

1.3 の補足

最初に λ の値を十分大きくして全ての β_j を 0 に設定した後、 λ の値を下げながら、座標降下法を実行することを考える。簡単のため、各 $j = 1, \dots, p$ について $\sum_{i=1}^N x_{i,j}^2 = 1$ であっ

て、 $\sum_{i=1}^N x_{i,j} y_i$ の値が全て異なると仮定する。このとき、全ての j に対して $\beta_j = 0$ であるよ

うな λ の値は $\lambda = \max_{1 \leq j \leq p} \left| \frac{1}{N} \sum_{i=1}^N x_{i,j} y_i \right|$ で与えられる；

β_j を 1 つ選び、そのほかの β_k は 0 として固定する。このとき、(1.10) 式から L の劣微分を 0 にするような β_j を求めると

$$\begin{aligned} & -\frac{1}{N} \sum_{i=1}^N x_{i,j} \left(y_i - \sum_{k=1}^p x_{i,k} \beta_k \right) + \lambda \begin{cases} 1 & \beta_j > 0 \\ [-1, 1] & \beta_j = 0 \\ -1 & \beta_j < 0 \end{cases} \\ & = -\frac{1}{N} \sum_{i=1}^N x_{i,j} y_i + \frac{1}{p} \sum_{i=1}^N x_{i,j}^2 \beta_j + \lambda \begin{cases} 1 & \beta_j > 0 \\ [-1, 1] & \beta_j = 0 \\ -1 & \beta_j < 0 \end{cases} \quad (\because \beta_k = 0 \ (k \neq j)) \\ & = -\frac{1}{N} \sum_{i=1}^N x_{i,j} y_i + \frac{1}{N} \beta_j + \lambda \begin{cases} 1 & \beta_j > 0 \\ [-1, 1] & \beta_j = 0 \ni 0 \\ -1 & \beta_j < 0 \end{cases} \quad \left(\because \sum_{i=1}^N x_{i,j}^2 = 1 \right) \\ & \therefore \beta_j = N \mathcal{S}_\lambda \left(\frac{1}{N} \sum_{i=1}^N x_{i,j} y_i \right) = 0 \quad \left(\because \lambda \geq \max_{1 \leq j \leq p} \left| \frac{1}{N} \sum_{i=1}^N x_{i,j} y_i \right| \right) \end{aligned}$$

となるからである。もし λ をその値より小さくすると、ある 1 つの j で $p_j = \frac{1}{N} \sum_{i=1}^N x_{i,j} y_i$ としたときに

$$\beta_j = N \mathcal{S}_\lambda(p) = N \begin{cases} p - \lambda & (p > \lambda) \\ p + \lambda & (p < \lambda) \end{cases}$$

となるので

$$\begin{aligned} \frac{1}{N} \left| \sum_{i=1}^N x_{i,j} \left(y_i - \sum_{k=1}^p x_{i,k} \beta_k \right) \right| &= \frac{1}{N} \left| \sum_{i=1}^N x_{i,j} y_i - \sum_{k=1}^N x_{i,j}^2 \beta_j \right| \\ &= \frac{1}{N} |N\lambda| = \lambda \end{aligned}$$

が成立する.

Ridge

1.1 節において, 行列 $X^T X$ が正則であるという仮定の下で二乗誤差 $\|y - X\beta\|$ を最小にする β が $\hat{\beta} = (X^T X)^{-1} X^T y$ となることを導いた.

その後 $N < p$ の場合には $X^T X$ が正則でないことを示したが, $N \geq p$ であって $X^T X$ が正則であっても, 行列式が小さければ信頼区間が大きくなるなど不都合が生じる. このような問題を避けるため, 定数 $\lambda \geq 0$ を用いて二乗誤差に β のノルムの λ 倍を加えた

$$L := \frac{1}{N} \|y - X\beta\|^2 + \lambda \|\beta\|^2$$

を最小にする方法がよく用いられる. この方法を Ridge と呼ぶ. 上式を最小にする β を求めるために, L を β で微分すると

$$\frac{\partial L}{\partial \beta} = -\frac{2}{N} X^T (y - X\beta) + 2\lambda\beta$$

となる. (ベクトル微分の公式 $\partial/\partial \mathbf{x} (A\mathbf{x} - \mathbf{b})^T (A\mathbf{x} - \mathbf{b}) = 2A^T (A\mathbf{x} - \mathbf{b})$ を用いた.) さらに $X^T X + N\lambda I$ が正則であれば, $\frac{\partial L}{\partial \beta} = 0$ となる $\hat{\beta}$ は

$$\begin{aligned} 0 &= -\frac{2}{N} X^T (y - X\beta) + 2\lambda\beta \\ &= \frac{2}{N} X^T y - \frac{2}{N} (X^T X + N\lambda) \hat{\beta} \\ (X^T X + N\lambda) \hat{\beta} &= X^T y \\ \hat{\beta} &= (X^T X + N\lambda)^{-1} X^T y \end{aligned}$$

となることがわかる. ここで, $\lambda > 0$ ならば $X^T X + N\lambda$ が正則になることがわかる. 証明は以下の通り

Proof. まず, $(X^T X)^T = X^T X$ が成立するので $X^T X$ は対称行列となる. さらに任意の $\mathbf{x} \in \mathbb{R}^p$ に対して

$$\begin{aligned} \mathbf{x}^T (X^T X) \mathbf{x} &= (\mathbf{x}^T X^T) (X \mathbf{x}) \\ &= (\mathbf{x} X)^T (\mathbf{x} X) \end{aligned}$$

となり, 最右辺は $\mathbf{x} X$ 自身の内積を示しているので $\mathbf{x}^T (X^T X) \mathbf{x} \geq 0$ であること, つまり $X^T X$ が非負定値行列であることがわかる. さらに $X^T X$ は対称行列であるから, ある直交行列 P と対角行列 Λ を用いて

$$X^T X = P^{-1} \Lambda P$$

と表すことができる．今 e_i を \mathbb{R}^p の標準基底とすると，

$$\begin{aligned}(P^{-1}e_i)^T(X^T X)(P^{-1}e_i) &= ((P^{-1}e_i)^T P^{-1})\Lambda(P(P^{-1}e_i)) \\ &= (P(P^{-1}e_i)\Lambda(e_i) \quad (\because P^{-1} = P^T) \\ &= e_i^T \Lambda e_i = \mu_i \geq 0\end{aligned}$$

となることがわかる．ただし μ_i は Λ の (i, i) 成分．したがって任意の λ に対して $\mu_i \geq 0$ であることがわかり，各 μ_i は $X^T X$ の固有値であることから $X^T X$ の全ての固有値が非負であることがわかる．さらに， $X^T X - n\lambda$ の固有値を t とすると

$$\begin{aligned}\det|(X^T X + N\lambda) - tI| &= \det|X^T X - (t - N\lambda)I| = 0 \\ &\Rightarrow t - N\lambda = \mu_i \geq 0 \quad \forall i \\ &\Leftrightarrow t = N\lambda + \mu_i > 0 \quad \forall i\end{aligned}$$

が成立するので $X^T X - n\lambda$ の固有値が全て正であることがわかり，これより $X^T X$ の行列式が 0 でないことがわかる．したがって $X^T X$ が正則であることがわかる．

□