frac{d_j^2}{d_j^2 + \lambda} u_j^\top \mathbf{y}$. This means that a greater amount of shrinkage is applied to the coordinates of basis vectors with smaller d_j^2 .