

The Elements of Statistical Learning - Chapter 3 Exercises

Exercise 3.1

Show that the F statistic (3.13) for dropping a single coefficient from a model is equal to the square of the corresponding z -score (3.12).

Solution

Without loss of generality, assume that the smaller model has had the final feature x_p removed. Let $\hat{\mathbf{y}}$ denote the least squares approximation for the larger model and $\hat{\mathbf{y}}^*$ that for the smaller. We need to show that

$$\frac{RSS_0 - RSS_1}{RSS_1/(N - p - 1)} = \frac{\hat{\beta}_p^2}{\hat{\sigma}^2 v_p}$$

where v_j is the j th diagonal element of $(\mathbf{X}^T \mathbf{X})^{-1}$.

First note that by definition $\hat{\sigma}^2 = RSS_0/(N - p - 1)$. Moreover, since $\hat{\mathbf{y}}$ is the projection of \mathbf{y} onto the column space of \mathbf{X} , their difference is orthogonal to any element of the columns space. In particular, it is orthogonal to $\hat{\mathbf{y}} - \mathbf{y}^*$, so

$$\begin{aligned} \|\mathbf{y} - \hat{\mathbf{y}}\|^2 &= \|\mathbf{y} - \hat{\mathbf{y}}\| + \|\mathbf{y} - \mathbf{y}^*\|^2 \\ \Rightarrow RSS_0 - RSS_1 &= \|\mathbf{y} - \hat{\mathbf{y}}\|^2 - \|\mathbf{y} - \mathbf{y}^*\|^2 = \|\hat{\mathbf{y}} - \mathbf{y}^*\|^2. \end{aligned}$$

Now let $\mathbf{x}_0, \dots, \mathbf{x}_p$ denote the orthogonal basis of the column space of \mathbf{X} obtained from $\mathbf{x}_0, \dots, \mathbf{x}_p$ using the Gram-Schmidt process (Algorithm 3.1). The least squares estimates $\hat{\mathbf{y}}$ and \mathbf{y}^* are the projections of \mathbf{y} onto the column space of \mathbf{X} and $\text{span}\{\mathbf{x}_j \mid 0 \leq j \leq p - 1\} = \text{span}\{\mathbf{x}_j \mid 0 \leq j \leq p - 1\}$ respectively. Since the x_i are orthogonal, this implies that

$$\hat{\mathbf{y}} - \mathbf{y}^* = \frac{\langle \mathbf{x}_p, \mathbf{y} \rangle}{\langle \mathbf{x}_p, \mathbf{x}_p \rangle} \mathbf{x}_p = \hat{\beta}_p \mathbf{x}_p.$$

Putting these elements together, it just remains to show that $v_p = \|\mathbf{x}_p\|^{-2}$. But, if $\mathbf{X} = \mathbf{QR}$ is the QR-decomposition of \mathbf{X} then $(\mathbf{X}^T \mathbf{X})^{-1} = \mathbf{R}^{-1}(\mathbf{R}^T)^{-1}$. Since \mathbf{R} is upper-triangular, the p th diagonal element of \mathbf{R}^{-1} is $R_{pp}^{-1} = \|\mathbf{x}_p\|^{-1}$ and the claim follows.

Exercise 3.2

Given data on two variables X and Y , consider fitting a cubic polynomial regression model $f(X) = \sum_{j=0}^3 \beta_j X^j$. In addition to plotting the fitted curve, you would like a 95% confidence band about the curve. Consider the following two approaches:

- At each point x_0 , form a 95% confidence interval for the linear function $a^T \hat{\beta} = \sum_{j=0}^3 \hat{\beta}_j x_0^j$.
- Form a 95% confidence interval for β as in (3.15), which in turn generates confidence intervals for $f(x_0)$.

How do these approaches differ? Which band is likely to be wider? Conduct a small simulation experiment to compare the two methods.

Solution

Construction of Confidence Intervals

Let $p = 3$ and $\alpha = 0.05$ and write $\mathbf{x} = (1, x_0, x_0^2, x_0^3)$.

- We have $\hat{\beta} \sim \mathcal{N}(\hat{\beta}, \hat{\sigma}^2 \mathbf{I})$, so

$$\mathbf{x}^T (\hat{\beta} - \beta) \sim \mathcal{N}(0, \mathbf{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x} \hat{\sigma}^2).$$

Let $\mathbf{v} = \mathbf{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x} \in \mathbf{R}$. Then

$$\frac{\mathbf{x}^T (\hat{\beta} - \beta)}{\hat{\sigma} \sqrt{v}} \sim t_{N-p-1},$$

where $\hat{\sigma}$ is the unbiased estimate for σ on p.47. Therefore, a $100(1 - \alpha)\%$ confidence interval for $f(x_0) = \mathbf{x}^T \beta$ has endpoints

$$\hat{f}(x_0) \pm \hat{\sigma} \sqrt{\mathbf{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}} t_{N-p-1, \alpha/2}$$

where $\hat{f}(x_0) = \mathbf{x}^T \hat{\beta}$ and $t_{N-p-1, \alpha/2}$ is the $\frac{1}{2}$ th percentile of a T distribution with $N - p - 1$ degrees of freedom.

- By the argument on p.49, an approximate 95% confidence set for $f(x_0)$ is the set of $\mathbf{x}^T \beta$ such that $\hat{\beta}$ lies in $C = \{ \hat{\beta} \mid \|\mathbf{X}(\hat{\beta} - \hat{\beta})\|^2 \leq \hat{\sigma}^2 \chi^2_p \},$

where $\chi^2_p = \chi^2_{p+1, \alpha}$ is the α th percentile of a chi-squared distribution with $p + 1$ degrees of freedom.

First note that C is an ellipsoid in \mathbf{R}^{p+1} . This implies that the restriction of the linear function $\mathbf{x}^T \beta$ to C achieves its maximum and minimum on the boundary ∂C of C and takes every value in between. In particular, $\{\mathbf{x}^T \beta \mid \beta \in C\}$ is an interval and the endpoints of this interval are the maximum and minimum of $\mathbf{x}^T \beta$ subject to the constraint $\hat{\beta} \in \partial C$, or equivalently

$$\|\mathbf{X}(\hat{\beta} - \hat{\beta})\|^2 = \hat{\sigma}^2 \chi^2_p.$$

We solve this problem using Lagrange multipliers. Let

$$\mathcal{L}(\beta, \lambda) = \mathbf{x}^T \beta - \lambda (\|\mathbf{X}(\hat{\beta} - \hat{\beta})\|^2 - \hat{\sigma}^2 \chi^2_p).$$

This has gradient $\nabla \mathcal{L} = (\frac{\partial \mathcal{L}}{\partial \beta}, \frac{\partial \mathcal{L}}{\partial \lambda})$ with

$$\frac{\partial \mathcal{L}}{\partial \beta} = \mathbf{x} - 2\lambda \mathbf{X}^T (\mathbf{X} \hat{\beta} - \mathbf{X} \hat{\beta}), \quad \frac{\partial \mathcal{L}}{\partial \lambda} = \|\mathbf{X}(\hat{\beta} - \hat{\beta})\|^2 - \hat{\sigma}^2 \chi^2_p.$$

Any solution to our optimisation problem will have $\nabla \mathcal{L} = 0$. Setting the first partial derivative equation to zero gives

$$\mathbf{X}^T \mathbf{X}(\hat{\beta} - \hat{\beta}) = \frac{1}{2\lambda} \mathbf{x} \Rightarrow \hat{\beta} - \hat{\beta} = \frac{1}{2\lambda} \mathbf{X}^T \mathbf{x}.$$

Setting $\frac{\partial \mathcal{L}}{\partial \lambda} = 0$ and substituting this in give

$$\left\| \frac{1}{2\lambda} \mathbf{X}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x} \right\|^2 = \hat{\sigma}^2 \chi^2_p \Rightarrow \frac{1}{2\lambda} = \pm \frac{\hat{\sigma} \chi_p}{\|\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}\|}.$$

So, since we know the optimisation problem has a maximum and a minimum, they must occur at

$$\hat{\beta} \pm \hat{\sigma} \frac{\mathbf{X}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}}{\|\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}\|} \chi_p.$$

Therefore, the confidence interval for $f(x_0)$ has endpoints

$$\hat{f}(x_0) \pm \hat{\sigma} \frac{\mathbf{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}}{\|\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}\|} \chi_{p+1, \alpha}.$$

Comparison

The key difference is that the first method gives a $100(1 - \alpha)\%$ confidence interval at a particular x_0 . That is, for a fixed x_0 the probability that $f(x_0)$ lies inside the confidence bands is $1 - \alpha$. The second method starts with a $100(1 - \alpha)\%$ confidence set for β and so is valid for all x_0 simultaneously. More precisely, there is a probability of $1 - \alpha$ that $\hat{\beta}$ lies in the confidence set and thus $f(x_0)$ lies in the confidence interval for all x_0 .

A less important distinction is that first method is exact (given the assumptions), whereas the second uses the approximation of t_{N-p-1} by the standard normal distribution $p + 1$ times. The T distribution is less peaked than the standard normal, $t_{\epsilon, \alpha} > z_\alpha$ for $\alpha < 0.5$, but for N sufficiently large the difference will be insignificant.

Both of these differences imply that the chi-squared confidence interval should be wider.

Simulation

The following function from a custom module calculates and prints the confidence bands from the two methods against the regression function $f(x_0)$. We use a simulated data set of size N with X sampled from a normal distribution with mean `xmean` and standard deviation `xstdev`, and show the $100(1 - \alpha)\%$ confidence bands for x_0 in the range `plot_range` (`=xmean ± 2 xstdev` by default).

```
In [1]: from confband import run_sim

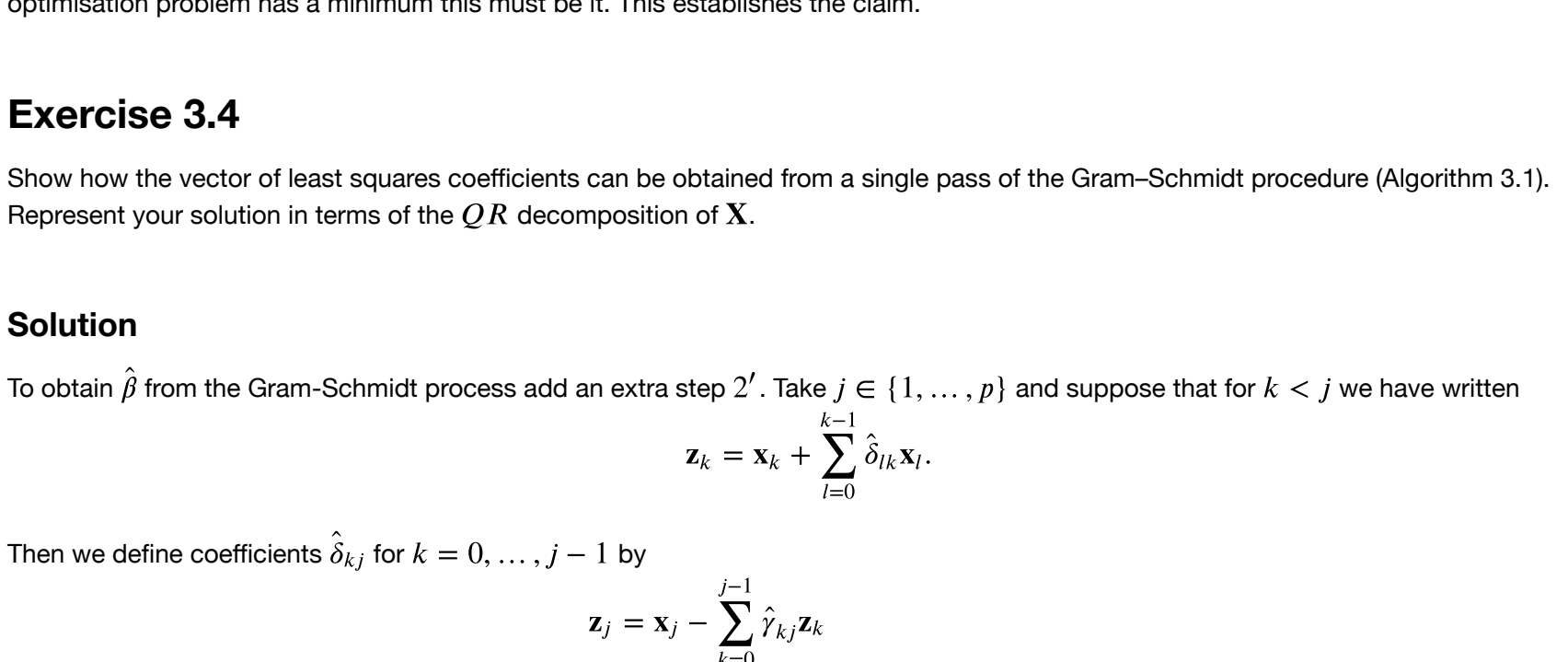
beta = [-1, -3, 1, 1]
N = 100
sigma = 1
xmean = 0
xstdev = 0.5
alpha = 0.05
plot_range = None

run_sim(beta, alpha, N, sigma, xmean, xstdev, plot_range)
```

```
Beta: [-1, -3, 1, 1]
Betahat: [-1.055 -2.678 1.113 0.949]
```

```
RSS: 90.9
Standard error: 0.973
```

```
Endpoints of 95.0% confidence interval at xmean:
T method: [-1.23, -0.921]
Chi-square method: [-1.419, -0.692]
```



Exercise 3.3

Gauss-Markov theorem:

- Prove the Gauss-Markov theorem: the least squares estimate of a parameter $a^T \hat{\beta}$ has variance no bigger than that of any other linear unbiased estimate of $a^T \hat{\beta}$ (Section 3.2.2).
- The matrix inequality $\mathbf{B} \leq \mathbf{A}$ holds if $\mathbf{A} - \mathbf{B}$ is positive semidefinite. Show that if $\hat{\mathbf{V}}$ is the variance-covariance matrix of the least squares estimate of $\hat{\beta}$ and $\hat{\mathbf{V}}$ is the variance-covariance matrix of any other linear unbiased estimate, then $\hat{\mathbf{V}} \leq \hat{\mathbf{V}}$.

Solution

- Let $\hat{\theta} = a^T \hat{\beta}$. The least squares estimator $\hat{\beta}$ of β satisfies $\hat{\beta} \sim \mathcal{N}(\mathbf{X}(\hat{\beta}, \sigma^2) (\mathbf{X}^T \mathbf{X})^{-1})$ so the least squares estimate $\hat{\theta} = a^T \hat{\beta}$ of θ has $\text{Var}(\hat{\theta}) = a^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X} a$.

Suppose $\hat{\theta} = c^T \mathbf{y}$ is a linear unbiased estimator of θ . Then

$$\mathbf{y} \sim \mathcal{N}(\mathbf{X} \hat{\beta}, \sigma^2 \mathbf{I}_N) \Rightarrow \hat{\theta} \sim \mathcal{N}(c^T \mathbf{X} \hat{\beta}, \sigma^2 c^T c).$$

Since $\hat{\theta}$ is unbiased for θ ,

$$\mathbb{E}(\hat{\theta}) = \theta \Rightarrow c^T \mathbf{X} \hat{\beta} = a^T \hat{\beta} \Rightarrow (c^T \mathbf{X} - a^T) \hat{\beta} = 0.$$

But this must hold for every $\hat{\beta} \in \mathbf{R}^{p+1}$, so $c^T \mathbf{X} = a^T$.

We have reduced the problem to showing that if $c^T \mathbf{X} = a^T$ then $c^T c \geq a^T (\mathbf{X}^T \mathbf{X})^{-1} a$. We will establish this by showing that $a^T (\mathbf{X}^T \mathbf{X})^{-1} a$ is a global minimum for $\|c\|^2$ subject to the constraint $c^T \mathbf{X} = a^T$.

First, observe that a minimum exists since the solution set is non-empty (consider $c = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T a$), bounded below by 0, and closed. We will find it using Lagrange multipliers.

Let $\lambda \in \mathbf{R}^{p+1}$ be variables and define

$$\mathcal{L}(c, \lambda) = \|c\|^2 = \mathbf{X}^T (\mathbf{X} c - a).$$

This satisfies

$$\frac{\partial \mathcal{L}}{\partial c} = 2c - \mathbf{X} \lambda, \quad \frac{\partial \mathcal{L}}{\partial \lambda} = \mathbf{X}^T c - a.$$

At an extremum of our optimisation problem both of these are zero, so

$$\begin{aligned} c &= \frac{1}{2} \mathbf{X} \lambda \Rightarrow \frac{1}{2} \mathbf{X}^T \mathbf{X} \lambda = a \\ &\Rightarrow \lambda = 2(\mathbf{X}^T \mathbf{X})^{-1} a \\ &\Rightarrow c = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} a \\ &\Rightarrow \hat{\theta} = c^T \mathbf{y} = \hat{\theta}. \end{aligned}$$

Since the Lagrange multiplier has a unique solution, this must be the global minimum of the constrained optimisation problem. In fact we have proved something slightly stronger: that $\hat{\theta}$ is the unique unbiased estimator with minimal variance.

- Our proof is analogous to that in part (a). The covariance matrix of $\hat{\beta}$ is $\text{Var}(\hat{\beta}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$. If $\hat{\beta} = \mathbf{A} \mathbf{y}$ is another unbiased estimator for $\hat{\beta}$ then $\hat{\beta} \sim \mathcal{N}(\mathbf{A} \mathbf{X} \hat{\beta}, \sigma^2 \mathbf{A} \mathbf{A}^T)$. Since $\hat{\beta}$ is unbiased, $\mathbf{A} \mathbf{X} \hat{\beta} = \hat{\beta}$ for all $\hat{\beta} \in \mathbf{R}^{p+1}$ and so $\mathbf{A} \mathbf{X} = \mathbf{I}_{p+1}$. Thus we have reduced the problem to showing that if \mathbf{A} is a $(p + 1) \times N$ matrix with $\mathbf{A} \mathbf{X} = \mathbf{I}_{p+1}$ then $\mathbf{A} \mathbf{A}^T - (\mathbf{X}^T \mathbf{X})^{-1}$ is positive semi-definite, or equivalently $\mathbf{v}^T \mathbf{A} \mathbf{A}^T \mathbf{v} \geq \mathbf{v}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{v}$ for all $\mathbf{v} \in \mathbf{R}^{p+1} \setminus \{0\}$.

Again we treat this as a constrained optimisation problem and employ Lagrange multipliers. Fix $\mathbf{v} \in \mathbf{R}^{p+1}$. The set $\{\|\mathbf{A}^T \mathbf{v}\|^2 \mid \mathbf{A} \mathbf{X} = \mathbf{I}_{p+1}\}$ is non-empty (take $\mathbf{A} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$), bounded below by zero, and closed so has a minimum.

Let \mathbf{A} be a $(p + 1) \times (p + 1)$ matrix of variables and define

$$\mathcal{L}(\mathbf{A}, \lambda) = \|\mathbf{A}^T \mathbf{v}\|^2 - \lambda \cdot (\mathbf{A} \mathbf{X} - \mathbf{I}_{p+1}),$$

where \cdot denotes the dot product. This has

$$\frac{\partial \mathcal{L}}{\partial \mathbf{A}} = 2\mathbf{v} \mathbf{v}^T \mathbf{A} - \mathbf{A} \mathbf{X}^T, \quad \frac{\partial \mathcal{L}}{\partial \lambda} = \mathbf{A} \mathbf{X} - \mathbf{I}_{p+1}.$$

At an extremum of the constrained optimisation problem both of these are zero, so

$$\begin{aligned} 2\mathbf{v} \mathbf{v}^T \mathbf{A} \mathbf{X} &= \mathbf{A} \mathbf{X}^T \mathbf{X} \Rightarrow 2\mathbf{v} \mathbf{v}^T = \mathbf{A} \mathbf{X}^T \mathbf{X} \\ &\Rightarrow c = \mathbf{v} \mathbf{v}^T (\mathbf{X}^T \mathbf{X})^{-1} \\ &\Rightarrow \mathbf{v} \mathbf{v}^T \mathbf{A} = \mathbf{v} \mathbf{v}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T. \end{aligned}$$

Multiplying each term on the right by its transpose yields

$$\begin{aligned} \mathbf{v} \mathbf{v}^T \mathbf{A} \mathbf{A}^T \mathbf{v} \mathbf{v}^T &= \mathbf{v} \mathbf{v}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{v} \mathbf{v}^T \\ [\mathbf{v}^T \mathbf{A} \mathbf{A}^T \mathbf{v}] \mathbf{v} \mathbf{v}^T &= [\mathbf{v}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{v}] \mathbf{v} \mathbf{v}^T \\ \mathbf{v}^T \mathbf{A} \mathbf{A}^T \mathbf{v} &= \mathbf{v}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{v}, \end{aligned}$$

where the second line holds as both terms in square brackets are scalars and so commute with everything. Since the constrained optimisation problem has a minimum this must be it. This establishes the claim.

Exercise 3.4

Show how the vector of least squares coefficients can be obtained from a single pass of the Gram-Schmidt procedure (Algorithm 3.1). Represent your solution in terms of the QR decomposition of \mathbf{X} .

Solution

To obtain $\hat{\beta}$ from the Gram-Schmidt process add an extra step $2'$. Take $j \in \{1, \dots, p\}$ and suppose that for $k < j$ we have written

$$\mathbf{z}_k = \mathbf{x}_k + \sum_{i=0}^{k-1} \hat{\gamma}_{ki} \mathbf{z}_i.$$

Then we define coefficients $\hat{\delta}_{kj}$ for $k = 0, \dots, j - 1$ by

$$\begin{aligned} \mathbf{z}_j &= \mathbf{x}_j - \sum_{k=0}^{j-1} \hat{\gamma}_{kj} \mathbf{z}_k \\ &= \mathbf{x}_j - \sum_{k=0}^{j-1} \left(\mathbf{x}_k + \sum_{i=0}^k \hat{\delta}_{ki} \mathbf{x}_i \right) \\ &= \mathbf{x}_j + \sum_{i=0}^{j-1} \hat{\delta}_{ji} \mathbf{x}_i. \end{aligned}$$

Let Δ be the matrix with $\Delta_{ij} = \hat{\delta}_{ji}$ for $i < j$, ones on the diagonal, and zeros elsewhere. By construction,

$$\mathbf{Z} = \mathbf{X} \Delta \Rightarrow \Delta = \mathbf{I}^{-1}.$$

By (3.32), $\hat{\beta} = \mathbf{I}^{-1} \mathbf{Z}^T \mathbf{y} = \Delta \mathbf{Z}^T \mathbf{y}$ so we can calculate $\hat{\beta}$ explicitly:

$$\hat{\beta}_j = \mathbf{z}_j + \sum_{k=j+1}^p \hat{\delta}_{jk} \mathbf{z}_k.$$

Exercise 3.5

Consider the ridge regression problem (3.41). Show that this problem is equivalent to the problem

$$\hat{\beta}^r = \text{argmin}_{\beta} \left\{ \sum_{j=1}^N \left[y_j - \hat{\beta}_j^r - \sum_{j=1}^p (x_{ij} - \hat{\beta}_j^r)^2 \right]^2 + \sum_{j=1}^p (\hat{\beta}_j^r)^2 \right\}.$$

Give the correspondence between $\hat{\beta}^r$ and the original $\hat{\beta}$ in (3.41). Characterize the solution to this modified criterion. Show that a similar result holds for the lasso.

Solution

Our solution is the same for both ridge regression and lasso. Since

$$\hat{\beta}_i^r = \sum_{j=1}^p (x_{ij} - \hat{\beta}_j^r) \hat{\beta}_i^r = \left(\hat{\beta}_i^r - \sum_{j=1}^p x_{ij} \hat{\beta}_j^r \right) + \sum_{j=1}^p x_{ij} \hat{\beta}_j^r,$$

the two minimisation problems are equivalent with

$$\hat{\beta}_i^r = \hat{\beta}_i + \sum_{j=1}^p \hat{\gamma}_{ij} \hat{\beta}_j = \frac{1}{N} \sum_{i=1}^N \hat{\gamma}_{ij} \hat{\beta}_j$$

and $\hat{\beta}_j^r = \hat{\beta}_j$ for $j \neq 0$.

At the minimum, the derivative with respect to $\hat{\beta}_i^r$ of the expression in braces is zero, so

$$\sum_{j=1}^N \left[y_j - \hat{\beta}_i^r - \sum_{j=1}^p (x_{ij} - \hat{\beta}_j^r) \hat{\beta}_i^r \right] = 0 \Rightarrow \left(\sum_{j=1}^N \hat{\gamma}_{ij} \right) - N \hat{\beta}_i^r = 0 = 0$$

and thus the solution to the modified criterion is $\hat{\beta}_0^r = \bar{y}$, $\hat{\beta}_j^r = \hat{\beta}_j$ for $j > 0$.

Exercise 3.6

Show that the ridge regression estimate is the mean (and mode) of the posterior distribution, under a Gaussian prior $\beta \sim \mathcal{N}(0, \tau^2 \mathbf{I})$, and Gaussian sampling model $y \sim \mathcal{N}(\mathbf{X} \beta, \sigma^2 \mathbf{I})$. Find the relationship between the regularization parameter λ in the ridge formula, and the variances τ^2 and σ^2 .

Solution

The posterior distribution for β has density function

$$\begin{aligned} f_{\beta|y, \dots, y_p}(\beta) &\propto f_{y_1, \dots, y_p|y_1, \dots, y_p}(\beta) f_{\beta}(\beta) \\ &= \left(\prod_{i=1}^N f_{y_i|\beta}(y_i) \right) f_{\beta}(\beta) \\ &\propto \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mathbf{x}_i^T \beta)^2 - \frac{1}{2\tau^2} \|\beta\|^2 \right) \\ &= \exp \left(-\frac{1}{2\sigma^2} \left(\|\mathbf{y} - \mathbf{X} \beta\|^2 + \frac{\sigma^2}{\tau^2} \|\beta\|^2 \right) \right). \end{aligned}$$

So the mode of this distribution is

$$\text{argmax}_{\beta} (f_{\beta|y, \dots, y_p}(\beta)) = \text{argmin}_{\beta} \left(\|\mathbf{y} - \mathbf{X} \beta\|^2 + \frac{\sigma^2}{\tau^2} \|\beta\|^2 \right).$$

This is the ridge regression estimate with $\lambda = \sigma^2/\tau^2$. Moreover, since the posterior is Gaussian (p.64) the mean equals the mode.

Exercise 3.7

Assume $y_i \sim \mathcal{N}(\hat{\beta}_0 + x_i^T \hat{\beta}, \sigma^2)$, $i = 1, 2, \dots, N$, and the parameters $\hat{\beta}_j$, $j = 1, \dots, p$, are each distributed as $\mathcal{N}(0, \tau^2)$, the (independently) of one another. Assuming σ^2 and τ^2 are known, and $\hat{\beta}_0$ is not governed by a prior (or has a flat improper prior), show that the (marginally) log-posterior density of β is proportional to $\sum_{i=1}^N (y_i - \hat{\beta}_0 - \sum_{j=1}^p x_{ij} \hat{\beta}_j)^2 + \lambda \sum_{j=1}^p \hat{\beta}_j^2$ where $\lambda = \sigma^2/\tau^2$.

Solution

This follows from the solution to exercise 3.6.

Exercise 3.8

Consider the QR decomposition of the uncentered $N \times (p + 1)$ matrix \mathbf{X} (whose first column is all ones), and the SVD of the $N \times p$ centered matrix $\tilde{\mathbf{X}}$. Show that Q_2 and U span the same subspace, where Q_2 is the sub-matrix of Q with the first column removed. Under what circumstances will they be the same, up to sign flips?

Solution

For $j = 1, \dots, p$, the Gram-Schmidt algorithm implies

$$\begin{aligned} \mathbf{z}_j &= \mathbf{x}_j + \sum_{i=0}^{j-1} \hat{\gamma}_{ji} \mathbf{z}_i \\ &\Rightarrow \mathbf{z}_j - \hat{\gamma}_{j0} \mathbf{z}_$$

Exercise 3.18

Read about conjugate gradient algorithms (Murray et al., 1981, for example), and establish a connection between these algorithms and partial least squares.

Solution

TODO

Exercise 3.19

Show that $\|\hat{\beta}^{ridge}\|$ increases as its tuning parameter $\lambda \rightarrow 0$. Does the same property hold for the lasso and partial least squares estimates? For the latter, consider the “tuning parameter” to be the successive steps in the algorithm.

Solution

Ridge: Let

$$R(\hat{\beta}; \lambda) = \|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2 + \lambda \|\hat{\beta}\|^2.$$

Take $\lambda > \lambda' > 0$ and let $\hat{\beta}$ and $\hat{\beta}'$ be the ridge approximations with tuning parameters λ and λ' , respectively (so they minimise $R(\hat{\beta}; \lambda)$ and $R(\hat{\beta}; \lambda')$). Then by definition

$$R(\hat{\beta}; \lambda) \leq R(\hat{\beta}'; \lambda) \quad \text{and} \quad R(\hat{\beta}'; \lambda') \leq R(\hat{\beta}; \lambda').$$

Subtract the second inequality from the first yields

$$\begin{aligned} R(\hat{\beta}; \lambda) - R(\hat{\beta}; \lambda') &\leq R(\hat{\beta}'; \lambda) - R(\hat{\beta}'; \lambda') \\ (\lambda - \lambda') \|\hat{\beta}\|^2 &\leq (\lambda - \lambda') \|\hat{\beta}'\|^2 \\ \|\hat{\beta}\| &\leq \|\hat{\beta}'\|, \end{aligned}$$

so the ridge approximation increases in L^2 norm as λ decreases.

Lasso: The same proof as above shows that the lasso approximation to $\hat{\beta}$ increases in L^1 norm as λ decreases. However, this doesn't necessarily mean that it increases in L^2 norm. For example suppose that $\beta_j = 2$ and assume the least squares approximation has $\hat{\beta}_1, \hat{\beta}_2 > 0$. If lasso modification of the LAR path went from $(0, 0)$ to $(\hat{\beta}_1 + \hat{\beta}_2, 0)$ to $(\hat{\beta}_1, \hat{\beta}_2)$ then the L^1 norm of the coefficient vector would always be increasing as λ decreases, but the L^2 norm would decrease on the last segment.

It remains to give an example to show this behaviour is possible. Take $p = 2$ and $N = 3$ - the minimal values for an example. Take

$$\mathbf{x}_1 = \begin{pmatrix} -\sqrt{\frac{1}{2}} \\ 0 \\ \sqrt{\frac{1}{2}} \end{pmatrix}, \quad \mathbf{x}_2 = \begin{pmatrix} -\sqrt{\frac{4}{3}} \\ -\sqrt{\frac{1}{3}} \\ \sqrt{\frac{1}{3}} \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix}.$$

Note that all three vectors have mean zero and the first two have unit variance. Since

$$\langle \mathbf{x}_1, \mathbf{y} \rangle = \sqrt{6} > \sqrt{2} = \langle \mathbf{x}_2, \mathbf{y} \rangle,$$

the first step is in the direction of \mathbf{x}_1 . This ends at $\frac{\sqrt{2}}{11}(\sqrt{3} + 5)\mathbf{x}_1$ since at this point the residual \mathbf{r} has equal correlation with \mathbf{x}_1 and \mathbf{x}_2 ($\langle \mathbf{x}_j, \mathbf{r} \rangle = 5\sqrt{6}$ for $j = 1, 2$). We then continue in a straight line in direction $\begin{pmatrix} -1 \\ 1 \\ 5 \end{pmatrix}$ to the least squares solution

$$\hat{\beta} = \frac{\sqrt{2}}{11} \begin{pmatrix} \sqrt{3} \\ 5 \end{pmatrix}.$$

So although the coefficient path increases in L^1 norm it decreases in L^2 norm along the final stretch.

PLS: Let $\hat{\beta}^{(m)}$ denote the PLS approximation after m steps. We claim that $\|\hat{\beta}^{(m)}\|$ increases as m increases. For $m = 0, \dots, p$ let $\mathbf{X}^{(m)}$ denote the matrix whose j th column is $\mathbf{x}_j^{(m)}$. The proof relies heavily on the following three facts:

- $\mathbf{z}_m = \mathbf{X}^{(m-1)} \hat{\phi}_m$;
- $\hat{\phi}_m = (\mathbf{X}^{(m-1)})^T \mathbf{y}$;
- $\mathbf{X} \hat{\phi}_m = \mathbf{X}^{(r-1)} \hat{\phi}_m$ for any $l \leq m$.

The first two of these are by definition. The last was essentially proved in the course of Exercise 3.15, but we quickly recall the proof here. By the equivalent characterisation of $\hat{\phi}_m$ in Exercise 3.15, $\mathbf{X} \hat{\phi}_m$ is orthogonal to $\mathbf{x}_1, \dots, \mathbf{x}_{l-1}$. So if π is the projection onto the subspace they span then $(1 - \pi)(\mathbf{X} \hat{\phi}_m) = \mathbf{X} \hat{\phi}_m$. But by construction $(1 - \pi)(\mathbf{x}_l) = \mathbf{x}_l^{(l-1)}$ so

$$\mathbf{X} \hat{\phi}_m = (1 - \pi)(\mathbf{X} \hat{\phi}_m) = (1 - \pi) \left(\sum_{j=1}^m \hat{\phi}_{m,j} \mathbf{x}_j \right) = \sum_{j=1}^m \hat{\phi}_{m,j} \mathbf{x}_j^{(l-1)} = \mathbf{X}^{(l-1)} \hat{\phi}_m.$$

Now we continue to the proof. First note that

$$\mathbf{y}^{(m)} = \sum_{i=1}^m \hat{\theta}_i \alpha_i = \sum_{i=1}^m \hat{\theta}_i \mathbf{X}^{(i-1)} \hat{\phi}_i = \sum_{i=1}^m \hat{\theta}_i \mathbf{X} \hat{\phi}_i = \mathbf{X} \left(\sum_{i=1}^m \hat{\theta}_i \hat{\phi}_i \right),$$

so $\hat{\beta}^{(m)} = \sum_{i=1}^m \hat{\theta}_i \hat{\phi}_i$. We claim that $\hat{\theta}_i \geq 0$ and $\langle \hat{\phi}_i, \hat{\phi}_m \rangle \geq 0$ for all $i \leq m$. Since

$$\|\hat{\beta}^{(m)}\|^2 = \|\hat{\beta}^{(m-1)}\|^2 + \sum_{i=1}^{m-1} \hat{\theta}_i \hat{\theta}_m \langle \hat{\phi}_i, \hat{\phi}_m \rangle + \hat{\theta}_m^2 \|\hat{\phi}_m\|^2$$

this will establish that $\|\hat{\beta}^{(m)}\|$ increases with m .

The first claim holds because by definition $\hat{\theta}_i = \langle \mathbf{y}, \mathbf{x}_i \rangle / \|\mathbf{x}_i\|^2$ and

$$\begin{aligned} \langle \mathbf{y}, \mathbf{x}_i \rangle &= \langle \mathbf{y}, \mathbf{X}^{(i-1)} \hat{\phi}_i \rangle \\ &= \langle \mathbf{y}, \mathbf{X}^{(i-1)} (\mathbf{X}^{(i-1)})^T \mathbf{y} \rangle \\ &= \langle (\mathbf{X}^{(i-1)})^T \mathbf{y}, (\mathbf{X}^{(i-1)})^T \mathbf{y} \rangle \\ &= \|\hat{\phi}_i\|^2 \geq 0. \end{aligned}$$

For the second claim observe that if $l \leq m$ then

$$\begin{aligned} \langle \hat{\phi}_i, \hat{\phi}_m \rangle &= \langle (\mathbf{X}^{(i-1)})^T \mathbf{y}, \hat{\phi}_m \rangle \\ &= \langle \mathbf{y}, \mathbf{X}^{(i-1)} \hat{\phi}_m \rangle \\ &= \langle \mathbf{y}, \mathbf{X}^{(m-1)} \hat{\phi}_m \rangle \\ &= \langle (\mathbf{X}^{(m-1)})^T \mathbf{y}, \hat{\phi}_m \rangle \\ &= \|\hat{\phi}_m\|^2 \geq 0 \end{aligned}$$

as required.

Exercise 3.20

Consider the canonical-correlation problem (3.67). Show that the leading pair of canonical variates u_1 and v_1 solve the problem

$$\max_{\substack{c^T(\mathbf{X}^T \mathbf{X})c=1 \\ c^T(\mathbf{Y}^T \mathbf{Y})c=1}} u^T(\mathbf{Y}^T \mathbf{X})c,$$

a generalised SVD problem. Show that the solution is given by $u_1 = (\mathbf{Y}^T \mathbf{X})(\mathbf{X}^T \mathbf{X})^{-\frac{1}{2}} u_1^*$ and $v_1 = (\mathbf{X}^T \mathbf{X})^{-\frac{1}{2}} v_1^*$, where u_1^* and v_1^* are the leading left and right singular vectors in

$$(\mathbf{Y}^T \mathbf{Y})^{-\frac{1}{2}} (\mathbf{Y}^T \mathbf{X})(\mathbf{X}^T \mathbf{X})^{-\frac{1}{2}} = \mathbf{U}^* \mathbf{D}^* \mathbf{V}^{*T}.$$

Show that the entire sequence $u_m, v_m, m = 1, \dots, \min(K, p)$ is also given by (3.67).

Solution

The leading pair of canonical variates u_1 and v_1 must maximise:

$$\text{Corr}(\mathbf{Y}u, \mathbf{X}v) = \frac{u^T \mathbf{Y}^T \mathbf{X} v}{(u^T \mathbf{Y}^T \mathbf{Y} u)(v^T \mathbf{X}^T \mathbf{X} v)}.$$

Since any solution vectors u and v can be scaled so that the two quantities in the denominator equal one, the problem is equivalent to the generalised SVD problem given.

Take Cholesky decompositions $\mathbf{X}^T \mathbf{X} = \mathbf{K}^T \mathbf{K}$ and $\mathbf{Y}^T \mathbf{Y} = \mathbf{L}^T \mathbf{L}$ where \mathbf{K} and \mathbf{L} are square and invertible. We will prove a slightly modified version of the problem: $u_1 = \mathbf{L}^{-1} u_1^*$ and $v_1 = \mathbf{K}^{-1} v_1^*$, where the starred vectors are left and right singular vectors of $\mathbf{M} = (\mathbf{L}^T)^{-1} \mathbf{Y}^T \mathbf{X} \mathbf{K}^{-1}$. One can oneself that these two statements are essentially the same.

We begin by performing a change of basis by writing $u^* = \mathbf{L}u$ and $v^* = \mathbf{K}v$. In this basis our goal is to maximise

$$u^{*T} \mathbf{Y}^T \mathbf{X} v^* = (\mathbf{L}u)^T (\mathbf{L}^T)^{-1} \mathbf{Y}^T \mathbf{X} \mathbf{K}^{-1} (\mathbf{K}v) = u^T \mathbf{M} v^*$$

subject to $\|u^*\| = \|v^*\| = 1$.

First note that this problem has a maximum since $u^T \mathbf{M} v^*$ is a continuous function on the compact set $\{u^*, v^* \mid \|u^*\| = \|v^*\| = 1\}$. We will find it using Lagrange multipliers. Let

$$\mathcal{L}(u^*, v^*, \lambda, \mu) = u^T \mathbf{M} v^* - \lambda(u^{*T} u^* - 1) - \mu(v^{*T} v^* - 1).$$

At a maximum of the constrained optimisation problem,

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial u^*} &= \mathbf{M} v^* - 2\lambda u^* = 0 \\ \frac{\partial \mathcal{L}}{\partial v^*} &= \mathbf{M}^T u^* - 2\mu v^* = 0 \\ \frac{\partial \mathcal{L}}{\partial \lambda} &= u^{*T} u^* - 1 = 0 \\ \frac{\partial \mathcal{L}}{\partial \mu} &= v^{*T} v^* - 1 = 0. \end{aligned}$$

Multiplying the first and second equations on the left by u^{*T} and v^{*T} respectively yields

$$u^{*T} \mathbf{M} v^* = 2\lambda \quad \text{and} \quad v^{*T} \mathbf{M}^T u^* = 2\mu,$$

which implies that $\lambda = \mu$. Moreover, this implies that u^* and v^* are a pair of singular vectors for \mathbf{M} with singular value 2λ . To attain the global maximum of $u^{*T} \mathbf{M} v^* = 2\lambda$ they are must be the first singular vectors.

Now take $1 < m \leq \min(K, p)$ and suppose that we have shown that u_1^*, \dots, u_{m-1}^* and v_1^*, \dots, v_{m-1}^* are the $m-1$ first singular vectors of \mathbf{M} . The pair (u_m^*, v_m^*) are the solution to the following constrained optimisation problem: maximise $u_m^{*T} \mathbf{M} v_m^*$ subject to $\|u_m^*\| = \|v_m^*\| = 1$ and

$$u_m^{*T} u_k^* = u_m^{*T} v_k^* = v_m^{*T} v_k^* = v_m^{*T} u_k^* = 0$$

for $k < m$.

We approach this problem much the same as above. First note that there is a maximum since we have a continuous function on a compact set. Now consider the function

$$\begin{aligned} \mathcal{L} &= u_m^{*T} \mathbf{M} v_m^* - \lambda(u_m^{*T} u_m^* - 1) - \mu(v_m^{*T} v_m^* - 1) \\ &\quad + \sum_{k=1}^{m-1} (\alpha_k^X u_m^{*T} u_k^* - \alpha_k^Y u_m^{*T} v_k^* + \alpha_k^Y v_m^{*T} u_k^* - \alpha_k^X v_m^{*T} v_k^*) \end{aligned}$$

Setting $\frac{\partial \mathcal{L}}{\partial u_m^*} = \frac{\partial \mathcal{L}}{\partial v_m^*} = 0$ yields

$$\sum_{k=1}^m \alpha_k^{X^*} \mathbf{M} v_k^* = \sum_{k=1}^m \alpha_k^{Y^*} u_k^* \quad \text{and} \quad \sum_{k=1}^m \alpha_k^{Y^*} \mathbf{M} u_k^* = \sum_{k=1}^m \alpha_k^{X^*} v_k^*,$$

where we have written $\alpha_k^{X^*} = \alpha_k^{X^*} = 1$, $\alpha_k^{Y^*} = 2\lambda$, and $\alpha_k^{X^*} = 2\mu$ for notational convenience. Setting the derivatives with respect to the parameters equal to zero yields the original constraints of the problem.

Multiplying the two equations on the left by u_m^{*T} and v_m^{*T} respectively yields

$$u_m^{*T} \mathbf{M} v_m^* = 2\lambda \quad \text{and} \quad v_m^{*T} \mathbf{M}^T u_m^* = 2\mu,$$

so $\lambda = \mu$. On the other hand, multiplying the equations by u_k^{*T} and v_k^{*T} yields

$$u_k^{*T} \alpha_k^{X^*} u_k^* \mathbf{M} v_k^* = \alpha_k^{Y^*} v_k^* \quad \text{and} \quad \alpha_k^{Y^*} v_k^{*T} \mathbf{M}^T u_k^* = \alpha_k^{X^*} u_k^*$$

so $\alpha_k^{Y^*} = v_k^{*T} \alpha_k^{X^*} u_k^*$ and $\alpha_k^{X^*} = u_k^{*T} \alpha_k^{Y^*} v_k^*$, where $u_k^* = (u_1^*, v_1^*, \dots, u_{m-1}^*, v_{m-1}^*)$ are pairs of singular vectors, the equations from setting the partial derivatives of \mathcal{L} equal to zero yield

$$\mathbf{M} u_m^* = 2\lambda u_m^* \quad \text{and} \quad \mathbf{M}^T v_m^* = 2\lambda v_m^*.$$

So (u_m^*, v_m^*) are a pair of singular vectors for \mathbf{M} . To maximise $u_m^{*T} \mathbf{M} v_m^*$ = 2λ subject to the orthogonality conditions, the must by the m th pair of singular vectors.

Exercise 3.21

Show that the solution to the reduced-rank regression problem (3.68), with Σ estimated by $\mathbf{Y}^T \mathbf{Y}/N$, is given by (3.69). *Hint:* Transform \mathbf{Y} to $\mathbf{Y}^* = \mathbf{Y} \Sigma^{-\frac{1}{2}}$, and solve in terms of the canonical vectors u^* . Show that $\mathbf{U}_m = \Sigma^{-\frac{1}{2}} \mathbf{U}_m^*$, and a generalized inverse is $\mathbf{U}_m^- = \mathbf{U}_m^{*T} \Sigma^{\frac{1}{2}}$.

Exercise 3.22

Exercise 3.23

Consider a regression problem with all variables and response having mean zero and standard deviation one. Suppose also that each variable has identical absolute correlation with the response:

$$\frac{1}{N} |\langle \mathbf{x}_j, \mathbf{y} \rangle| = \lambda, \quad j = 1, \dots, p.$$

Let $\hat{\beta}$ be the least-squares coefficient of \mathbf{y} on \mathbf{X} , and let $\mathbf{u}(\alpha) = \alpha \hat{\beta}$ for $\alpha \in [0, 1]$ be the vector that moves a fraction α towards the least squares fit \mathbf{u} . Let RSS be the residual sum-of-squares from the full least squares fit.

(a) Show that

$$\frac{1}{N} |\langle \mathbf{x}_j, \mathbf{y} - \mathbf{u}(\alpha) \rangle| = \frac{1}{N} |(1 - \alpha) \lambda|, \quad j = 1, \dots, p$$

and hence the correlations of each \mathbf{x}_j with the residuals remain equal in magnitude as we progress toward \mathbf{u} .

(b) Show that these correlations are all equal to

$$\lambda(\alpha) = \frac{1 - \alpha}{\sqrt{(1 - \alpha)^2 + \frac{\alpha^2 - \alpha}{N} \cdot RSS}} \cdot \lambda,$$

and hence they decrease monotonically to zero.

(c) Use these results to show that the LAR algorithm in Section 3.4.4 keeps the correlations tied and monotonically decreasing, as claimed in (3.55).

Solution

(a) Take $j \in \{0, \dots, p\}$ and let \mathbf{e}_j denote the coordinate vector with a 1 in the j th position and 0s elsewhere. Then

$$\begin{aligned} \langle \mathbf{x}_j, \alpha \hat{\beta} \rangle &= \langle \mathbf{x}_j, \alpha \mathbf{X} \hat{\beta} \rangle \\ &= \alpha \langle \mathbf{X} \mathbf{e}_j, \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \rangle \\ &= \alpha \langle \mathbf{e}_j, \mathbf{X}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \rangle \\ &= \alpha \langle \mathbf{e}_j, \mathbf{X}^T \mathbf{y} \rangle \\ &= \alpha \langle \mathbf{X} \mathbf{e}_j, \mathbf{y} \rangle \\ &= \alpha \langle \mathbf{x}_j, \mathbf{y} \rangle \end{aligned}$$

so

$$\frac{1}{N} |\langle \mathbf{x}_j, \mathbf{y} - \mathbf{u}(\alpha) \rangle| = \frac{1}{N} |(1 - \alpha) \langle \mathbf{x}_j, \mathbf{y} \rangle| = (1 - \alpha) \lambda$$

as required.

(b) It will be useful in part (c) to drop the assumption that \mathbf{y} has unit standard deviation and show that \mathbf{x}_j and \mathbf{y} have absolute correlation

$$\lambda(\alpha) = \frac{1 - \alpha}{\sqrt{(1 - \alpha)^2 + \frac{\alpha^2 - \alpha}{N} \cdot RSS}} \cdot \lambda.$$

Since $\tilde{\mathbf{y}} = 0$, the variance of \mathbf{y} is $\|\mathbf{y}\|^2/N$ and so reduces to the formula in the question if \mathbf{y} has standard deviation one.

Recall that the correlation between two variables is equal to their covariance over the product of their standard deviations. The standard deviation of each \mathbf{x}_j is 1 so it suffices to show that the sample standard deviation of $\mathbf{y} - \mathbf{u}(\alpha)$ is equal to the denominator of the expression. Since $\tilde{\mathbf{y}} = 0$ by assumption and $\mathbf{u}(\alpha) = \alpha \sum_{j=1}^p \lambda_j \mathbf{x}_j$, this is this is equivalent to showing

$$\frac{1}{N} \|\mathbf{y} - \alpha \tilde{\mathbf{y}}\|^2 = (1 - \alpha)^2 + \frac{\alpha^2(2 - \alpha)}{N} RSS.$$

First note that $\mathbf{y} - \tilde{\mathbf{y}}$ is orthogonal to the column space of \mathbf{X} so

$$\|\mathbf{y} - \tilde{\mathbf{y}}\|^2 = \langle \mathbf{y} - \tilde{\mathbf{y}}, \mathbf{y} \rangle - \langle \mathbf{y} - \tilde{\mathbf{y}}, \tilde{\mathbf{y}} \rangle = \langle \mathbf{y} - \tilde{\mathbf{y}}, \mathbf{y} \rangle.$$

Using this, we have

$$\begin{aligned} \|\mathbf{y} - \alpha \tilde{\mathbf{y}}\|^2 &= \|\alpha(\mathbf{y} - \tilde{\mathbf{y}}) + (1 - \alpha)\mathbf{y}\|^2 \\ &= \alpha^2 \|\mathbf{y} - \tilde{\mathbf{y}}\|^2 + 2\alpha(1 - \alpha) \langle \mathbf{y} - \tilde{\mathbf{y}}, \mathbf{y} \rangle + (1 - \alpha)^2 \|\mathbf{y}\|^2 \\ &= \alpha^2 \|\mathbf{y} - \tilde{\mathbf{y}}\|^2 + 2\alpha(1 - \alpha) \|\mathbf{y} - \tilde{\mathbf{y}}\|^2 + (1 - \alpha)^2 \|\mathbf{y}\|^2, \end{aligned}$$

which establishes the formula for $\lambda(\alpha)$.

Finally, one can show that

$$\lambda'(\alpha) = \frac{-RSS\lambda}{N \left(\frac{(1 - \alpha)^2}{N} \|\mathbf{y}\|^2 + \frac{\alpha^2 - \alpha}{N} \cdot RSS \right)^{3/2}} < 0$$

so that $\lambda(\alpha)$ decreases monotonically to zero.

(c) Start with standardised predictors with mean zero and unit variance.

Take $j \in \{1, \dots, p\}$ and assume that at the start of step j we have a residual vector \mathbf{r}_j which has mean zero and identical absolute correlation λ_j with $\mathbf{x}_1, \dots, \mathbf{x}_j$. Let $\mathbf{u}_j(\alpha)$ for $\alpha \in [0, 1]$ be the vector that moves a fraction α towards the least squares fit of $\mathbf{x}_1, \dots, \mathbf{x}_j$ to \mathbf{r}_j (for the first step we set $\mathbf{r}_1 = \mathbf{y} - \tilde{\mathbf{y}}$ and assume wlog that \mathbf{x}_1 has the greatest absolute correlation with \mathbf{r}_1).

Take $\alpha_j \in [0, 1]$ minimal such that a predictor \mathbf{x}_k with $k > j$ has correlation with $\mathbf{r}_j - \mathbf{u}_j(\alpha_j)$ equal to

$$\lambda_j(\alpha_j) = \frac{(1 - \alpha_j) \lambda_j}{\sqrt{\frac{1}{N} ((1 - \alpha_j)^2 \|\mathbf{r}_j\|^2 + \alpha_j(2 - \alpha_j) RSS)}}.$$

where RSS denotes the residual sum of squares between \mathbf{r}_j and its least squares fit (if there is no such k then we're done). Wlog assume $k = j + 1$ and set $\mathbf{r}_{j+1} = \mathbf{r}_j - \mathbf{u}_j(\alpha_j)$ and $\lambda_{j+1} = \lambda_j(\alpha_j)$. By (b) we are now in the right position to start step $j + 1$.

Assuming that \mathbf{X} has full column rank this process will terminate at the j th step where it reaches the least squares solution. This gives us a piecewise linear path $\mathbf{u}(\alpha)$ for $\alpha \in [0, \sum_{j=1}^p \alpha_j]$ with $\mathbf{u}(0) = \mathbf{y}$, $\mathbf{u}(\sum_{j=1}^p \alpha_j) = \hat{\beta}$ and

$$\mathbf{u} \left(\sum_{j=1}^k \alpha_j + t \right) = \mathbf{r}_k + \mathbf{u}_k(t)$$

for $t \in [0, \alpha_{j+1}]$. By part (a) the correlations stay tied and decreasing along this path

Exercise 3.24

LAR directions. Using the notation around equation (3.55) on page 74, show that the LAR direction makes an equal angle with each of the predictors in \mathcal{A}_k .

Solution

The angle between \mathbf{u}_k and \mathbf{x}_j is

$$\arccos \left(\frac{|\langle \mathbf{x}_j, \mathbf{u}_k \rangle|}{\|\mathbf{x}_j\| \|\mathbf{u}_k\|} \right)$$

so the smallest angle corresponds to the greatest value of $\frac{|\langle \mathbf{x}_k, \mathbf{u}_k \rangle|}{\|\mathbf{x}_k\| \|\mathbf{u}_k\|}$, which is the greatest absolute correlation. But by exercise 3.23 elements of \mathcal{A}_k have equal absolute correlation with \mathbf{u}_k , which is maximal among \mathbf{x}_j .

Exercise 3.25

LAR ahead (Efron et al., 2004, Sec. 2). Starting at the beginning of the k th step of the LAR algorithm, derive expressions to identify the next variable to enter the active set at step $k + 1$, and the value of α at which this occurs (using the notation around equation (3.55) on page 74).

Solution

Take $\mathbf{x}_k \in \mathcal{A}_k$ and $\mathbf{x}_j \notin \mathcal{A}_k$. We wish to identify the values of α such that the two vectors have the same absolute correlation with $\mathbf{r}_k - \alpha \mathbf{u}_k$. This is the case if and only if

$$\frac{|\langle \mathbf{x}_k, \mathbf{r}_k - \alpha \mathbf{u}_k \rangle|}{\|\mathbf{r}_k - \alpha \mathbf{u}_k\|} = \frac{|\langle \mathbf{x}_k, \mathbf{$$