

**Solution 1**

- (a) If  $B$  is invertible, then  $B^{-1/2}$  exists and the identity is immediate, so if  $C = B^{-1/2}AB^{-1/2}$  and  $(A + \alpha B)^{-1}Av = \eta v$ , where  $v \neq 0$ , then

$$B^{-1/2}(C + \alpha I)^{-1}B^{-1/2}Av = \eta v \implies (C + \alpha I)^{-1}CB^{1/2}v = \eta B^{1/2}v$$

so  $\eta$  is an eigenvalue of  $(C + \alpha I)^{-1}C$  with eigenvector  $w = B^{1/2}v$ . This implies that

$$Cw = \eta(C + \alpha I)w, \quad \text{so} \quad Cw = \frac{\alpha\eta}{1 - \eta}w = \eta'w,$$

say, where  $\eta = 1$  would lead to a contradiction. This yields

$$\frac{\alpha\eta}{1 - \eta} = \eta' \implies \eta = \frac{\eta'}{\alpha + \eta'}.$$

- (b) If  $A$  is invertible,

$$\begin{aligned} A^{1/2}(A + \alpha B)^{-1}A &= A^{1/2}(A^{1/2}A^{1/2} + \alpha B)^{-1}A \\ &= (I + \alpha A^{-1/2}BA^{-1/2})^{-1}A^{-1/2}A \\ &= (I + \alpha A^{-1/2}BA^{-1/2})^{-1}A^{1/2}, \end{aligned}$$

and then an argument similar to that above gives

$$\eta = \frac{1}{1 + \alpha\eta''},$$

where  $\eta''$  is an eigenvalue of  $A^{-1/2}BA^{-1/2}$ .

**Solution 2**

- (a) We have  $y \sim (X\beta, \sigma^2 I_n) \sim (UD\gamma, \sigma^2 I_n)$ , where  $\gamma = V^T\beta$ , and

$$\hat{\beta} = (X^T X)^{-1}X^T y = (VD^T U^T U D V^T)^{-1}VD^T U^T y = V(D^T D)^{-1}D^T U^T y,$$

with a similar calculation giving  $\hat{\gamma} = (D^T D)^{-1}D^T U^T y$ , so  $\hat{\beta} = V\hat{\gamma}$  (surprise!)

- (b) As  $\gamma = V\beta$  and  $V$  is orthogonal,

$$Q = (\hat{\beta} - \beta)^T(\hat{\beta} - \beta) = (V\hat{\gamma} - V\gamma)^T(V\hat{\gamma} - V\gamma) = (\hat{\gamma} - \gamma)^T V^T V (\hat{\gamma} - \gamma) = (\hat{\gamma} - \gamma)^T(\hat{\gamma} - \gamma),$$

as required. Having  $y \sim (UD\gamma, \sigma^2 I_n)$  implies that

$$\hat{\gamma} = \text{diag}(d_1^{-1}, \dots, d_p^{-1}, 0, \dots, 0)U^T y,$$

so

$$\text{var}(\hat{\gamma}) = \sigma^2 \text{diag}(d_1^{-2}, \dots, d_p^{-2}).$$

This will be large if at least one of the  $d_r$  is small, and then there is at least one direction in which  $\gamma$ , i.e.,  $v^T\beta$  for some  $v_{p \times 1}$ , is extremely poorly determined.

(c) Under the normal model,  $\hat{\gamma}_1, \dots, \hat{\gamma}_p$  are independent  $\mathcal{N}(\gamma_r, \sigma^2/d_r^2)$  variables, so  $\hat{\gamma}_r - \gamma_r \stackrel{D}{=} \sigma Z_r/d_r$ , giving

$$Q = (\hat{\gamma} - \gamma)^T(\hat{\gamma} - \gamma) \stackrel{D}{=} \sum_{r=1}^p \sigma^2 Z_r^2/d_r^2,$$

and as  $E(Z_r^2) = 1$  and  $\text{var}(Z_r^2) = 2$  we get  $E(Q) = \sigma^2 \sum_{r=1}^p 1/d_r^2$  and  $\text{var}(Q) = 2\sigma^4 \sum_{r=1}^p 1/d_r^4$ .

### Solution 3

(a) We have

$$(y - X\beta)^T(y - X\beta) + \lambda\beta^T\beta = y^T y - 2y^T X\beta + \beta^T(X^T X + \lambda I_p)\beta,$$

and differentiation with respect to  $\beta$  gives first and second derivatives

$$-2yX^T + 2(X^T X + \lambda I_p)\beta, \quad 2(X^T X + \lambda I_p).$$

The second derivative matrix is positive definite for any  $\lambda > 0$ , and setting the first derivative to zero gives

$$\hat{\beta}_\lambda = (X^T X + \lambda I_p)^{-1} X^T y.$$

(b) Setting  $X = UDV^T$  gives  $X^T y = VD^T U^T y = \sum_{d_j > 0} v_j d_j u_j^T y$  and

$$(X^T X + \lambda I_p)^{-1} = (VD^T DV^T + \lambda I_p)^{-1} = \{V(D^T D + \lambda I_p)V^T\}^{-1} = VS_\lambda V^T,$$

where  $S_\lambda = \text{diag}(d_1^2 + \lambda, \dots, d_r^2 + \lambda)^{-1}$  exists because all its elements are positive. Hence

$$\hat{\beta}_\lambda = (X^T X + \lambda I_p)^{-1} X^T y = VS_\lambda V^T (VD^T U^T) y = \sum_{d_j > 0} \frac{d_j}{d_j^2 + \lambda} u_j^T y \times v_j.$$

Likewise

$$\hat{y}_\lambda = X\hat{\beta}_\lambda = H_\lambda y = UDS_\lambda D^T U^T y = \sum_{d_j > 0} u_j \times \frac{d_j^2}{d_j^2 + \lambda} u_j^T y.$$

Both  $\hat{\beta}_\lambda$  and  $\hat{y}_\lambda$  shrink towards zero as  $\lambda$  increases, with the strongest shrinkage for those vectors  $v_j$  and  $u_j$  for which  $d_j$  is smallest.

(c) We have

$$\text{edf}_\lambda = \text{tr}(H_\lambda) = \text{tr}\{(UDV^T)VS_\lambda V^T(UDV^T)^T\} = \text{tr}\{U^T UDS_\lambda D^T\} = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda},$$

and

$$E(\hat{\beta}_\lambda) = VS_\lambda V^T VD^T U^T UDV^T \beta = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda} v_j^T \beta \times v_j,$$

$$\text{var}(\hat{\beta}_\lambda) = VS_\lambda D^T U^T \text{cov}(y) \{VS_\lambda D^T U^T\}^T = \sigma^2 V \text{diag} \left\{ \frac{d_1^2}{(d_1^2 + \lambda)^2}, \dots, \frac{d_p^2}{(d_p^2 + \lambda)^2} \right\} V^T,$$

so the bias is

$$E(\hat{\beta}_\lambda) - \beta = \sum_{r=1}^p \frac{d_r^2}{d_r^2 + \lambda} v_r^T \beta \times v_r - \sum_{r=1}^p v_r v_r^T \beta = - \sum_{r=1}^p \frac{\lambda}{d_r^2 + \lambda} v_r^T \beta \times v_r$$

### Solution 4

- (a) Let  $H = 1_n(1_n^T 1_n)^{-1} 1_n^T$  correspond to regression on a column of ones, and note that  $(I_n - H)1_n = 0$ ,  $Hy = 1_n \bar{y}$  and  $HX = 1_n \bar{x}^T$ , where  $\bar{x}$  contains the means of the columns of  $X$ . Then we can set  $y_* = (I_n - H)y$  and  $X_* = (I_n - H)X$ , so

$$y - \beta_0 1_n - X\beta = (I_n - H)y + Hy - \beta_0 1_n - (I_n - H)X\beta + HX\beta = y_* - (\gamma - \bar{y})1_n - X_*\beta,$$

where  $\gamma = \beta_0 - \bar{x}^T \beta$ . The interpretation of  $\beta$  remains the same; only the intercept  $\gamma$  has changed.

- (b) The equality implies that we can write

$$\|y - \beta_0 1_n - X\beta\|_2^2 = \|y_* - (\gamma - \bar{y})1_n - X_*\beta\|_2^2 = \|y_* - X_*\beta\|_2^2 + \|(\gamma - \bar{y})1_n\|_2^2,$$

because

$$(y_* - X_*\beta)^T (\gamma - \bar{y})1_n = (\gamma - \bar{y})(y - X\beta)^T (I - H)^T 1_n = (\gamma - \bar{y})(y - X\beta)^T (I - H)1_n = 0.$$

Hence

$$\min_{\beta_0, \beta} \|y - \beta_0 1_n - X\beta\|_2^2 + \lambda p(\beta) = \min_{\gamma, \beta} \|y_* - X_*\beta\|_2^2 + \|(\gamma - \bar{y})1_n\|_2^2 + \lambda p(\beta),$$

which gives  $\hat{\gamma} = \bar{y}$  and  $\hat{\beta}_\lambda$  as the solution to the second minimisation problem, as required. Hence provided  $y$  and the columns of  $X$  are centered, the intercept need not be included.

- (c) Including  $\beta_0$  in  $\beta$  would mean that as  $\lambda$  increases,  $\hat{\beta}_\lambda \rightarrow 0$ , i.e., shrinkage would apply also to the intercept, which depends on the units used for measuring  $y$ . Hence a change from measuring temperature in  $^\circ C$  to  $^\circ F$  would lead to different conclusions about the effects of the covariates, which is clearly undesirable.

Expressed in algebra, we would have a column  $1_n$  in  $X$  if  $\beta$  contains the intercept, and then if the intercept is the first column of  $X$ , we have

$$y - X\beta \mapsto ay + b1_n - X\beta = a(y - X\beta_*), \quad \beta \mapsto \beta_* = \{\beta - (b, 0, \dots, 0)^T\}/a.$$

In this new parametrisation we have  $\hat{\beta}_{*,\lambda} \rightarrow 0$  as  $\lambda \rightarrow \infty$ , corresponding to the estimate of  $\beta_0$  tending to  $b$  rather than to 0, and this would affect all the other parameter estimates.