

· (Deviations from Means) Let $S = X_1^T X_1 - n\bar{X}\bar{X}^T = Z^T Z$ and $Z = X_1 - 1\bar{X}^T$ then $\hat{\beta}_0 = \bar{y} - \bar{X}^T \hat{b} = \frac{1}{n} 1^T Y - \bar{X}^T \hat{b}$, $\hat{b} = S^{-1}(X_1^T Y - n\bar{y}\bar{X}) = (Z^T Z)^{-1} Z^T Y$ and $\text{var}(\hat{b}) = (Z^T Z)^{-1} \sigma^2$, $\text{var}(\hat{\beta}_0) = \frac{\sigma^2}{n} + \bar{X}^T \text{var}(\hat{b}) \bar{X}$, $\text{cov}(\hat{\beta}_0, \hat{b}) = -\bar{X}^T \text{var}(\hat{b})$.

· $\frac{\ln-r(X)\hat{\sigma}^2}{\sigma^2} = \frac{SSE}{\sigma^2} = \frac{1}{\sigma^2} Y^T (I - H) Y \sim \chi^2_{n-r(X)}$, $\frac{Y^T Y}{\sigma^2} \sim \chi^2(n, \frac{\beta^T X^T X \beta}{2\sigma^2})$, $\frac{SSR}{\sigma^2} \sim \chi^2(r(X), \frac{\beta^T X^T X \beta}{2\sigma^2})$, $\frac{MSR}{MSE} \sim F[r(X), n-r(X), \frac{\beta^T X^T X \beta}{2\sigma^2}]$.

· $SST_m = SST - SSM = Y(I - \frac{1}{n} 11^T) Y$, $\frac{SSM}{\sigma^2} \sim \chi^2(1, \frac{(1^T X \beta)^2}{2n\sigma^2})$, $SSR_m = SSR - SSM = \hat{b}^T Z^T Y = \hat{b}^T (Z^T Z) \hat{b}$, $\frac{SST_m}{\sigma^2} \sim \chi^2(n-1, \frac{\beta^T X^T X \beta - \frac{1}{n} (1^T X \beta)^2}{2\sigma^2})$, $\frac{SSR_m}{\sigma^2} \sim \chi^2(r-1, \frac{b^T (Z^T Z) b}{2\sigma^2})$, $\frac{MSM}{MSE} \sim F[1, n-r(X), \frac{(1^T X \beta)^2}{2n\sigma^2}]$, $\frac{MSR_m}{MSE} \sim F[r(X)-1, n-r(X), \frac{b^T (Z^T Z) b}{2\sigma^2}]$.

· (General linear hypothesis) The general hypothesis we consider is $H_0: K^T \beta = m$ where β is the $(k+1)$ -dimensional vector of parameters of the model, K^T is any full row rank matrix of size $s \times (k+1)$ and m is a $s \times 1$ vector of specified constants. Define $Q = (K^T \beta - m)^T [K^T (X^T X)^{-1} K]^{-1} (K^T \beta - m)$, $\frac{Q}{\sigma^2} \sim \chi^2(s, \frac{1}{2\sigma^2} (K^T \beta - m)^T [K^T (X^T X)^{-1} K]^{-1} (K^T \beta - m))$, $F(H) = \frac{Q/s}{SSE/(n-r(X))} \sim F(s, n-r(X), (K^T \beta - m)^T [K^T (X^T X)^{-1} K]^{-1} (K^T \beta - m))$.

· (Estimation under the null hypothesis) When H_0 is true, $\hat{\beta}$ is derived so as to minimize the least squares objective function subject to the constraint $K^T \beta = m$, $\tilde{\beta} = \hat{\beta} - (X^T X)^{-1} K (K^T (X^T X)^{-1} K)^{-1} (K^T \hat{\beta} - m)$ and $SSE_{H_0} = SSE + Q$.

Models not of Full Rank

· (Properties of $\beta^0 = GX^T Y$) (1) $E(\beta^0) = GX^T X \beta$. (2) $\text{Var}(\beta^0) = GX^T X G^T \sigma^2$. (3) $\hat{y} = X \beta^0 = XGX^T Y$. (4) $E(\hat{y}) = XGX^T X \beta = X \beta$. (5) $SSE = y^T (I - XGX^T) y$. (6) $SSR = y^T XGX^T y = (\beta^0)^T X^T y$. (7) $\hat{\sigma}^2 = \frac{SSE}{n-r(X)}$ is an unbiased estimator of σ^2 . (8) $SSR_m = y^T (XGX^T - \frac{11^T}{n}) y$.

· (Distributional Properties) (1) $\beta^0 \sim N(GX^T X \beta, GX^T X G^T \sigma^2)$. The covariance matrix is singular. (2) $\frac{SSE}{\sigma^2} \sim \chi^2_{n-r(X)}$, $\frac{SSR}{\sigma^2} \sim \chi^2(r(X), \frac{1}{\sigma^2} \beta^T X^T X \beta)$, $\frac{MSR}{MSE} \sim F(r(X), n-r(X), \frac{1}{2\sigma^2} \beta^T X^T X \beta)$.

· (Identifiability) Formally, the parameter β is identifiable if $f(\beta_1) = f(\beta_2)$ implies that $\beta_1 = \beta_2$ for any β_1 and β_2 . More generally, the vector-valued function $g(\beta)$ is identifiable if $f(\beta_1) = f(\beta_2)$ implies that $g(\beta_1) = g(\beta_2)$.

[Proposition] In a linear model for which X is of full rank, β is identifiable.

[Proposition] A function $g(\beta)$ is identifiable if and only if $g(\beta) = (h \circ f)(\beta)$ for some function h .

· (Estimable functions) If a vector t exists such that $t^T E(y) = q^T \beta$, then $q^T \beta$ is said to be estimable. Linear combinations of estimable functions are estimable.

· $E(t^T y) = t^T X \beta = q^T \beta \Rightarrow t^T X = q^T$ for some t . This is equivalent to saying that q is in the row space of X .

[Theorem] (Gauss-Markov Theorem) The best linear unbiased

$$u_t = \lambda u_t + w_t + \varepsilon u_t \quad t=1, \dots, T$$

$$w_k = (w_{k1}, \dots, w_{kT})^T, \quad \varepsilon w_k = (\varepsilon w_{k1}, \dots, \varepsilon w_{kT})^T \quad \xi = (\xi_1^T, \dots, \xi_T^T, \xi_{T+1}^T, \dots, \xi_{T+q}^T)$$

estimator of the estimable function $q^T \beta$ is $q^T \beta^0$. (Sketch of proof) $\text{var}(q^T \beta^0) = t^T XGX^T t \sigma^2 = q^T G q \sigma^2$. Suppose $k^T y$ is another linear unbiased estimator of $q^T \beta$, $\text{cov}(q^T \beta^0, k^T y) = \text{cov}(q^T GX^T y, k^T y) = q^T G q \sigma^2$. $\text{var}(q^T \beta^0 - k^T y) = \text{var}(q^T \beta^0) + \text{var}(k^T y) - 2\text{cov}(q^T \beta^0, k^T y) = \text{var}(k^T y) + q^T G q \sigma^2 - 2q^T G q \sigma^2 = \text{var}(k^T y) - \text{var}(q^T \beta^0) \geq 0$.

[Theorem] (Test of Estimability) The linear function $q^T \beta$ is

estimable if and only if $q^T H = q^T$ where $H = GX^T X$.

· (Testable hypothesis) A testable hypothesis is $H_0: K^T \beta = m$ such that $K^T = T^T X$ for some matrix T^T of order $r \times n$. The matrix K^T of size $r \times p$ is always of full-row rank.

· (Hypothesis testing) $K^T \beta^0 - m \sim N(K^T \beta - m, K^T G K \sigma^2)$, $Q = (K^T \beta^0 - m)^T (K^T G K)^{-1} (K^T \beta^0 - m)$. Then

$$\frac{Q}{\sigma^2} \sim \chi^2(r, (K^T \beta - m)^T (K^T G K)^{-1} (K^T \beta - m) / 2\sigma^2),$$

$$F(H) = \frac{Q/r}{SSE/(n-r(X))} \sim F(r, n-r(X), (K^T \beta - m)^T (K^T G K)^{-1} (K^T \beta - m) / 2\sigma^2).$$

The Bias-Variance Trade-off

· (The bias-variance trade-off) Consider the expected prediction error

$$\text{PE}(z_0) = E_{Y|Z=z_0}[(Y - \hat{f}(Z))^2 | Z = z_0] = \sigma_\varepsilon^2 + \text{Bias}^2(\hat{f}(z_0)) + \text{Var}(\hat{f}(z_0)) = \sigma_\varepsilon^2 + \text{MSE}(\hat{f}(z_0)).$$

Ridge Regression

· (Ridge regression) Minimize $(Y - Z\beta)^T (Y - Z\beta)$ s.t. $\sum_{j=1}^p \beta_j^2 \leq t$.

$\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} \|Y - Z\beta\|_2^2 + \lambda \|\beta\|_2^2$. Its solution may have smaller average PE than $\hat{\beta}^{\text{ls}}$. Note that $\text{PRSS}(\beta)_{\ell_2}$ is convex, and hence has a unique solution. $\hat{\beta}_\lambda^{\text{ridge}} = (Z^T Z + \lambda I_p)^{-1} Z^T Y$.

· By convention, (1) Z is assumed to be standardized (mean = 0, variance = 1) and (2) Y is assumed to be centered.

· ($\hat{\beta}_\lambda^{\text{ridge}}$ is biased). Let $R = Z^T Z$ and Z has full column rank.

$$\hat{\beta}_\lambda^{\text{ridge}} = (Z^T Z + \lambda I_p)^{-1} Z^T Y = (R + \lambda I_p)^{-1} R (R^{-1} Z^T Y) = [R(I_p + \lambda R^{-1})]^{-1} R \hat{\beta}^{\text{ls}} = (I_p + \lambda R^{-1})^{-1} \hat{\beta}^{\text{ls}}.$$

Thus, $E(\hat{\beta}_\lambda^{\text{ridge}}) = E[(I_p + \lambda R^{-1})^{-1} \hat{\beta}^{\text{ls}}] = (I_p + \lambda R^{-1})^{-1} \beta \neq \beta$ ($\forall \beta \neq 0, \lambda \neq 0$) Note that 1 is not an eigenvalue of the matrix $(I_p + \lambda R^{-1})^{-1}$ as $\lambda \neq 0$ and R is positive definite.

· (SVD) By the singular value decomposition, $Z = UDV^T$ where U is an $n \times p$ orthogonal matrix, D is $p \times p$ diagonal matrix and V^T is a $p \times p$ orthogonal matrix. $\hat{\beta}_\lambda^{\text{ridge}} = V \text{diag}(\frac{d_j}{d_j^2 + \lambda}) U^T Y$, by using

$$Z^T Z = (UDV^T)^T (UDV^T) = V D^T U^T U D V^T = V D^2 V^T.$$

· (Orthonormal Z) Z is orthonormal, then $Z^T Z = I_p$. Let $\hat{\beta}^{\text{ls}}$ denote

$$\text{the LS solution, then } \hat{\beta}_\lambda^{\text{ridge}} = \frac{1}{1+\lambda} \hat{\beta}^{\text{ls}} \text{ and } \hat{\beta}_j^{\text{ridge}} = \frac{\hat{\beta}_j^{\text{ls}}}{1+\lambda}.$$

The optimal choice of λ minimizing the expected prediction error is $\lambda^* = \frac{\sigma^2}{\sum_{j=1}^p \beta_j^2}$ where $\beta = (\beta_1, \beta_2, \dots, \beta_p)$ is the true coefficient vector.

· Variance of the ridge regression estimate is $\text{Var}(\hat{\beta}^{\text{ridge}}) = \sigma^2 W_\lambda (Z^T Z)^{-1} W_\lambda = \sigma^2 (Z^T Z + \lambda I_p)^{-1} Z^T Z (Z^T Z + \lambda I_p)^{-1}$ where $W_\lambda = (Z^T Z + \lambda I_p)^{-1} Z^T Z$. Recall that $\hat{\beta}^{\text{ridge}} = W_\lambda \hat{\beta}^{\text{ls}}$.

· The bias of the ridge regression estimate is

$$\text{bias}(\hat{\beta}^{\text{ridge}}) = -\lambda W_\lambda (Z^T Z)^{-1} \beta = -\lambda (Z^T Z + \lambda I_p)^{-1} \beta.$$

$$u_k = \begin{pmatrix} u_{k1} \\ u_{k2} \\ \vdots \\ u_{kp} \end{pmatrix} = \begin{pmatrix} \lambda u_{t,1} & \lambda u_{t,2} & \dots & \lambda u_{t,p} \\ \lambda u_{t,1} & \lambda u_{t,2} & \dots & \lambda u_{t,p} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda u_{t,1} & \lambda u_{t,2} & \dots & \lambda u_{t,p} \end{pmatrix} \begin{pmatrix} u_{t,1} \\ u_{t,2} \\ \vdots \\ u_{t,p} \end{pmatrix} = \begin{pmatrix} \lambda u_{t,1} & \lambda u_{t,2} & \dots & \lambda u_{t,p} \\ \lambda u_{t,1} & \lambda u_{t,2} & \dots & \lambda u_{t,p} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda u_{t,1} & \lambda u_{t,2} & \dots & \lambda u_{t,p} \end{pmatrix} \begin{pmatrix} u_{t,1} \\ u_{t,2} \\ \vdots \\ u_{t,p} \end{pmatrix}$$

· It can be shown that (1) the total variance $\sum_{j=1}^p \text{Var}(\hat{\beta}_j^{\text{ridge}})$ is a monotone decreasing sequence with respect to λ . (2) the total squared bias $\sum_{j=1}^p \text{Bias}^2(\hat{\beta}_j^{\text{ridge}})$ is a monotone increasing sequence with respect to λ .

· (Existence Theorem) There always exists a λ such that the MSE of $\hat{\beta}^{\text{ridge}}$ is less than the MSE of $\hat{\beta}^{\text{OLS}}$.

LASSO

· (LASSO) Minimize $(Y - Z\beta)^T (Y - Z\beta)$, s.t. $\sum_{j=1}^p |\beta_j| \leq t$. It is often convenient to rewrite the LASSO problem in the so-called Lagrangian form: minimize $(Y - Z\beta)^T (Y - Z\beta) + \lambda \|\beta\|_1$, where $\lambda \geq 0$ is the so-called tuning parameter. The resulting lasso estimate is denoted by $\hat{\beta}_\lambda^{\text{lasso}}$.

· There is usually a factor $1/(2n)$ or $1/n$ appearing in front of $(Y - Z\beta)^T (Y - Z\beta)$. This kind of standardization makes λ values comparable for different sample sizes (useful for cross-validation).

· We typically center the response and standardize the predictors so that each column is centered ($\frac{1}{n} \sum_{i=1}^n Z_{ij} = 0$) and has unit variance ($\frac{1}{n} \sum_{i=1}^n Z_{ij}^2 = 1$).

· (KKT condition) The necessary and sufficient conditions for a solution to LASSO take the form

$$-(z_j, Y - Z\beta) + \lambda s_j = 0, \quad j = 1, \dots, p$$

Here each s_j is an unknown quantity equal to $\text{sign}(\beta_j)$ if $\beta_j \neq 0$ and some value lying in $[-1, 1]$ otherwise.

· (Orthonormal Z) Consider an orthogonal design case with $Z^T Z = I$. The LASSO method is equivalent to: solve β_j 's componentwisely by solving: $\min_{\beta_j} (\beta_j - \hat{\beta}_j^{\text{ls}})^2 + \lambda |\beta_j|$. The solution to the above problem is

$$\hat{\beta}_j^{\text{lasso}} = \text{sgn}(\hat{\beta}_j^{\text{ls}}) (|\hat{\beta}_j^{\text{ls}}| - \frac{\lambda}{2})_+ = \begin{cases} \hat{\beta}_j^{\text{ls}} - \frac{\lambda}{2} & \text{if } \hat{\beta}_j^{\text{ls}} > \frac{\lambda}{2} \\ 0 & \text{if } |\hat{\beta}_j^{\text{ls}}| \leq \frac{\lambda}{2} \\ \hat{\beta}_j^{\text{ls}} + \frac{\lambda}{2} & \text{if } \hat{\beta}_j^{\text{ls}} < -\frac{\lambda}{2} \end{cases}$$

Non-Normal Case

· Assume that $Y = X\beta + \epsilon$, $f_\epsilon(\cdot)$ is unknown. The score function is defined as $\nabla_\beta \log f_\epsilon(Y - X\beta) = (-X)^T (\dot{f}_\epsilon(Y - X\beta) / f_\epsilon(Y - X\beta))$. Let $g_n(\beta) = \sum_i (\dot{f}_\epsilon(y_i - x_i^T \beta) / f_\epsilon(y_i - x_i^T \beta)) x_i$. Then $g_n(\beta) = 0$ gives $\hat{\beta}_n$. $\sqrt{n}(\hat{\beta}_n - \beta_0) = -(\frac{1}{n} \sum_i \dot{g}_i(\beta_0))^{-1} (\frac{1}{\sqrt{n}} \sum_i g_i(\beta_0))$. Then

$\sqrt{n}(\hat{\beta}_n - \beta_0) \xrightarrow{d} N(0, A^{-1} B (A^{-1})^T)$, where $A = E \dot{g}(\beta_0)$ (Hessian) and $B = \text{Var}(g(\beta))$. To ensure consistency, there should be $Eg(\beta_0) = 0$.

Matrix Derivative

Formula: $df = \text{tr} \left(\frac{\partial f}{\partial X}^T dX \right)$

- $d(XY) = (dX)Y + XdY$; $d(X^T) = (dX)^T$; $d \text{tr}(X) = \text{tr}(dX)$
- $dX^{-1} = -X^{-1} dX X^{-1}$.
- $d|X| = \text{tr}(X^\# dX)$, $X^\#$ represents the adjugate matrix of X . When X is invertible $d|X| = |X| \text{tr}(X^{-1} dX)$.
- $d(X \odot Y) = dX \odot Y + X \odot dY$, \odot represents the element-wise multiplication of matrices X and Y of the same size.
- $d\sigma(X) = \sigma'(X) \odot dX$, $\sigma(X) = [\sigma(X_{ij})]$ is an element-wise scalar function operation, $\sigma'(X) = [\sigma'(X_{ij})]$ is the element-wise derivative.

$$\eta = \beta + \Pi \eta + \Gamma F(\xi) + \delta$$

$$\eta = \beta + \Lambda \delta G(\xi) + \delta$$