

In all the questions that follow,  $X$  is an  $n$  by  $p$  design matrix with full column rank and  $P$  is the orthogonal projection on to the column space of  $X$ . We will assume that  $n - p \geq 2$ . The vector  $Y \in \mathbb{R}^n$  will be a vector of responses and we will define  $\hat{\beta} := (X^T X)^{-1} X^T Y$ ,  $\tilde{\sigma}^2 := \|(I - P)Y\|^2 / (n - p)$  and  $\hat{\varepsilon} := Y - X\hat{\beta}$ .

1. Show that writing  $P_0$  for the orthogonal projection on to  $X_0$ , a matrix composed of a (proper and non-empty) subset of the columns of  $X$ , we have

$$\|(P - P_0)Y\|^2 = \|(I - P_0)Y\|^2 - \|(I - P)Y\|^2.$$

2. Data are available on weights of two groups of three rats at the beginning of a fortnight and at its end. During the fortnight, one group was fed normally, and the other was given a growth inhibitor. The weights of the  $k^{\text{th}}$  rat in the  $j^{\text{th}}$  group before and after the fortnight are  $X_{jk}$  and  $y_{jk}$  respectively. The  $y_{jk}$  are taken as realisations of random variables  $Y_{jk}$  that follow the model  $Y_{jk} = \alpha_j + \beta_j X_{jk} + \varepsilon_{jk}$ .

(a) Let  $Y$  be the vector of responses, so  $Y = (Y_{11}, Y_{12}, Y_{13}, Y_{21}, Y_{22}, Y_{23})^T$ , and similarly let  $\varepsilon$  be the vector of random errors. Write down the model above in the form  $Y = X\theta + \varepsilon$ , giving the design matrix  $X$  explicitly.

(b) The model is to be reparametrised in such a way that it can be specialised to (i) two parallel lines for the two groups, (ii) two lines with the same intercept, (iii) one common line for both groups, just by setting parameters to zero. Give one design matrix that can be made to correspond to (i), (ii) and (iii), just by dropping columns, specifying which columns are to be dropped for which cases.

3. Suppose the design matrix  $X$  consists of just a single variable and a column of 1's representing an intercept term (as the first column). Show that the leverage,  $p_i$ , of the  $i^{\text{th}}$  observation satisfies

$$p_i = \frac{1}{n} + \frac{(X_{i2} - \bar{X}_2)^2}{\sum_{k=1}^n (X_{k2} - \bar{X}_2)^2},$$

where  $\bar{X}_2 := \frac{1}{n} \sum_{k=1}^n X_{k2}$ . *Hint: Why can we assume that the  $i^{\text{th}}$  component of the second column is  $X_{i2} - \bar{X}_2$  rather than  $X_{i2}$ ?*

4. Return to the brain sizes data studied in practical 3.

```
> file_path <-
+ "http://www.statslab.cam.ac.uk/~rds37/teaching/statistical_modelling/"
> BrainSize <- read.csv(paste(file_path, "BrainSize.csv", sep = ""))
> attach(BrainSize)
> BrainSizeLM2 <- lm(PIQ ~ MRI_Count + Height)
```

In this question we will plot a confidence ellipse for the coefficients for brain size and height. To do this, first install the `ellipse` package using

```
> install.packages("ellipse")
```

and selecting a mirror of your choice. Next load the package with `library(ellipse)`. Look at `?ellipse.lm` and plot a 95% confidence ellipse for the coefficients with

```
> plot(ellipse(BrainSizeLM2, c(2, 3)), type = "l")
```

Using `abline` add to the plot the end points of 95% confidence intervals for each of the coefficients in red, and also add in blue the sides of the confidence rectangle in question 7 of Example sheet 1. If you are using `Rstudio`, you can output a pdf of your plot by clicking on “Export” above the plot window. Now look at the correlation between the estimates of these coefficients using

```
> summary(BrainSizeLM2, correlation = TRUE)$correlation
```

and compare this to the correlation between the corresponding variables

```
> cor(Height, MRI_Count)
```

What do you notice? Explain.

5. Let  $f$  and  $g$  be two densities on  $\mathbb{R}$  with  $S_g := \{x : g(x) > 0\} \subseteq \{x : f(x) > 0\}$ . Show that the Kullback–Liebler divergence,

$$K(g, f) := \int_{S_g} [\log\{g(x)\} - \log\{f(x)\}]g(x)dx,$$

is non-negative.

6. Consider forward selection in the linear model  $Y = \beta_0 \mathbf{1}_n + X\beta + \varepsilon$ , where  $\mathbf{1}_n$  is an  $n$ -vector of 1's. At the 0<sup>th</sup> stage, only the intercept term is in the model. Now suppose that the design matrix for the model fitted in the  $k$ <sup>th</sup> stage for  $k < p$  is  $X^{(k)} := (\mathbf{1}_n \ X_{j_1} \ \cdots \ X_{j_k})$ , where  $X_j$  denotes the  $j$ <sup>th</sup> column of  $X$ . Show that the next variable to enter the model is  $X_{j^*}$  where

$$j^* = \operatorname{argmax}_{j \neq j_1, \dots, j_k} \frac{|(X_j^\perp)^T Y|}{\|X_j^\perp\|}.$$

Here  $X_j^\perp$  denotes the orthogonal projection of  $X_j$  onto the orthogonal complement of the column space of  $X^{(k)}$ .

7. One of the data sets in the *Modern Applied Statistics in S-Plus* (MASS) library is `hills`. You can find out about the data with

```
> library(MASS)
> ?hills
> pairs(hills)
```

The data contain one known error in the winning time. Identify this error (think carefully!) and subtract an hour from the winning time. *Hint: You can examine the plots and identify observations for which the response and covariates satisfy certain inequalities e.g.*

```
> hills[(hills$time > 50) & (hills$dist < 10), ]
```

Can you see any reason why we might want to consider taking logarithms of the variables? Explain why we should include an intercept term if we do choose to take logarithms.

Explore at least two linear models for this data, and give estimates with standard errors for your preferred model. Predict the record time for a hypothetical 5.3 mile race with a 1100ft climb, giving a 95% prediction interval.

8. (a) Let  $A$  be a  $p \times p$  non-singular matrix and let  $b \in \mathbb{R}^p$ . Prove that if  $b^T A^{-1} b \neq 1$ , then  $A - bb^T$  is invertible with inverse given by

$$(A - bb^T)^{-1} = A^{-1} + \frac{A^{-1}bb^T A^{-1}}{1 - b^T A^{-1}b}.$$

- (b) Consider a linear model  $Y = X\beta + \varepsilon$  with  $\text{Var}(\varepsilon) = \sigma^2 I$ , and let  $x_i^T$  denote the  $i^{\text{th}}$  row of  $X$ . Further, let  $X_{(-i)}$  denote the  $(n-1) \times p$  matrix obtained by deleting the  $i^{\text{th}}$  row of  $X$ , and suppose that this matrix has full column rank and that the leverage score of the  $i^{\text{th}}$  observation,  $p_i$ , is less than 1. By noting that

$$X^T X = \sum_{i=1}^n x_i x_i^T,$$

prove that writing  $\hat{\beta}_{(-i)}$  for the OLS estimate of  $\beta$  when the  $i^{\text{th}}$  observation has been removed, the difference

$$\text{Var}(\hat{\beta}_{(-i)}) - \text{Var}(\hat{\beta})$$

is positive semi-definite. Here  $\hat{\beta}$  is the usual OLS estimate of  $\beta$  based on all  $n$  observations.

- (c) Show that

$$\hat{\beta} - \hat{\beta}_{(-i)} = \frac{1}{1 - p_i} (X^T X)^{-1} x_i (Y_i - x_i^T \hat{\beta}),$$

and hence deduce that the Cook's distance  $D_i$  of the observation  $(Y_i, x_i)$  satisfies

$$D_i = \frac{1}{p} \left( \frac{p_i}{1 - p_i} \right) \hat{\eta}_i^2,$$

where  $\hat{\eta}_i = (Y_i - x_i^T \hat{\beta}) / (\tilde{\sigma} \sqrt{1 - p_i})$  is the  $i^{\text{th}}$  studentised fitted residual.

9. (a) (Continuation) The *externally studentised residual* of the  $i^{\text{th}}$  observation may be defined as

$$\tilde{\eta}_i := \frac{\hat{\varepsilon}_i}{\tilde{\sigma}_{(-i)} \sqrt{1 - p_i}},$$

where  $\tilde{\sigma}_{(-i)}$  is the equivalent of  $\tilde{\sigma}$  but calculated omitting the  $i^{\text{th}}$  observation, so

$$\tilde{\sigma}_{(-i)}^2 = \frac{1}{n - p - 1} \|Y_{(-i)} - X_{(-i)} \hat{\beta}_{(-i)}\|^2,$$

where  $Y_{(-i)}$  is the response  $Y$  without the  $i^{\text{th}}$  component. Show that  $\tilde{\eta}_i \sim t_{n-p-1}$ . *Hint: It may help to first show that*

$$\hat{\varepsilon}_i = (1 - p_i)(Y_i - x_i^T \hat{\beta}_{(-i)}).$$

How can we construct a hypothesis test based on  $\tilde{\eta}_i$  to test whether the  $i^{\text{th}}$  observation is an outlier?

- (b) Another dataset in the MASS package is `mammals` which gives the body and brain masses of 68 mammals. Fit a linear model to the log-transformed data and then apply your hypothesis test to check whether the observation corresponding to humans is an outlier. The function `rstudent` that calculates externally studentised residuals may be of help. What is the  $p$ -value you obtain?

10. Show that

$$p_i + \frac{\hat{\varepsilon}_i^2}{\|(I - P)Y\|^2} \leq 1,$$

so if  $p_i$  is close to 1, the  $i^{\text{th}}$  residual is forced to be close to 0. *Hint: Use the Cauchy-Schwarz inequality and the fact that  $I - P$  is an orthogonal projection.*