

1. The ridge and lasso regression can be viewed as special cases of penalized least-squares

$$\min_{\beta} Q(\beta) = \frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \sum_{j=1}^p p_{\lambda}(|\beta_j|),$$

where p_{λ} is some penalization function. Consider an orthogonal design matrix, i.e., $\mathbf{X}^T \mathbf{X} = n\mathbf{I}_p$.

- (a) Show that minimizing $Q(\beta)$ is equivalent to minimizing

$$\sum_{j=1}^p \left\{ \frac{1}{2} (\hat{\beta}_j - \beta_j)^2 + p_{\lambda}(|\beta_j|) \right\},$$

where $\hat{\beta}$ is the ordinary least-square estimate.

- (b) Recall that in the case of orthogonal design matrix, lasso gives $\hat{\beta}^{\text{lasso}} = \text{sign}(\hat{\beta})(|\hat{\beta}| - \lambda)_+$. As a result, the bias of lasso estimate is approximately λ for large true value β . To address the bias, one idea is to use a penalty that tapers off as β becomes larger in absolute value, e.g., smoothly clipped absolute deviations (SCAD) penalty ($\gamma > 2$):

$$p_{\lambda, \gamma}(x) = \begin{cases} \lambda|x|, & \text{if } |x| \leq \lambda, \\ \frac{2\gamma\lambda|x| - x^2 - \lambda^2}{2(\gamma-1)}, & \text{if } \lambda < |x| \leq \gamma\lambda, \\ \frac{\lambda^2(\gamma+1)}{2}, & \text{if } |x| > \gamma\lambda. \end{cases}$$

SCAD coincides with the lasso until $|x| = \lambda$, then it is quadratic, after which it remains constant. Use part (a), show that for the SCAD penalty, the solution under orthogonal design matrix is

$$\hat{\beta}^{\text{SCAD}} = \begin{cases} \text{sign}(\hat{\beta})(|\hat{\beta}| - \lambda)_+, & \text{if } |\hat{\beta}| \leq 2\lambda, \\ \text{sign}(\hat{\beta}) \left[(\gamma - 1)|\hat{\beta}| - \gamma\lambda \right] / (\gamma - 2), & \text{if } 2\lambda < |\hat{\beta}| \leq \gamma\lambda, \\ \hat{\beta}, & \text{if } |\hat{\beta}| > \gamma\lambda. \end{cases}$$

2. Let

$$\hat{\beta}_{\lambda} = \underset{\beta}{\text{argmin}} g(\beta|\hat{\beta}, \lambda) = \underset{\beta}{\text{argmin}} \left\{ \frac{1}{2} (\hat{\beta} - \beta)^2 + p_{\lambda}(|\beta|) \right\},$$

for some penalty function $p_{\lambda}(\cdot)$. Following the convention, let $p'_{\lambda}(0) = p'_{\lambda}(0+)$ the right derivative at 0. Assume that $p_{\lambda}(\cdot)$ is nondecreasing and continuously differentiable on $[0, \infty)$, and the function $-\beta - p'_{\lambda}(\beta)$ is strictly unimodal on $[0, \infty)$.

- (a) **[Sparsity]** Show that if $t_0 = \min_{\beta \geq 0} \{\beta + p'_{\lambda}(\beta)\} > 0$, then $\hat{\beta}_{\lambda} = 0$ when $|\hat{\beta}| \leq t_0$.
 (b) **[Unbiasedness]** Show that if $p'_{\lambda}(\beta) = 0$ for $|\beta| > t_1$, then $\hat{\beta}_{\lambda} = \hat{\beta}$ when $|\hat{\beta}| \leq t_1$ for large t_1 .

3. Given data $\{(X_i, Y_i), i = 1, 2, \dots, n\}$ where $X_i \in \mathbb{R}^1$, consider the model

$$Y = f(X) + \epsilon,$$

where ϵ is a Gaussian noise with zero-mean and variance σ^2 and

$$f(x) = \sum_{i=-D \cdot 2^D}^{D \cdot 2^D} \theta_i \phi_{\xi_i}(x), \quad \text{with } \phi_{\xi_i}(x) = e^{-(x-\xi_i)^2}.$$

That is, we model Y as the sum of $2 \cdot D \cdot 2^D$ basis functions $\phi_{\xi_i}(\cdot)$, which are placed equally in the interval $[-D, D]$ with gaps 2^{-D} . Apparently the number of parameters is huge and we will definitely overfit the model. We learn the parameters using ridge regression with $\lambda = \sigma^2 \cdot 2^D$, that is, the cost function is given by

$$J(\boldsymbol{\theta}) = \|\mathbf{X}\boldsymbol{\theta} - \mathbf{Y}\|^2 + \sigma^2 2^D \|\boldsymbol{\theta}\|^2.$$

Let $D \rightarrow \infty$, show that the prediction can be computed as

$$\hat{f}(x) = \mathbf{k}(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{Y},$$

where $\mathbf{k} = (k_1, \dots, k_N)$, $\mathbf{K} = \{K_{ij}\}_{0 \leq i, j \leq N}$ with

$$k_i = k(x, X_i), \quad k_{i,j} = k(X_i, X_j)$$

and

$$k(x, y) = \sqrt{\frac{\pi}{2}} e^{-(x-y)^2/2}.$$