

# Applied Data Analysis

## Exercise Sheet 4

### Exercise 14

Consider a normal linear model

$$\mathbf{Y} = B\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

as given in Definition I.4.3 with  $n \geq \text{rank}(B) = p \geq 2$ , (unknown) parameter vector  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)' \in \mathbb{R}^p$  and error term  $\boldsymbol{\varepsilon} \sim N_n(0, \sigma^2 I_n)$ , where  $\sigma > 0$  is also unknown.

(a) Explain, how to test the following hypotheses

$$(i) \quad H_0: \beta_1 = \beta_2 \quad \longleftrightarrow \quad H_1: \beta_1 \neq \beta_2,$$

$$(ii) \quad H_0: \beta_1 = \dots = \beta_p \quad \longleftrightarrow \quad H_1: \beta_j \neq \beta_k \text{ for at least one pair } (j, k) \text{ with } j \neq k,$$

by using general linear hypothesis testing and nested model comparison.

(b) For  $\alpha \in (0, 1)$ , construct two confidence regions for  $\beta_1 - \beta_2$  of confidence level  $1 - \alpha$ , one based on the F-distribution and one based on the t-distribution.

**Hint:** Confidence regions can be constructed by inverting a statistical hypothesis test procedure in the sense of constructing the non-rejecting region corresponding to that statistical test.

(c) If for the testing problem (a),(ii)  $H_0$  is rejected, one might ask, which pairs in

$$M := \{(j, k) \in \{1, \dots, p\} \mid j \neq k\}$$

have significantly different effects.

To answer this question, the multiple pairwise comparison problem

$$H_0^{(j,k)} : \beta_j = \beta_k \quad \longleftrightarrow \quad H_1^{(j,k)} : \beta_j \neq \beta_k$$

for all  $m := \binom{p}{2}$  pairs  $(j, k) \in M$  could be considered.

Show for  $\alpha \in (0, 1)$  that if each single hypothesis  $H_0^{(j,k)}$ ,  $(j, k) \in M$ , is tested on significance level  $\frac{\alpha}{m}$ , the multiple pairwise comparison tests all together hold a family-wise error rate of  $\alpha$  (so-called *Bonferroni correction*).

To formalize this task, for  $\boldsymbol{\beta} \in \mathbb{R}^p$ , let  $I_0(\boldsymbol{\beta}) \subseteq M$  denote the subset of all pairs  $(j, k)$  for which  $H_0^{(j,k)}$  is true (i.e.  $\beta_j = \beta_k$ ) and for  $(j, k) \in M$ , let  $A^{(j,k)}$  denote the event that the corresponding testing procedure decides in favour for  $H_1^{(j,k)}$ .

Then, you have to show for  $\boldsymbol{\beta} \in \mathbb{R}^p$ :

$$P_{\boldsymbol{\beta}} \left( \bigcup_{(j,k) \in I_0(\boldsymbol{\beta})} A^{(j,k)} \right) \leq \alpha$$

(with  $P_{\boldsymbol{\beta}}$  denoting the underlying probability distribution corresponding to  $\boldsymbol{\beta}$ ).

## Exercise 15

Consider the model of simple linear regression as given in Example I.4.6 with  $n \geq 2 = \text{rank}(B)$  and corresponding sums of squares as given in Theorem I.5.11:

$$\text{SST} = \sum_{i=1}^n (Y_i - \bar{Y})^2, \quad \text{SSR} = \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_i - \bar{Y})^2, \quad \text{SSE} = \sum_{i=1}^n (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2.$$

According to Definition I.5.43, the corresponding coefficient of determination is defined by

$$R^2 := \frac{\text{SSR}}{\text{SST}} \stackrel{\text{Th. I.5.11,(1)}}{=} \frac{\text{SST} - \text{SSE}}{\text{SST}}.$$

Show:

$$(a) \quad (\hat{\beta}_0, \hat{\beta}_1)' = \left( \bar{Y} - \frac{s_{\mathbf{x}\mathbf{Y}}}{s_{\mathbf{x}\mathbf{x}}} \bar{x}, \frac{s_{\mathbf{x}\mathbf{Y}}}{s_{\mathbf{x}\mathbf{x}}} \right)',$$

$$(b) \quad \text{SSE} = (n-1) s_{\mathbf{Y}\mathbf{Y}} (1 - r_{\mathbf{x}\mathbf{Y}}^2),$$

$$(c) \quad R^2 = r_{\mathbf{x}\mathbf{Y}}^2,$$

using the usual notations of the corresponding sample variances, sample covariance and Pearson correlation coefficient

$$s_{\mathbf{x}\mathbf{x}} := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \quad s_{\mathbf{Y}\mathbf{Y}} := \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2,$$

$$s_{\mathbf{x}\mathbf{Y}} := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y}), \quad r_{\mathbf{x}\mathbf{Y}} := \frac{s_{\mathbf{x}\mathbf{Y}}}{\sqrt{s_{\mathbf{x}\mathbf{x}} s_{\mathbf{Y}\mathbf{Y}}}}.$$

## Exercise 16

Consider the one factorial analysis of variance model in effect representation as given in Definition I.6.15 with the side condition

$$\sum_{i=1}^p n_i \alpha_i = 0.$$

For the corresponding sum of squares  $\text{SSR} = \sum_{i=1}^p n_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2$  (defined on slide 165), show:

$$\text{E}(\text{SSR}) = (p-1) \sigma^2 + \sum_{i=1}^p n_i \alpha_i^2.$$

## Exercise 17

Consider the following two classes of univariate probability distributions:

(i) The class of all univariate (regular) normal distributions  $\{\mathcal{N}(\mu, \sigma^2) \mid \mu \in \mathbb{R}, \sigma > 0\}$ ,

(ii) the class of all Poisson distributions  $\{\mathcal{P}(\mu) \mid \mu \in (0, \infty)\}$ .

(a) Show that each of these two classes is a member of the (univariate) Exponential Dispersion Family (EDF) by finding suitable representations of the corresponding density functions or probability mass functions, respectively, as in Definition II.2.3.

In particular, for each class, confirm the representations for the natural parameter  $\theta$  and the functions  $b$ ,  $a$  and  $c$  given in Example II.2.5.

(b) Using the results of (a), confirm the expressions for the mean  $\text{E}(Y)$  and for the variance  $\text{Var}(Y)$  given in Example II.2.5 for each of the two classes.