

Linear transformations of covariates in regression:

LS in original variables:

$$\textcircled{1} \quad \| \tilde{Y} - X\beta \|_2^2 = \sum_{i=1}^n (Y_i - x_i^T \beta)^2,$$

x_i^T is the i th row of the design matrix X .

$(p+1) \times (p+1)$

$\textcircled{2}$ using linearly transformed variables: $B = X \cdot A$

$$\| Y - B \cdot b \|_2^2$$

$n \times (p+1)$

$$\textcircled{1} \quad \| Y - X\beta \|_2^2$$

$$X \cdot I \cdot \beta = \underbrace{X \cdot A}_B \cdot \underbrace{A^{-1} \cdot \beta}_b$$

\Rightarrow if $\hat{\beta}$ minimizes $\textcircled{1}$, then $\hat{b} = A^{-1} \cdot \hat{\beta}$

minimizes $\textcircled{2}$

if \hat{b} minimizes $\textcircled{2}$, then $\hat{\beta} = A \cdot \hat{b}$

is the minimizer of $\textcircled{1}$.

Case 1: recentering covariates

$$X = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix} = \left(\underbrace{\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, X_1, \dots, X_p}_{\text{columns of } X} \right);$$

$$A_c = \begin{pmatrix} 1 & -\bar{x}_1 & -\bar{x}_2 & \dots & -\bar{x}_p \\ & 1 & & & 0 \\ & & \ddots & & \\ 0 & & & 1 & \\ & & & & 1 \end{pmatrix};$$

\bar{x}_j is the mean of X_j .

A_c is of full rank.

$$X \cdot A_c[:, 1] = \tilde{X}_1 - \bar{x}_1 \cdot \mathbf{1}$$

Case 2: ^(rc) scaling of features.

Idea: divide X_j by its sample std. dev., to get \tilde{X}_j that has unit (sample) variance.

$$A_s = \begin{pmatrix} 1 & & & & \\ & 1/\sqrt{s_1^2} & & & 0 \\ & & \ddots & & \\ 0 & & & 1/\sqrt{s_p^2} & \end{pmatrix};$$

s_j^2 is the sample variance for X_j .

Case 3: centering and scaling.

$$\tilde{X} = (X \cdot A_c) \cdot A_s.$$

Lasso & ridge regression w/o the intercept.

let $\tilde{Y} = Y - \bar{y} \cdot \frac{1}{n}$
the sample mean of Y .

let $\tilde{X} = (\tilde{X}_1, \dots, \tilde{X}_p)$, where

$$\tilde{X}_j = (X_j - \bar{x}_j \cdot \frac{1}{n}) / \sqrt{s_j^2}$$

equivalent criteria are

$$\| \tilde{Y} - \tilde{X} \tilde{\beta} \|_2^2 + \lambda \sum_{j=1}^p |\tilde{\beta}_j|^q,$$

note: no intercept

$q = 1$: lasso

$q = 2$: ridge regression

Singular value decomposition. (SVD)

X be $n \times p$ matrix, $n \geq p$. E.g., X_j 's centered and scaled; no intercept.

SVD of X is $X = U \cdot D \cdot V^T$, such that $U^T U = I_p$, $V^T V = I_p$. D is a diagonal matrix.

$D_{11} \geq D_{22} \geq \dots \geq D_{pp} \geq 0$; the d_j 's are called singular values.

of nonzero d_j 's is the rank of X .

Often, PCA is defined in terms of SVD (spectral decomp.) of $(X^T X)$.

original LS criterion is

$$\textcircled{1} \quad \| Y - X\beta \|_2^2$$

From svd, $X = U \cdot \underbrace{D \cdot V^T}_A$

Alternatively, one
can take $A = V^T$

$$\textcircled{2} \quad \| Y - \underbrace{U b}_{\sum_{j=1}^p u_j b_j} \|_2^2 \quad : \text{ no dimension reduction}$$

$$\| Y - \sum_{j=1}^M u_j b_j \|_2^2 : \text{ principal components regression.}$$

Dimension reduction: $M < p$.

We use only a subset of u_j 's.