

**KTH Matematik** 

# SF2930 Regression Analysis Exam Generator Version for the re-exam 2023-06-07

Timo Koski

# Innehåll

1	Introduction         1.1       Instructions         1.2       Linear Algebra for Warm-Up	<b>3</b> 3 3			
<b>2</b>	Matrix Tricks and Linear Regression				
3	Covariance Matrices, Random Vectors, Minimum Mean Square Estimation				
4	Linear Regression4.1Simple Linear Regression4.2Multiple Linear Regression4.3Estimability and the Gauss-Markov Theorem4.4Normal Linear Regression	9 9 15 20 23			
<b>5</b>	Properties of the Hat Matrix and Diagnostics of Regression Models 2				
6	Choice of Regression Models 3				
7	Generalized Linear Regression				
8	The Woodbury Matrix Identity & Ridge Regression				

39
39
39
40
40
40
41
41
42
42
42
43
43
44
44
45
45
45
46

## 1 Introduction

## 1.1 Instructions

This is the final version. Errors and typos and bad formulations will sought for and corrected, whenever detected.

In the Exam the students will be provided with two documents:

- 1. First: an A4 sheet pointing out by problem number (e.g., 4.11) those problems in this generator to be solved in the Exam hall.
- 2. Second: The complete Exam generator (incl. Appendices, that is section section 9) WILL BE AVAI-LABLE IN THE EXAM. When solving a problem in the exam, one may refer to formulas and statement in problems in the generator NOT included in the exam without deriving or proving them.

The formulas in section 9 can be used in the solutions of any exam problems without proof or derivation.

## 1.2 Linear Algebra for Warm-Up

The problems in section are not very likely to appear as exam questions, but can be useful in more serious candidates for exam assignments.

## Problem 1.1.

The ON-basis of  $\mathbb{R}^n$  is

$$\mathfrak{E}_{1} = \begin{pmatrix} 1 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{pmatrix} \qquad \qquad \mathfrak{E}_{j} = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \qquad \qquad \mathfrak{E}_{n} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \vdots \\ 1 \end{pmatrix}. \tag{1}$$

Let A be any  $n \times n$  matrix. Then

$$\mathfrak{E}_i^T A \tag{2}$$

is the ith row of A.

• Express in words the meaning of the following expressions:

 $A\mathfrak{E}_i, \quad \mathfrak{E}_i^T A\mathfrak{E}_i \quad \mathfrak{E}_i^T A\mathfrak{E}_j \quad \mathfrak{E}_j^T A\mathfrak{E}_i. \tag{3}$ 

#### Problem 1.2.

Let  $A^T = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n)$  where  $\mathbf{a}_i$  is the *i*th column vector of the  $k \times n$  matrix  $A^T$ . Thus  $\mathbf{a}_i^T$  is the *i*th row vector of A (an  $n \times k$  matrix), i.e.,

$$A = \begin{pmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_n^T \end{pmatrix}.$$

Let B be a  $k \times k$  matrix. Check that

$$ABA^{T} = \begin{pmatrix} \mathbf{a}_{1}^{\mathsf{T}}B\mathbf{a}_{1} & \mathbf{a}_{1}^{\mathsf{T}}B\mathbf{a}_{2} & \cdots & \mathbf{a}_{1}^{\mathsf{T}}B\mathbf{a}_{n} \\ \mathbf{a}_{2}^{\mathsf{T}}B\mathbf{a}_{1} & \mathbf{a}_{2}^{\mathsf{T}}B\mathbf{a}_{2} & \cdots & \mathbf{a}_{2}^{\mathsf{T}}B\mathbf{a}_{n} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{a}_{n}^{\mathsf{T}}B\mathbf{a}_{1} & \mathbf{a}_{n}^{\mathsf{T}}B\mathbf{a}_{2} & \cdots & \mathbf{a}_{n}^{\mathsf{T}}B\mathbf{a}_{n} \end{pmatrix}.$$
(4)

## 2 Matrix Tricks and Linear Regression

### Problem 2.1.

Consider the  $k \times k$  matrix

$$S := \frac{1}{n-1} X_c^T X_c.$$

$$\tag{5}$$

where  $X_c$  is given in (64).

- a) Is S positive definite? Justify your answer. Aid: Slide 22/81 in Lecture 3.
- b) Check that

$$s_{jk} = \mathfrak{E}_{j}^{T} S \mathfrak{E}_{k} = \frac{1}{n-1} \sum_{i=1}^{n} \left( x_{ij} - \bar{x}_{j} \right) \left( x_{ik} - \bar{x}_{k} \right), \tag{6}$$

(c.f. (3)), where  $\mathfrak{E}_j$  and  $\mathfrak{E}_k$  are now column vectors in the ON-basis of  $\mathbb{R}^k$ , see (1).

 $s_{jk}$  in (6) is an estimate of the covariance between the regressor j and regressor number k in the population underlying the rows in the matrix  $X_R = (\mathbf{x}_1^T \dots \mathbf{x}_n^T)^T$ , see also (63). Hence the matrix S is called the **sample covariance matrix**.

## Problem 2.2.

The centering matrix is

$$C_{ce} = \mathbb{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T.$$

See section 9.1 and (97) for the detailed definition.

a) Check that  $C_{ce}$  is a projection matrix, c.f. section 9.3.

b) Show that

$$\mathbf{y}^T C_{ce} C_{ce} \mathbf{y} = \sum_{i=1}^n \left( y_i - \bar{y} \right)^2 \tag{7}$$

c) Is the centering matrix invertible? Your answer should contain an explanation, not merely one of Yes or No.

### Problem 2.3.

A training set of observed responses  $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$  and the corresponding  $n \times (k+1)$  data matrix X are available. We fit a multiple least squares model to the training set and obtain the LSE  $\hat{\boldsymbol{\beta}}$ . Then

$$\hat{\mathbf{y}} = X\hat{\boldsymbol{\beta}} = H\mathbf{y}$$

is the predictor vector, where H is the hat matrix and  $\hat{\mathbf{y}} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n)^T$ .

a)  $\hat{e}_i = y_i - \hat{y}_i, i = 1, ..., n$  are the LS residuals. Show that

$$\sum_{i=1}^{n} \hat{e}_i = 0$$

by studying

$$\hat{\mathbf{e}}^T \mathbf{1}_n,$$

where  $\hat{\mathbf{e}} = (\hat{e}_1, \hat{e}_2 \dots \hat{e}_n)^T$  is an  $n \times 1$  vector.

b)  $\hat{y}_i$  are the LS-predictors of  $y_i$ , respectively, for  $i = 1, \ldots, n$ . Show that

$$\frac{1}{n}\sum_{1=}^{n}\hat{y}_{i}=\bar{y}$$

by studying

$$\hat{\mathbf{y}}^T \mathbf{1}_n,$$

where  $\hat{\mathbf{y}} = (\hat{y}_1, \hat{y}_2 \dots \hat{y}_n)^T$  is an  $n \times 1$  vector.

c) Show that

$$\sum_{i=1}^{n} (y_i - \bar{y})^2 = \min_{\beta} \|\mathbf{y} - \mathbf{1}_n\beta\|^2$$

where  $\mathbf{y} = (y_1, y_2 \dots y_n)^T$  is an  $n \times 1$  vector and  $\beta$  is a scalar.

## Problem 2.4.

 $\mathbf{x}$  and  $\mathbf{y}$  are two column vectors in  $\mathbb{R}^n$ . Then we define the (sample) correlation coefficient cor  $(\mathbf{x}, \mathbf{y})$  by

$$\operatorname{cor}\left(\mathbf{x},\mathbf{y}\right) := \frac{\mathbf{x}^{T} C_{ce} \mathbf{y}}{\sqrt{\mathbf{x}^{T} C_{ce} \mathbf{x}} \sqrt{\mathbf{y}^{T} C_{ce} \mathbf{y}}}$$
(8)

A training set of observed responses  $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$  and the corresponding  $n \times (k+1)$  data matrix X are available. We fit a multiple least squares model to the training set and obtain the LSE  $\hat{\boldsymbol{\beta}}$ . Then

$$\hat{\mathbf{y}} = X\hat{\boldsymbol{\beta}} = H\mathbf{y}$$

is the predictor vector, where H is the hat matrix. Show by means of (8) that

$$\operatorname{cor}\left(\mathbf{y}, \hat{\mathbf{y}}\right) = \sqrt{\frac{\mathrm{SS}_{R}}{\mathrm{SS}_{T}}}$$

Aid: Use the hat matrix formula for  $\hat{\mathbf{y}}$ .

#### Problem 2.5.

a) Set, see (96) in section 9.1 for details,

$$P_{\mathbf{1}_n} = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T.$$

Check that  $P_{\mathbf{1}_n}$  is a projection matrix, see section 9.3.

b) Show that the range space, see (98), of  $P_{\mathbf{1}_n}$  is

$$\mathcal{R}(P_{\mathbf{1}_n}) = \left\{ \mathbf{v} \in \mathbb{R}^n | \mathbf{v} = \hat{v} \mathbf{1}_n \right\},\$$

where

$$\hat{v} = \frac{\mathbf{v}^T \mathbf{1}_n}{\parallel \mathbf{1}_n \parallel^2}$$

In other words,  $\mathcal{R}(P_{\mathbf{1}_n})$  is the subspace of vectors proportional to the vector  $\mathbf{1}_n$ .

c) We have that the centering matrix, see (97) in section 9.1 is also written as

$$C_{ce} = \mathbb{I}_n - P_{\mathbf{1}_n}.$$

We know that  $C_{ce}$  is a projection matrix. What is the range space of  $C_{ce}$ ? c) How are the range spaces of  $P_{\mathbf{1}_n}$  and  $C_{ce}$  related to each other?

# 3 Covariance Matrices, Random Vectors, Minimum Mean Square Estimation

## Problem 3.1.

An  $n \times 1$  random vector **X** has expectation  $\boldsymbol{\mu}_{\mathbf{X}} = E[\mathbf{X}]$ , and An  $n \times 1$  random vector **Y** has expectation  $\boldsymbol{\mu}_{\mathbf{Y}} = E[\mathbf{Y}]$ . Define their cross-covariance matrix as

$$\operatorname{Cov}(\mathbf{X}, \mathbf{Y}) := E\left[ (\mathbf{X} - \boldsymbol{\mu}_{\mathbf{X}}) (\mathbf{Y} - \boldsymbol{\mu}_{\mathbf{Y}})^T \right].$$

It is easily seen that (you need not prove this)

$$\operatorname{Cov}(\mathbf{Y}, \mathbf{X}) = \operatorname{Cov}(\mathbf{X}, \mathbf{Y})^{T}.$$
(9)

Clearly

$$C_{\mathbf{X}} = \operatorname{Cov}(\mathbf{X}, \mathbf{X}). \tag{10}$$

a) Show that

$$\operatorname{Cov}(\mathbf{X}, \mathbf{Y}) = E\left[\mathbf{X}\mathbf{Y}^{T}\right] - \boldsymbol{\mu}_{\mathbf{X}}\boldsymbol{\mu}_{\mathbf{Y}}^{T}.$$

b) A and B are matrices of suitable dimensions. Show that

$$Cov(A\mathbf{X}, B\mathbf{Y}) = ACov(\mathbf{X}, \mathbf{Y})B^{T}.$$
(11)

c) A, B and C are matrices of suitable dimensions. **Z** is an  $n \times 1$  random vector. Show that

$$\operatorname{Cov}(A\mathbf{X}, B\mathbf{Y} + C\mathbf{Z}) = \operatorname{Cov}(A\mathbf{X}, B\mathbf{Y}) + \operatorname{Cov}(A\mathbf{X}, C\mathbf{Z}).$$
(12)

## Problem 3.2.

**X** is an  $n \times 1$  random vector with  $\boldsymbol{\mu}$  as mean vector and  $\Sigma$  as a positive definite covariance matrix. The **Mahalanobis distance** between a vector  $\mathbf{x} \in \mathbb{R}^n$  and the mean  $\boldsymbol{\mu}$  is defined as

$$d_M(\mathbf{x}, \boldsymbol{\mu}; \boldsymbol{\Sigma}) := \sqrt{\left(\mathbf{x} - \boldsymbol{\mu}\right)^T \boldsymbol{\Sigma}^{-1} \left(\mathbf{x} - \boldsymbol{\mu}\right)}.$$
(13)

a) Find

$$E\left[d_{M}^{2}\left(\mathbf{X},\boldsymbol{\mu};\Sigma\right)
ight].$$

Aid: Use  $E\left[\left(\mathbf{X}-\boldsymbol{\mu}\right)^T \Sigma^{-1} \left(\mathbf{X}-\boldsymbol{\mu}\right)\right] = \operatorname{Tr} E\left[\left(\mathbf{X}-\boldsymbol{\mu}\right) \Sigma^{-1} \left(\mathbf{X}-\boldsymbol{\mu}\right)^T\right]$  and and a linear change of variable a factorization  $\Sigma = AA^T$ , where A is invertible. See the statements in section 9.12.

b) If  $\mathbf{v} \in \mathbb{R}^n$  and A is a real, symmetric, positive-definite  $n \times n$  matrix, then the set

$$\mathbb{D}(\mathbf{v},h) = \{\mathbf{x} \in \mathbb{R}^n | (\mathbf{x} - \mathbf{v})^\mathsf{T} A(\mathbf{x} - \mathbf{v}) \le h\}$$

is an ellipsoid with radius h centered at **v**. The eigenvectors of A are the principal axes of  $\mathbb{D}$ . Consider the Mahalanobis ellipsoid

$$\mathbb{D}_M(\boldsymbol{\mu}, h) = \{ \mathbf{x} \in \mathbb{R}^n | (\mathbf{x} - \boldsymbol{\mu})^\mathsf{T} \Sigma^{-1} \ (\mathbf{x} - \boldsymbol{\mu}) \le h \}$$
(14)

Assume that  $\mathbf{X} \sim N_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Find  $h_{0.05}$  such that

$$P\left(\mathbf{X} \in \mathbb{D}_M(\boldsymbol{\mu}, h_{0.05})\right) = 0.95.$$

Aid: See section 9.12.

#### Problem 3.3.

X is  $n \times k$  and rank  $X = r \le n < k$ . The BMP - inverse of X is

$$X^{-} = \left(X^{T}X\right)^{-}X^{T}.$$
(15)

Consider the OLS

 $\mathbf{Y} = X\boldsymbol{\beta}_* + \boldsymbol{\epsilon},$ 

where rank  $X = r \leq n \leq k$ .  $\hat{\boldsymbol{\beta}}^{\dagger} = (X^T X)^T \mathbf{y}$ . We have  $\| \mathbf{x} \| = \sqrt{\mathbf{x}^T \mathbf{x}} = \sqrt{\sum_{i=1}^n x_i^2}$ . Show that

$$\frac{1}{n} \mathbf{E} \parallel X \widehat{\boldsymbol{\beta}}^{\dagger} - X \boldsymbol{\beta}_* \parallel^2 = \frac{\sigma^2}{n} r.$$

You can use the facts that  $((X^T X)^{-})^T = ((X^T X)^{T})^{-}$  and  $\operatorname{Tr} X X^{-} = r$  without proving these.

#### Problem 3.4.

This assignment deals with Linear Minimal Mean Square Estimation. We have two  $n \times 1$  random vectors **X** and **Y** with zero expectation vectors. We consider a situation, where **Y** is a hidden and **X** is an observed measurement. We know the cross covariance matrix  $Cov(\mathbf{Y}, \mathbf{X}) = E[\mathbf{Y}\mathbf{X}^T]$  and the covariance matrix  $C_{\mathbf{X}}$ , where we assume that  $C_{\mathbf{X}}^{-1}$  exists. The goal is to find a linear estimator  $A\mathbf{X}$  of **Y** such that the mean square error

$$MSE = E \parallel \mathbf{Y} - A\mathbf{X} \parallel^2$$

is minimized. We have

$$\mathbf{E} \parallel \mathbf{Y} - A\mathbf{X} \parallel^{2} = \mathbf{E} \left[ (\mathbf{Y} - A\mathbf{X})^{T} (\mathbf{Y} - A\mathbf{X}) \right] = \mathrm{Tr} \mathbf{E} \left[ (\mathbf{Y} - A\mathbf{X}) (\mathbf{Y} - A\mathbf{X})^{T} \right]$$

We find the linear minimal mean square estimator by the following steps.

a) Check that if

$$A = \operatorname{Cov}(\mathbf{Y}, \mathbf{X}) \operatorname{C}_{\mathbf{X}}^{-1}, \tag{16}$$

then we have the orthogonality condition

$$\operatorname{Cov}(\mathbf{Y} - A\mathbf{X}, \mathbf{X}) = \operatorname{E}\left[\left(\mathbf{Y} - A\mathbf{X}\right)\mathbf{X}^{T}\right] = \mathbf{0}_{n}.$$
(17)

b) Let B be any  $n \times n$  matrix. Write

$$(\mathbf{Y} - B\mathbf{X})(\mathbf{Y} - B\mathbf{X})^{T} = ((\mathbf{Y} - A\mathbf{X}) + (A - B)\mathbf{X})(\mathbf{Y} - A\mathbf{X}) + (A - B)\mathbf{X})^{T}$$

where A is given by (16). Then we get by a direct expansion that

$$(\mathbf{Y} - B\mathbf{X})(\mathbf{Y} - B\mathbf{X})^T = (\mathbf{Y} - A\mathbf{X})(\mathbf{Y} - A\mathbf{X})^T$$

$$-(A-B)\mathbf{X}(\mathbf{Y}-A\mathbf{X})^{T} - (\mathbf{Y}-A\mathbf{X})((A-B)\mathbf{X})^{T} + (A-B)\mathbf{X}((A-B)\mathbf{X})^{T}$$

(you need not re-check this, in case you think this is correct). Then show that

$$E\left[\left(\mathbf{Y} - B\mathbf{X}\right)\left(\mathbf{Y} - B\mathbf{X}\right)^{T}\right] = E\left[\left(\mathbf{Y} - A\mathbf{X}\right)\left(\mathbf{Y} - A\mathbf{X}\right)^{T}\right] + E\left[\left(A - B\right)\mathbf{X}\left((A - B)\mathbf{X}\right)^{T}\right].$$
(18)

Aid: You will need (9) in the form

$$\operatorname{Cov}((A - B)\mathbf{X}, (\mathbf{Y} - A\mathbf{X})) = \operatorname{Cov}(\mathbf{Y} - A\mathbf{X}), (A - B)\mathbf{X}))^{T}.$$

and certain rules of covariance matrices of linear transformations.

c) Now show that

$$\mathbf{E} \parallel \mathbf{Y} - B\mathbf{X} \parallel^2 = \mathbf{E} \parallel \mathbf{Y} - A\mathbf{X} \parallel^2 + \mathbf{E} \parallel (A - B)\mathbf{X} \parallel^2$$

and draw the conclusion that A as in (16) gives the linear minimal mean square estimator  $A\mathbf{X}$ .

#### Problem 3.5.

We continue with Minimal Linear Mean Square Estimation. **X** and **Y** are  $n \times 1$  random vectors with zero expectation vectors. We have the cross covariance matrix  $\text{Cov}(\mathbf{Y}, \mathbf{X}) = E[\mathbf{Y}\mathbf{X}^T]$  and the covariance matrix  $C_{\mathbf{X}}$ , where we assume that  $C_{\mathbf{X}}^{-1}$  exists. By the preceding assignment above we have that

$$A = \operatorname{Cov}(\mathbf{Y}, \mathbf{X}) C_{\mathbf{X}}^{-1}, \tag{19}$$

minimizes the mean square error of estimating  $\mathbf{Y}$  by a linear map of  $\mathbf{X}$ .

Find that the **Linear Minimal Mean Square Error** (LMME) is

LMME = E 
$$\| \mathbf{Y} - A\mathbf{X} \|^2$$
 = Tr  $[C_{\mathbf{Y}} - Cov(\mathbf{Y}, \mathbf{X})C_{\mathbf{X}}^{-1}Cov(\mathbf{X}, \mathbf{Y})].$  (20)

Aid: Start with

$$\mathbf{E} \parallel \mathbf{Y} - A\mathbf{X} \parallel^{2} = \mathbf{E} \left[ \left( \mathbf{Y} - A\mathbf{X} \right)^{T} \left( \mathbf{Y} - A\mathbf{X} \right) \right] = \mathrm{Tr} \, \mathbf{E} \left[ \left( \mathbf{Y} - A\mathbf{X} \right) \left( \mathbf{Y} - A\mathbf{X} \right)^{T} \right],$$

expand, compute the expectations and use the rules for traces in section 9.4. A Solution to this Problem is found in Section 9.16.

## 4 Linear Regression

## 4.1 Simple Linear Regression

Problem 4.1.



In the left hand field you see a plot of a training set with predictor x = body height (in feet and inches) and y = score in mathematics exams in an US elementary school and a least squares regression line fitted to it. One might conclude that taller students perform better in mathematics. In the right hand field we see the same data and regression line after an additional analysis. Explain what is found now! What is the phenomenon known as?

#### Problem 4.2.

In (109) we have the LSE of the regression coefficient as

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n \left(y_i - \overline{y}\right) x_i}{\sum_{i=1}^n \left(x_i - \overline{x}\right) x_i}.$$

Rewrite this as

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n \left(y_i - \overline{y}\right) \left(x_i - \overline{x}\right)}{\sum_{i=1}^n \left(x_i - \overline{x}\right)^2}$$

by means of the rules in section 9.8.

## Problem 4.3.

 $\hat{\beta}_0$  and  $\hat{\beta}_1$  are the LSE of  $\beta_0$  and  $\beta_1$  in simple linear regression.  $\beta_0^*, \beta_1^*$  denote the true poulation values of the intercept and the regression coefficient.

a) Show that

$$\operatorname{Cov}\left(\hat{\beta}_{0},\hat{\beta}_{1}\right) = \frac{-\bar{x}\sigma^{2}}{\sum_{i=1}^{n}\left(x_{i}-\bar{x}\right)^{2}}$$
(21)

What does this mean in terms of the estimated regression line? Aid(1): the definition of Cov is

$$\operatorname{Cov}\left(\hat{\beta}_{0},\hat{\beta}_{1}\right)=E\left[\left(\hat{\beta}_{0}-\beta_{0}^{*}\right)\left(\hat{\beta}_{1}-\beta_{1}^{*}\right)\right]$$

Aid (2): You can use the means and variances in (110) even if we are not assuming normal regression.

b) Let us assume normal linear regression. What is the joint distribution of bivariate random variable

$$\left(\hat{\beta}_0, \hat{\beta}_1\right)$$

Recapitulate explicitly the mean vector and the covariance matrix.

## Problem 4.4.

We consider simple linear regression in the equivalent centered form:

$$Y_i = \alpha + \beta_1 \left( x_i - \bar{x} \right) + \varepsilon_i, \quad i = 1, \dots, n$$
(22)

where  $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$  and

$$\alpha = \beta_0 + \beta_1 \bar{x}.\tag{23}$$

In the matrix form this involves

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, X_c = \begin{pmatrix} 1 & x_1 - \bar{x} \\ 1 & x_2 - \bar{x} \\ \vdots \\ 1 & x_n - \bar{x} \end{pmatrix}, \boldsymbol{\beta}_c = \begin{pmatrix} \alpha \\ \beta_1 \end{pmatrix}.$$

The least square criterion (cost function) to be minimized is

$$Q\left(\boldsymbol{\beta}_{c}\right) = \|\mathbf{y} - X_{c}\boldsymbol{\beta}_{c}\|^{2}.$$

Hence we know that

$$\hat{\boldsymbol{\beta}}_c = (X_c^T X_c)^{-1} X_c^T \mathbf{y}.$$

a) Show by a direct computation that

$$X_c^T X_c = \left(\begin{array}{cc} n & 0\\ 0 & \sum_{i=1}^n \left(x_i - \bar{x}\right)^2 \end{array}\right).$$

b) It is now clear how to find  $(X_c^T X_c)^{-1}$ , if we assume that not all  $x_i$  are equal. Show under this assumption that the LSE of the parameters of the centered model are

$$\hat{\boldsymbol{\beta}}_{c} = \begin{pmatrix} \hat{\alpha} \\ \hat{\beta}_{1} \end{pmatrix} = \begin{pmatrix} \bar{y} \\ \frac{\sum_{i=1}^{n} (x_{i} - \bar{x}) y_{i}}{\sum_{i=1}^{n} (x_{i} - \bar{x})^{2}} \end{pmatrix}.$$
(24)

c) Derive the predictor equation for  $\hat{y}_c$  in the centered simple linear regression. Is this equivalent to the predictor  $\hat{y}$  in the uncentered model?

## Problem 4.5.

Let  $\hat{\alpha}$  and  $\hat{\beta}_1$  be given as in (24).

a) Show that

$$\operatorname{Cov}\left(\hat{\alpha},\hat{\beta}_{1}\right)=0\tag{25}$$

b) Suppose that  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$  in (22) is a normal random vector, i.e.,

$$\boldsymbol{\varepsilon} \sim N_n \left( \mathbf{0}_n, \sigma^2 \mathbb{I}_n \right).$$
 (26)

What does (25) imply in this case?

#### Problem 4.6.

Sometimes, with regards to content, it can be reasonable to assume that the regression line passes through the origin ( $\beta_0 = 0$ ). The corresponding regression model is

$$y_i = \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n$$

with the following assumptions:

$$E[\varepsilon_i] = 0, \quad i = 1, \dots, n$$
  

$$Var(\varepsilon_i) = \sigma^2, \quad i = 1, \dots, n$$
  

$$\{\varepsilon_i \mid i = 1, \dots, n\} \text{ stochastically independent}$$
  

$$\varepsilon_i \sim N(0, \sigma^2), \quad i = 1, \dots, n$$

- a) Derive the LS-estimator  $\hat{\beta}_1$  for  $\beta_1$ . Is  $\hat{\beta}_1$  also the ML-estimator?
- b) Show that the residuals in general do not sum up to zero. However, why does this hold in the linear regression model including the intercept  $\beta_0$ ? Explain the difference.

#### Problem 4.7.

Suppose that the components of the random vector  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$  are independent with means = 0, but that

$$C_{\varepsilon} = \sigma^{2} \begin{pmatrix} \frac{1}{w_{1}} & 0 & 0 & \dots & 0\\ 0 & \frac{1}{w_{2}} & 0 & \dots & 0\\ 0 & \ddots & \vdots & \dots & 0\\ 0 & 0 & 0 & \dots & \frac{1}{w_{n}} \end{pmatrix}.$$
 (27)

where the  $w_i = w(x_i)$  are the values of some known positive weight function evaluated at the regressor values.

We consider simple theoretic linear regression in a centered form:

$$E[Y_i] = \beta_0 + \beta_1 (x_i - \bar{x}(w)), \quad i = 1, 2, \dots, n,$$

where

$$\bar{x}(w) = \frac{\sum_{i=1}^{n} w_i x_i}{\sum_{i=1}^{n} w_i}$$

is the weighted mean. We minimize the weighted LS-criterion

$$Q_W(\beta_0, \beta_1) = \sum_{1=1}^n w_i \left( y_i - (\beta_0 + \beta_1 \left( x_i - \bar{x}(w) \right) \right)^2$$

a) Check first that

$$\frac{\sum_{i=1}^{n} w_i \left( x_i - \bar{x}(w) \right)}{\sum_{i=1}^{n} w_i} = 0.$$
(28)

This can be useful in b), c) and/or d) below. You are allowed to use (28) there, even if you have failed to establish (28).

b) Check that the weighted LSE is

$$\hat{\boldsymbol{\beta}} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \begin{pmatrix} \frac{\sum_{i=1}^n w_i y_i}{\sum_{i=1}^n w_i} \\ \frac{\sum_{i=1}^n w_i (x_i - \bar{x}(w)) y_i}{\sum_{i=1}^n w_i (x_i - \bar{x}(w))^2} \end{pmatrix}.$$
(29)

c) Show that

$$\hat{\beta}_0 = \frac{\sum_{i=1}^n w_i Y_i}{\sum_{i=1}^n w_i}$$

and

$$\hat{\beta}_{1} = \frac{\sum_{i=1}^{n} w_{i} \left( x_{i} - \bar{x}(w) \right) Y_{i}}{\sum_{i=1}^{n} w_{i} \left( x_{i} - \bar{x}(w) \right)^{2}}$$

are unbiased.

d) Show that

$$\operatorname{Var}[\hat{\beta}_0] = \frac{\sigma^2}{\sum_{i=1}^n w_i}$$

and determine in similar manner  $\operatorname{Var}[\hat{\beta}_1]$ . Remember to justify your calculations.

## Problem 4.8.

Suppose that the components of the random vector  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$  are independent with means = 0, but that

$$C_{\varepsilon} = \sigma^2 \begin{pmatrix} \lambda_1 & 0 & \cdots & 0\\ 0 & \lambda_2 & \cdots & 0\\ \vdots & \vdots & \ddots & \vdots\\ 0 & 0 & \cdots & \lambda_n \end{pmatrix},$$
(30)

where  $\lambda_i > 0$ . We observe outcomes  $y_1, \ldots, y_n$  of the true model (a constant with noise)

$$\mathbf{Y} = \mathbf{1}_n \beta_0 + \boldsymbol{\varepsilon}.$$

Now we minimize the weighted least squares criterion

$$Q_{\lambda}(\beta_0) := \sum_{i=1}^{n} \frac{1}{\lambda_i} \left( y_i - \beta_0 \right)^2 \tag{31}$$

as a function of  $\beta_0$ .

a) Check that the weighted LSE is

$$\hat{\beta}_0 = \frac{1}{\sum_{i=1}^n \frac{1}{\lambda_i}} \sum_{i=1}^n \frac{1}{\lambda_i} y_i.$$
(32)

b) Show that

$$\hat{\beta}_0 = \frac{1}{\sum_{i=1}^n \frac{1}{\lambda_i}} \sum_{i=1}^n \frac{1}{\lambda_i} Y_i$$

is unbiased.

c) Show that

$$\operatorname{Var}\left[\hat{\beta}_{0}\right] = \frac{\sigma^{2}}{\sum_{i=1}^{n} \frac{1}{\lambda_{i}}}.$$

Remember to justify your calculations.

#### Problem 4.9.

The phenomenon of fatigue was observed by engineers in metal materials in the mid-19th century, due to the shortcomings of the railway wagons after a short period of work. In the 1850s, August Wöhler, a railway engineer, conducted the first experimental fatigue program, testing wagon shafts for failure under alternating stress. The developed loads were recorded together with the number of rotations until failure, so it was possible to formulate the first S-N diagram (S = Stress, N = Number of Cycles) or Wöhler-curve. This assignment deals with what is called the statistical recalibration of the Wöhler curve. Let

Y = Number of life cycles to failure, X = stress.

The linear part of the curve is stated in log-log terms as follows. We assume that Y, given X = x, has the lognormal distribution

$$(Y|X=x) \sim \operatorname{LN}\left(\beta_0 + \beta_1 \ln(x), \sigma^2\right)$$

where x is greater than some  $x_0 > 0$ . There is a recorded set of data, n pairs of stress and life-cycles,  $D_{tr} = \{x_i, y_i\}_{i=1}^n$ .

- a) Find simple linear regression equations for this training set, where  $\beta_0$  is the intercept and  $\beta_1$  is the slope.
- b) Write down the LSEs  $\hat{\beta}_0$  and  $\hat{\beta}_1$ .
- c) What are the distributions of  $\hat{\beta}_0$  and  $\hat{\beta}_1$ ?
- e) What is now  $\widehat{E}[Y|X = x]$ , your estimate of the expected number of life cycles given the stress x?

References:

- Barbosa, Joelton Fonseca and Correia, José AFO and Freire Junior, RCS and Zhu, Shun-Peng and De Jesus, Abílio MP: Probabilistic SN fields based on statistical distributions applied to metallic and composite materials: State of the art. Advances in Mechanical Engineering, 11, 8, 2019.
- 2. A positive random variable X is log-normally distributed, written as

$$X \sim \mathrm{LN}(\mu, \sigma^2),$$

if the natural logarithm of X is normally distributed with mean  $\mu$  and variance  $\sigma^2$ , i.e.,





Fig.1 Typical fatigue endurance test data illustrating deviations from linear S-N curve.

## 4.2 Multiple Linear Regression

In every problem in this section and elsewhere, unless explicitly stated otherwise, it is assumed that the design matrix X has full column rank = k + 1 (< n).

### Problem 4.10.

In the OLS model,

$$\hat{\mathbf{Y}} = X\hat{\boldsymbol{\beta}} = H\mathbf{Y}$$

is the predictor vector, where H is the hat matrix. Show that

$$\operatorname{Cov}\left(\mathbf{Y},\mathbf{Y}-\hat{\mathbf{Y}}\right) = \sigma^{2}\left(\mathbb{I}_{n}-H\right).$$

Aid: Recall (10) and (12).

## Problem 4.11.

In the OLS model,

$$\hat{\mathbf{Y}} = X\hat{\boldsymbol{\beta}} = H\mathbf{Y}$$

is the predictor vector, where H is the hat matrix. Show that

$$\operatorname{Cov}\left(\hat{\mathbf{Y}},\mathbf{Y}-\hat{\mathbf{Y}}\right)=\mathbf{0}_{n}.$$

Aid: Recall (10) and (12).

### Problem 4.12.

We have the OLS

$$\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \widehat{\boldsymbol{\varepsilon}} = \mathbf{Y} - \widehat{\mathbf{Y}},$$

and

$$\mathbf{e}_{LSE} = \mathbf{y} - X\widehat{\boldsymbol{\beta}} = \mathbf{y} - H\mathbf{y} = \mathbf{y} - \widehat{\mathbf{y}}.$$

Explain the differences of concept in the OLS model between  $\varepsilon$ ,  $\hat{\varepsilon}$  and  $\mathbf{e}_{LSE}$ .

## Problem 4.13.

a) The hat matrix H is defined by

$$H = X \left( X^T X \right)^{-1} X^T.$$

Check that H is a projection matrix, see section 9.3.

b) We have

$$\mathbf{e}_{LSE} = \mathbf{y} - H\mathbf{y}$$

Show now that

$$\mathbf{e}_{LSE}^T H \mathbf{y} = 0.$$

Aid: You may need the idempotency of H.

c) Why is *H* called the hat matrix ? Answer: It puts a hat on **y**, i.e.,

$$\hat{\mathbf{y}} = H\mathbf{y}$$

d) Show that

$$\|\mathbf{y}\|^2 = \|\hat{\mathbf{y}}\|^2 + \|\mathbf{e}_{LSE}\|^2.$$
(33)

How can this be interpreted geometrically?

## Problem 4.14.

A training set of observed responses  $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$  and the corresponding  $n \times (k+1)$  data matrix X are available from some source with the true model

$$\mathbf{Y} = X\boldsymbol{\beta}_* + \boldsymbol{\varepsilon}.$$

Here  $\boldsymbol{\varepsilon}$  has he mean vector  $= \mathbf{0}_n$  and covariance matrix  $\sigma^2 \mathbb{I}_n$ . We fit a multiple least squares model to the training set and obtain the LSE  $\hat{\boldsymbol{\beta}}$ .

The  $m \times 1$  random vector  $\mathbf{Y}_o$  is the response vector given by

$$\mathbf{Y}_o = X_o \boldsymbol{\beta}_* + \boldsymbol{\varepsilon}_o$$

where  $\boldsymbol{\varepsilon}_0$  has mean vector  $= \mathbf{0}_m$  and covariance matrix  $\sigma^2 \mathbb{I}_m$ .  $\boldsymbol{\varepsilon}_0$  is independent of  $\boldsymbol{\varepsilon}$ .  $X_o$  is a  $m \times (k+1)$  matrix of full column rank.

Before observing  $\mathbf{Y}_o$ , we want to predict  $\mathbf{Y}_o$  by the predictor

$$\hat{\mathbf{Y}}_o = X_o \hat{\boldsymbol{\beta}}$$

- a) Show that  $E\left[\mathbf{Y}_{o}-\widehat{\mathbf{Y}}_{o}\right]=\mathbf{0}_{m}.$
- b) Find the covariance matrix  $C_{\varepsilon_{pr}}$  of the prediction error  $\varepsilon_{pr} := \mathbf{Y}_o \widehat{\mathbf{Y}}_o$ . Answer :  $C_{\varepsilon_{pr}} = \sigma^2 \left( X_o (X^T X)^{-1} X_o^T + \mathbb{I}_m \right)$ . Aid: You can have use of (54) even if no normal distribution is involved here.

#### Problem 4.15.

We consider the OLS model with the following generalizations:

1) <u>Correlated noise</u>: The covariance matrix  $C_{\boldsymbol{\varepsilon}}$  of  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$  is

$$C_{\varepsilon} = \sigma^2 \mathbb{C},\tag{34}$$

where  $\mathbb{C}$  is a known positive definite matrix.

2) Neither  $\sigma^2$  nor  $\mathbb{C}$  depend on X.

Questions:

a) Introduce a new random vector

$$\mathbf{Z} := \mathbb{C}^{-1/2} \mathbf{Y}$$

Why is this possible? Quote briefly some relevant facts from linear algebra.

b) Show that we have the multiple regression

$$\mathbf{Z} = \left(\mathbb{C}^{-1/2}X\right)\boldsymbol{\beta} + \boldsymbol{\varepsilon}^{\dagger},\tag{35}$$

where  $\boldsymbol{\varepsilon}^{\dagger}$  is a random vector satisfying  $E\left[\boldsymbol{\varepsilon}^{\dagger}\right] = \mathbf{0}_{n}, C_{\boldsymbol{\varepsilon}^{\dagger}} = \sigma^{2}\mathbb{I}_{n}.$ 

c) The response data in your training set are now  $z_i = C^{-1/2}y_i$ , i = 1, ..., n,

$$\mathbf{z} = \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{pmatrix}.$$

Show (justify your solution carefully) now that LSE using (35) is

$$\hat{\boldsymbol{\beta}} = \left( X^T \mathbb{C}^{-1} X \right)^{-1} X^T \mathbb{C}^{-1} \mathbf{y}.$$
(36)

d) The covariance matrix in (30) is a special case of (34). Check that you get (32) from (36).

## Problem 4.16.

In ridge regression, the ordinary LSE is replaced by  $\widehat{\boldsymbol{\beta}}_{\lambda}$  defined by

$$\widehat{\boldsymbol{\beta}}_{\lambda} := \operatorname{argmin}_{\boldsymbol{\beta}} \left( \|\mathbf{y} - X\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|^2 \right).$$

Using the rules of gradients for vector and matrix expressions found in Appendix 9.15, show that (p = k + 1)

$$\widehat{\boldsymbol{\beta}}_{\lambda} = \left( X^T X + \lambda \mathbb{I}_p \right)^{-1} X^T \mathbf{y}.$$
(37)

Whenever relevant, refer to the number of the rule you are evoking.

## Problem 4.17.

The k regressors of an ordinary LS model as defined in section 9.10

$$\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} \tag{38}$$

are **rescaled** as follows. The *i*th row in X is  $\mathbf{x}_i^{\mathsf{T}} = (1, x_{i1}, x_{i2}, \dots, x_{ik})$  and the rescaled row is  $\mathbf{z}_i^{\mathsf{T}} = (1, c_1 x_{i1}, c_2 x_{i2}, \dots, c_k x_{ik})$ , where  $c_i \neq 0$  for  $i = 1, 2, \dots, k$ . Let us define the  $(k+1) \times (k+1)$  diagonal matrix

$$D = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & c_1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & c_k \end{pmatrix}.$$

Then it follows that

$$Z = \begin{pmatrix} \mathbf{z}_1^{\mathsf{T}} \\ \mathbf{z}_2^{\mathsf{T}} \\ \vdots \\ \mathbf{z}_n^{\mathsf{T}} \end{pmatrix} = XD.$$
(39)

Next you regress the response data  $\mathbf{y} = (y_1, \dots, y_n)^T$  on Z and obtain

$$\widehat{\boldsymbol{\beta}}_{z} = \left(Z^{T} Z\right)^{-1} Z^{T} \mathbf{y} \tag{40}$$

a) Show that

 $\widehat{\boldsymbol{\beta}}_{z} = D^{-1}\widehat{\boldsymbol{\beta}},$ 

where  $\widehat{\boldsymbol{\beta}}$  is the LSE for regressing the response data  $\mathbf{y} = (y_1, \ldots, y_n)^T$  on X. Aid: Use (39) in (40) and simplify.

b) Show that the hat predictor is invariant w.r.t. the rescaling above in the sense that

 $\widehat{\mathbf{y}}_z = \widehat{\mathbf{y}}.$ 

Here  $\hat{\mathbf{y}}_z$  is the predictor of  $\mathbf{y}$  based on  $\hat{\boldsymbol{\beta}}_z$  and  $\hat{\mathbf{y}}$  is the predictor of  $\mathbf{y}$  based on  $\hat{\boldsymbol{\beta}}$ .

Problem 4.18.

$$\mathbf{e}_{LSE} = \mathbf{y} - X\widehat{\boldsymbol{\beta}} = \mathbf{y} - H\mathbf{y} = \mathbf{y} - \widehat{\mathbf{y}}$$

is the vector of LSE residuals.

$$\mathbf{e}_{LSE} = (\widehat{e}_1, \widehat{e}_2, \dots, \widehat{e}_n)^T.$$

a) Show that  $X^T \mathbf{e}_{LSE} = \mathbf{0}_k$ . When you look at the scalar product of the first row in  $X^T$  and  $\mathbf{e}_{LSE}$  this means

$$\sum_{i=1}^{n} \hat{e}_i = 0.$$
 (41)

b) Show now that

$$\frac{1}{n}\sum_{i=1}^{n}\widehat{y}_i = \bar{y}.$$
(42)

c) At a certain moment in the lectures it was obtained that

$$\mathbf{y}^T C_{ce} C_{ce} \mathbf{y} = \widehat{\mathbf{y}}^T C_{ce} \widehat{\mathbf{y}} + \widehat{\mathbf{y}}^T \mathbf{e}_{LSE} + \mathbf{e}_{LSE}^T \widehat{\mathbf{y}} + \mathbf{e}_{LSE}^T \mathbf{e}_{LSE}.$$

We know that  $\mathbf{y}^T C_{ce} C_{ce} \mathbf{y} = \sum_{i=1}^n (y_i - \bar{y})^2 = SS_T$  by (7) and section 9.14. Show no that

$$SS_{\rm R} + SS_{\rm Res} = \hat{\mathbf{y}}^T C_{ce} \hat{\mathbf{y}} + \hat{\mathbf{y}}^T \mathbf{e}_{LSE} + \mathbf{e}_{LSE}^T \hat{\mathbf{y}} + \mathbf{e}_{LSE}^T \mathbf{e}_{LSE}$$

Thus we have established the identity (117).

## 4.3 Estimability and the Gauss-Markov Theorem

The Gauss-Markov theorem states that LSE is best linear unbiased estimate of the regression coefficients. In this subsection the Gauss-Markov theorem is studied by an approach different from the one applied in Lecture 4 slides and in MVP. The material below is found in Chapter 2.8 of Julian R. Faraway: *Practical regression with R*, 2002, which is downloadable from canvas via the page SF2930 Course Plan. The proof is based on the notion of estimable linear combinations.

#### Problem 4.19.

An ordinary LS model as defined in section 9.10

$$\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} \tag{43}$$

where  $\boldsymbol{\beta}$  is  $(k+1) \times 1$ . We let  $\boldsymbol{c}$  be an arbitrary  $(k+1) \times 1$  vector in  $\mathbb{R}^{k+1}$ . Next we introduce the  $1 \times n$  vector  $\boldsymbol{a}^T$  by

$$\boldsymbol{a}^T := \boldsymbol{c}^T (X^T X)^{-1} X^T.$$
(44)

a) Show that

$$\mathbf{E}\left[\boldsymbol{a}^{T}\mathbf{Y}\right] = \boldsymbol{c}^{T}\boldsymbol{\beta}.$$
(45)

One says that  $c^T \beta$  is estimable, since there is an unbiased estimator  $a^T \mathbf{Y}$  of it.

b) Check that we also have

$$\mathbf{E}\left[\boldsymbol{a}^{T}\mathbf{Y}\right] = \boldsymbol{a}^{T}X\boldsymbol{\beta}.$$
(46)

#### Problem 4.20.

Here you need the notion of the column space or range  $\mathcal{R}(X)$  of a matrix X, as presented in the Appendix 9.2.

a) Verify that if (45) and (46) are valid for every  $\beta$ , then

$$\boldsymbol{c} \in \mathcal{R}\left(X^{T}\right)$$

b) X is an  $n \times (k+1)$  matrix with full column rank = k+1. Show now that

$$\boldsymbol{c} \in \mathcal{R}\left(X^T X\right). \tag{47}$$

Aid: In other words, you are through with this, once you have verified that

$$\mathcal{R}\left(X^{T}\right) \subset \mathcal{R}\left(X^{T}X\right).$$

This is the second part of the proof of the statement in Proposition 9.1. We give a stepwise guide for this.

1. Take any  $\mathbf{w} \in \mathcal{R}(X^T)$ . Thus there exists a  $\mathbf{b} \in \mathbb{R}^n$  such that

$$\mathbf{w} = X^T \mathbf{b}.$$

Define by use of the hat matrix the projected vector  $\widehat{\mathbf{b}} \in \mathcal{R}(X)$  by

$$\widehat{\mathbf{b}} = H\mathbf{b}.$$

Check now that

$$X^T \widehat{\mathbf{b}} = X^T \mathbf{b}. \tag{48}$$

2. Since  $\widehat{\mathbf{b}} \in \mathcal{R}(X)$ , there exists  $\mathbf{v}$  such that  $\widehat{\mathbf{b}} = X\mathbf{v}$ . Now compute

 $X^T X \mathbf{v}$ 

and draw the desired conclusion.

## Problem 4.21.

Now we prove the first part of the Gauss-Markov theorem. We are relying on (47), which you can invoke without a successfull proof thereof. We have seen in the first Problem of this subsection that  $\boldsymbol{c}^T\boldsymbol{\beta}$  is estimable, if X has full column rank. Let thus  $\boldsymbol{a}^T\mathbf{Y}$  be any unbiased estimator of  $\boldsymbol{c}^T\boldsymbol{\beta}$ . By the preceding problems, (45) and (46) hold and imply (47). By definition of  $\mathcal{R}(X^TX)$ , (47) implies that there exists a  $(k+1) \times 1$  vector  $\boldsymbol{z}$  such that

$$\boldsymbol{c} = \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{z}. \tag{49}$$

If  $\hat{\boldsymbol{\beta}}$  is the LSE of (43), then (49) implies

$$\boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}} = \boldsymbol{z}^{T}X^{T}X\widehat{\boldsymbol{\beta}} = \boldsymbol{z}^{T}X^{T}\mathbf{Y}.$$
(50)

Then we get

$$\operatorname{Var}\left[\boldsymbol{a}^{T}\mathbf{Y}\right] = \operatorname{Var}\left[\boldsymbol{a}^{T}\mathbf{Y} - \boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}} + \boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}\right] = \operatorname{Var}\left[\left(\boldsymbol{a}^{T}\mathbf{Y} - \boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}\right) + \boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}\right]$$
$$= \operatorname{Var}\left[\boldsymbol{a}^{T}\mathbf{Y} - \boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}\right] + \operatorname{Var}\left[\boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}\right] + 2\operatorname{Cov}\left(\boldsymbol{a}^{T}\mathbf{Y} - \boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}, \boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}\right).$$

where we used a familiar formula for the variance of a sum of two random variables from the first course.

a) Use to (50) to write

$$\operatorname{Cov}\left(\boldsymbol{a}^{T}\mathbf{Y}-\boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}},\boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}\right)=\operatorname{Cov}\left(\boldsymbol{a}^{T}\mathbf{Y}-\boldsymbol{z}^{T}X^{T}\mathbf{Y},\boldsymbol{z}^{T}X^{T}\mathbf{Y}\right)$$

Then show that

$$\operatorname{Cov}\left(\boldsymbol{a}^{T}\mathbf{Y}-\boldsymbol{z}^{T}X^{T}\mathbf{Y},\boldsymbol{z}^{T}X^{T}\mathbf{Y}\right)=0.$$

Aid: You may need the standard rule of covariance computation (11) and (47).

b) Hence we have

$$\operatorname{Var}\left[\boldsymbol{a}^{T}\mathbf{Y}\right] = \operatorname{Var}\left[\boldsymbol{a}^{T}\mathbf{Y} - \boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}\right] + \operatorname{Var}\left[\boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}\right].$$
(51)

Now draw the conclusion in the Gauss-Markov theorem.

c) We need to establish the uniqueness of  $\hat{\boldsymbol{\beta}}$  as the minimum variance unbiased linear estimator. This means that if  $\operatorname{Var} \left[ \boldsymbol{a}^T \mathbf{Y} \right] = \operatorname{Var} \left[ \boldsymbol{c}^T \hat{\boldsymbol{\beta}} \right]$ , then  $\boldsymbol{a}^T \mathbf{Y} = \boldsymbol{c}^T \hat{\boldsymbol{\beta}}$ . *Aid:* See what happens in (51), if  $\operatorname{Var} \left[ \boldsymbol{a}^T \mathbf{Y} \right] = \operatorname{Var} \left[ \boldsymbol{c}^T \hat{\boldsymbol{\beta}} \right]$ .

## Problem 4.22.

If X has full column rank, rank X = k + 1 < n, check that

$$X^+ = (X^T X)^{-1} X^T$$

is the BMP inverse of X by verifying the conditions MP1-MP4 in section 9.7.

## Problem 4.23.

In this problem X has full row rank n < k + 1.

If X has full row rank, rank X = n < k + 1, then

$$X^+ = X^T (XX^T)^{-1}$$

is the BMP inverse of X.

a) We know the normal equations:

$$X^T X \boldsymbol{\beta} = X^T \mathbf{y}. \tag{52}$$

For X with full row rank we define

$$\boldsymbol{\beta}^+ = X^+ \mathbf{y}.$$

Check now that  $\beta^+$  is a solution to (52).

b) We have the BMP predictor

$$\widehat{\mathbf{y}}^+ = X\boldsymbol{\beta}^+ = XX^+\mathbf{y}$$

Show that  $\widehat{\mathbf{y}}^+$  is an interpolation of the training set.

## Problem 4.24.

 $\mathbf{Y} = X\boldsymbol{\beta}_* + \boldsymbol{\varepsilon}$  $\mathbf{E}\left[\boldsymbol{\varepsilon}\right] = \mathbf{0}_n, C_{\boldsymbol{\varepsilon}} = \sigma^2 \mathbb{I}_n.$ 

and

$$\widehat{\mathbf{Y}} = H\mathbf{Y}, \widehat{\boldsymbol{\varepsilon}} = \mathbf{Y} - \widehat{\mathbf{Y}}.$$

Let the variance  $\sigma^2$  be estimated by

$$\widehat{\sigma^2} = \frac{1}{n-k-1}\widehat{\boldsymbol{\varepsilon}}^T\widehat{\boldsymbol{\varepsilon}}.$$

Show that  $\widehat{\sigma^2}$  is unbiased.

Aid: Insert  $\hat{\boldsymbol{\varepsilon}} = \mathbf{Y} - \hat{\mathbf{Y}}$  compute the required expectation and use the rules for traces in section 9.4.

## 4.4 Normal Linear Regression

#### Problem 4.25.

Prove in formal detail:

$$\mathbf{Y} \sim N_n \left( X \boldsymbol{\beta}_*, \sigma^2 \mathbb{I}_n \right).$$

## Problem 4.26.

a) Show that

b) Show that

The LSE of the regressor parameters is

$$\widehat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{Y}.$$

$$\widehat{\boldsymbol{\beta}} = \boldsymbol{\beta}_* + (X^T X)^{-1} X^T \boldsymbol{\varepsilon}.$$

$$\widehat{\boldsymbol{\beta}} \sim N_{k+1} \left( \boldsymbol{\beta}_*, \sigma^2 (X^T X)^{-1} \right).$$
(53)

Remember to justify carefully.

Problem 4.27.

 $\widehat{\mathbf{Y}} = H\mathbf{Y},$ 

where H is the hat matrix.

a) Show that

$$\widehat{\mathbf{Y}} \sim N_n \left( X \boldsymbol{\beta}_*, \sigma^2 H \right).$$
 (54)

b)

 $\widehat{\boldsymbol{\varepsilon}} = \mathbf{Y} - \widehat{\mathbf{Y}}.$ 

Show (with appropriate justification) that

$$\widehat{\boldsymbol{\varepsilon}} \sim N_n \left( \mathbf{0}_n, \sigma^2 \left( \mathbb{I}_n - H \right) \right) \tag{55}$$

Problem 4.28.

$$\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$
$$\boldsymbol{\varepsilon} \sim N_n \left( \mathbf{0}, \sigma^2 \mathbb{I}_n \right)$$

Thus

where

$$\mathbf{Y} \sim N_n \left( X \boldsymbol{\beta}, \sigma^2 \mathbb{I}_n \right)$$

and the p.d.f .is

$$f_{\mathbf{Y}}(\mathbf{y}) = \frac{1}{\sqrt{(\sigma^2 2\pi)^n}} e^{-\frac{1}{2\sigma^2} \|\mathbf{y} - X\boldsymbol{\beta}\|^2}$$

We have the -1 · loglikelihood function

$$l_{\mathbf{y}}\left(\boldsymbol{\beta},\sigma^{2}\right) := -\ln f_{\mathbf{Y}}\left(\mathbf{y}\right) = \frac{n}{2}\ln(2\pi) + \frac{n}{2}\ln(\sigma^{2}) + \frac{1}{2\sigma^{2}} \| \mathbf{y} - X\boldsymbol{\beta} \|^{2}$$
(56)

For  $j = 0, 1, \ldots, k$  it holds that

$$\frac{\partial}{\partial\beta_j} l_{\mathbf{y}} \left( \boldsymbol{\beta}, \sigma^2 \right) = \frac{1}{2\sigma^2} \frac{\partial}{\partial\beta_j} \parallel \mathbf{y} - X\boldsymbol{\beta} \parallel^2 = \frac{1}{2\sigma^2} \frac{\partial}{\partial\beta_j} Q \left( \boldsymbol{\beta} \right)$$

Hence the maximum likelihood estimate (MLE)  $\hat{\beta}_{ML}$  can be found by solving first w.r.t.  $\beta$  the equations

$$\frac{1}{2}\frac{\partial Q\left(\boldsymbol{\beta}\right)}{\partial\boldsymbol{\beta}}=\mathbf{0}_{k}$$

Since  $\mathbf{y}^T X \boldsymbol{\beta} = \boldsymbol{\beta}^T X^T \mathbf{y}$  we get

$$Q\left(\boldsymbol{\beta}\right) = \boldsymbol{\beta}^{T} X^{T} X \boldsymbol{\beta} - 2 \boldsymbol{\beta}^{T} X^{T} \mathbf{y} + \mathbf{y}^{T} \mathbf{y}.$$
(57)

a) Show now that

$$\frac{1}{2} \frac{\partial Q\left(\boldsymbol{\beta}\right)}{\partial \boldsymbol{\beta}} |_{\boldsymbol{\widehat{\beta}}_{ML}} = \mathbf{0}_k \Leftrightarrow X^T X \boldsymbol{\widehat{\beta}}_{ML} = X^T \mathbf{y}$$

Aid: Use the rules in section 9.15.

b) Hence, if X has full column rank,  $\hat{\beta}_{ML}$  equals the LSE  $\hat{\beta}$  and the LSE residuals are

$$\mathbf{e}_{LSE} = \mathbf{y} - X \widehat{\boldsymbol{\beta}}_{ML}.$$

Therefore

$$Q\left(\widehat{\boldsymbol{\beta}}_{ML}\right) = \parallel \mathbf{y} - X\widehat{\boldsymbol{\beta}}_{ML} \parallel^2 = \mathbf{e}_{LSE}^T \mathbf{e}_{LSE}.$$

When we substitute this in (56) we get

$$l_{\mathbf{y}}\left(\widehat{\boldsymbol{\beta}}_{ML},\sigma^{2}\right) = \frac{n}{2}\ln(2\pi) + \frac{n}{2}\ln(\sigma^{2}) + \frac{1}{2\sigma^{2}}\mathbf{e}_{LSE}^{T}\mathbf{e}_{LSE}.$$

Find now the MLE  $\hat{\sigma}_{ML}^2$  of  $\sigma^2$ .

c) Test the correctness of your result in b) by verifying that

$$l_{\mathbf{y}}\left(\widehat{\boldsymbol{\beta}}_{ML}, \widehat{\sigma}_{ML}^2\right) = \frac{n}{2}\ln(2\pi) + \frac{n}{2}\ln(\widehat{\sigma}_{ML}^2) + \frac{n}{2}$$
(58)

## Problem 4.29.

$$\mathbf{Y} = X \boldsymbol{\beta}_* + \boldsymbol{\varepsilon} \quad \boldsymbol{\varepsilon} \sim N_n \left( \mathbf{0}, \sigma^2 \mathbb{I}_n \right)$$

We have used a training data set

$$\mathcal{D}_{tr} = \left\{ \left( y_i, \mathbf{x}_i^T \right)_{i=1}^n \right\}$$

to find the LSE  $\widehat{\boldsymbol{\beta}}$ , which determines the hat matrix so that

$$\widehat{\mathbf{Y}} = H\mathbf{Y}, \quad \widehat{\boldsymbol{\varepsilon}} = \mathbf{Y} - \widehat{\mathbf{Y}}.$$

a) Show that

$$\widehat{\boldsymbol{\varepsilon}} = (\mathbb{I}_n - H)\boldsymbol{\varepsilon}. \tag{59}$$

b) Set

$$\mathbf{Z} := \frac{\boldsymbol{\varepsilon}}{\sigma}.\tag{60}$$

Check that  $\mathbf{Z} \sim N_n(\mathbf{0}_n, \mathbb{I}_n)$ , i.e.,  $\mathbf{Z}$  is a standard normal vector. Check that now that

$$\frac{\widehat{\boldsymbol{\varepsilon}}^T \widehat{\boldsymbol{\varepsilon}}}{\sigma^2} = \mathbf{Z}^T (\mathbb{I}_n - H) \mathbf{Z}.$$
(61)

Justify in detail. Aid: The formula in (59) may be useful here.

c) The unbiased estimate of variance  $\sigma^2$  is

$$\widehat{\sigma^2} = \frac{1}{n-k-1}\widehat{\boldsymbol{\varepsilon}}^T\widehat{\boldsymbol{\varepsilon}}$$

Show now that

$$(n-k-1)\frac{\hat{\sigma}^2}{\sigma^2} \sim \chi^2 (n-k-1).$$
 (62)

Aid: Check proposition 9.3 and the section 9.4 for rules on traces.

# 5 Properties of the Hat Matrix and Diagnostics of Regression Models

Let

$$X_{R} := \begin{pmatrix} x_{11} & \cdots & x_{1k} \\ x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots \\ x_{n1} & \cdots & x_{nk} \end{pmatrix}$$
(63)

so that the design matrix is partitioned as

$$X = \begin{pmatrix} \mathbf{1}_n & X_R \end{pmatrix}.$$

Note that

$$C_{ce}X_{R} = \left(X_{R} - \frac{1}{n}\mathbf{1}_{n}\mathbf{1}_{n}^{T}X_{R}\right) = \begin{pmatrix}x_{11} - \bar{x}_{1} & \cdots & x_{1k} - \bar{x}_{k}\\x_{21} - \bar{x}_{1} & \cdots & x_{2k} - \bar{x}_{k}\\\vdots & \vdots & \vdots\\x_{n1} - \bar{x}_{1} & \cdots & x_{nk} - \bar{x}_{k}\end{pmatrix},$$

where  $\bar{x}_{j} = \frac{1}{n} \sum_{i=1}^{n} x_{ij}, \ j = 1, ..., k$ . We set

$$X_{c} := \begin{pmatrix} x_{11} - \bar{x}_{1} & \cdots & x_{1k} - \bar{x}_{k} \\ x_{21} - \bar{x}_{1} & \cdots & x_{2k} - \bar{x}_{k} \\ \vdots & \vdots & \vdots \\ x_{n1} - \bar{x}_{1} & \cdots & x_{nk} - \bar{x}_{k} \end{pmatrix}$$
(64)

This is the matrix of centered regressor/covariate values.

The exam questions in this subsection are based on the material in the slides of Lecture 6. We provide a summary for convenience of reference.

Multiple regression in component form is:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i, \quad i = 1, \dots, n, \quad n > k+1.$$

The centered model is

$$Y_i = \alpha + \beta_1 \left( x_{i1} - \bar{x}_1 \right) + \dots + \beta_k \left( x_{ik} - \bar{x}_k \right) + \varepsilon_i, \quad i = 1, \dots, n$$
(65)

where  $\bar{x}_{j} = \frac{1}{n} \sum_{i=1}^{n} x_{ij}, \ j = 1, ..., k$ , and

$$\alpha = \beta_0 + \beta_1 \bar{x}_1 + \dots + \beta_k \bar{x}_k. \tag{66}$$

Now we can write the equations in (65) in matrix form using (64) as

$$\mathbf{Y} = (\mathbf{1}_n, X_c) \begin{pmatrix} \alpha \\ \boldsymbol{\beta}_R \end{pmatrix} + \boldsymbol{\varepsilon}.$$
 (67)

where

$$\boldsymbol{\beta}_{R} = \begin{pmatrix} \beta_{1} \\ \beta_{2} \\ \vdots \\ \beta_{k} \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_{1} \\ \varepsilon_{2} \\ \vdots \\ \varepsilon_{n} \end{pmatrix},$$

and  $X_c$  is given in (64).

It is shown in Lecture 6 that the LSE of the centered model are

$$\hat{\boldsymbol{\alpha}} = \bar{\boldsymbol{y}}$$

$$\hat{\boldsymbol{\beta}}_R = \left(\boldsymbol{X}_c^T \boldsymbol{X}_c\right)^{-1} \boldsymbol{X}_c^T \mathbf{y}$$
(68)

Next, advance from Lecture 6 to the exam questions.

First write (67) as

 $\mathbf{Y} = \mathbf{1}_n \alpha + X_c \boldsymbol{\beta}_R + \boldsymbol{\varepsilon}.$ 

Then the predictor (as a function of the training data) is

$$\hat{\mathbf{y}} = \mathbf{1}_n \hat{\alpha} + X_c \hat{\boldsymbol{\beta}}_R$$

Hence the corrresponding hat matrix is

$$H_c = X_c \left( X_c^T X_c \right)^{-1} X_c^T \tag{69}$$

## Problem 5.1.

In view of (68) we have

$$\hat{\mathbf{y}} = \mathbf{1}_n \bar{y} + X_c \hat{\boldsymbol{\beta}}_R$$

a) Do some simple steps of matrix calculus to derive the identity

$$\hat{\mathbf{y}} = \left(\frac{1}{n}\mathbf{1}_n\mathbf{1}_n^T + H_c\right)\mathbf{y}.$$
(70)

 $(Aid: \bar{y} = \frac{1}{n} \mathbf{1}_n^T \mathbf{y}).$ 

b) Show by (70) that the hat matrix of the non-centered model is

$$H = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T + X_c \left( X_c^T X_c \right)^{-1} X_c^T.$$
(71)

c) Show that the elements  $h_{ii}$  on the main diagonal of H are given as

$$h_{ii} = \frac{1}{n} + \left(\mathbf{x}_i - \bar{\mathbf{x}}\right)^T \left(X_c^T X_c\right)^{-1} \left(\mathbf{x}_i - \bar{\mathbf{x}}\right).$$
(72)

where

$$\mathbf{x}_i^T = (x_{i1}, x_{i2} \dots, x_{ik})$$

is the *i*th row in  $X_R$  and

$$\bar{\mathbf{x}}^T = (\bar{x}_1, \bar{x}_2 \dots, \bar{x}_k)$$

is the 1 × k vector of means of columns in  $X_R$ . *Aid*: The equation (72) follows from (71), when one takes  $A = X_c$  and  $B = (X_c^T X_c)^{-1}$  in (4).

d) Show that  $h_{ii} \ge \frac{1}{n}$ .

## Problem 5.2.

The hat matrix H plays an important part in the diagnostics of linear regression. This is due to the fact that it determines, as has been shown, the variances and covariances of  $\hat{\mathbf{Y}}$  and the residuals  $\hat{\boldsymbol{\varepsilon}}$  (recall (54) and (55).

a)

$$\widehat{\mathbf{y}} = \begin{pmatrix} \widehat{y}_1 \\ \vdots \\ \widehat{y}_n \end{pmatrix} = H \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}.$$

We have

 $\widehat{y}_i = \mathfrak{E}_i^T \widehat{\mathbf{y}}.$ 

Use (2) to write the right hand side of

$$\mathfrak{E}_i^T \widehat{\mathbf{y}} = \mathfrak{E}_i^T H \mathbf{y}$$

as a finite sum (display the sum) and compute by means of this sum the partial derivative

$$\frac{\partial \widehat{y}_i}{\partial y_i}$$

and interpret  $h_{ii}$  by this.

b) Show that (6.2)

$$h_{ii} = \mathbf{x}_i^T (X^T X)^{-1} \mathbf{x}$$

is the generic element in the main diagonal of the hat matrix H. Aid: The expression (4) may be useful.

c) In (72) we have shown that the elements  $h_{ii}$  on the main diagonal of H can also be written as

$$h_{ii} = \frac{1}{n} + \left(\mathbf{x}_i - \bar{\mathbf{x}}\right)^T \left(X_c^T X_c\right)^{-1} \left(\mathbf{x}_i - \bar{\mathbf{x}}\right).$$
(73)

In this expression the vector of means  $\bar{\mathbf{x}}$  can be regarded as a centroid. Check that for the simple centered linear regression model the elements of the hat matrix  $H_c$  are

$$h_{ij}(c) = \frac{(x_i - \bar{x})(x_j - \bar{x})}{S_{xx}}$$
 and  $h_{ii}(c) = \frac{(x_i - \bar{x})^2}{S_{xx}}$ 

Here you start with

$$X_c := \begin{pmatrix} x_1 - \bar{x} \\ x_2 - \bar{x} \\ \vdots \\ x_n - \bar{x} \end{pmatrix}, \tag{74}$$

where  $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$ . Check now using (73) that for the simple linear regression

$$\frac{1}{n} \le h_{ii} < 1$$

#### Problem 5.3.

This exam item deals with hidden extrapolation on pp. 107-110 of MPV. On page 110 we read the following:

read the following: prediction of estimation interses interpolation, while it this point new outside the RVH, extrapolation is required.

The diagonal elements  $h_{ii}$  of the hat matrix  $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$  are useful in detecting **hidden extrapolation**. The values of  $h_{ii}$  depend both on the Euclidean distance of the point  $\mathbf{x}_i$  from the centroid and on the density of the points in the RVH. In general, the point that has the largest value of  $h_{ii}$ , say  $h_{max}$ , will lie on the boundary of the RVH in a region of the x space where the density of the observations is relatively low. The set of points  $\mathbf{x}$  (not necessarily data points used to fit the model) that satisfy

$$\mathbf{x}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x} \le h_{\max}$$

is an ellipsoid enclosing all points inside the RVH (see Cook [1979] and Weisberg [1985]). Thus, if we are interested in prediction or estimation at the point  $\mathbf{x}'_0 = [1, x_{01}, x_{02}, \dots, x_{0k}]$ , the location of that point relative to the RVH is reflected by

$$h_{00} = \mathbf{x}_0' (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_0$$

a) We have used a training set to find the LSE  $\hat{\beta}$  regarded as *n* outcomes of the true model

$$\mathbf{Y} = X\boldsymbol{\beta}_* + \boldsymbol{\varepsilon}$$

We are predicting the value of the response Y at  $\mathbf{x}_o^T = (1, x_{o1} \dots, x_{ok})$  by

$$\hat{Y} = \mathbf{x}_0^T \hat{\boldsymbol{\beta}},$$

when

$$Y = \mathbf{x}_0^T \boldsymbol{\beta}_* + \varepsilon.$$

How does  $h_{\text{max}}$  in the quote of MPV above reflect the mean square error of the the random variable  $Y - \hat{Y}$ ?

b) Explain the problem of hidden extrapolation in predicting new responses and estimating the mean response at given point  $\mathbf{x}_0$  in the multiple linear regression. Justify your explanations by sketching the graph and explain how to detect this problem by using the properties of the hat matrix  $H = X(X^{\top}X)^{-1}X^{\top}$ .

### Problem 5.4.

Leverage is a technical term for the measurement of how far away a  $1 \times k$  regressor vector  $\mathbf{x}_i^T$  is from the rest of the regressor vectors in X, i.e., far away in  $\mathbb{R}^k$ . High-leverage points, if any, are outliers with respect to the RVH. This makes the predictor likely to pass close to a high leverage observation. Hence high-leverage points have the potential to cause large changes in the predictors if they are deleted: if this happens, then they are said to be **influential** points. The  $h_{ii}$ , i.e., the elements of the main diagonal of H, are used as the measures of leverage of the correspoding vector  $\mathbf{x}_i$ .

a) The distance related to  $h_{ii}$  reflects basically on the squared and sample based Mahalanobis distance. By the squared and sample based Mahalanobis distance between a vector  $\mathbf{x}_i$  in finite set of vectors and its centroid we mean the squared Mahalanobis distance in (13), when the sample covariance matrix S in (5) is used in place of  $\Sigma$ . We write

$$d_M^2\left(\mathbf{x}_i, \bar{\mathbf{x}}; S\right) = \left(\mathbf{x}_i - \bar{\mathbf{x}}\right)^T S^{-1} \left(\mathbf{x}_i - \bar{\mathbf{x}}\right).$$

Show that

$$(n-1)\left(h_{ii}-\frac{1}{n}\right) = d_M^2(\mathbf{x}_i, \bar{\mathbf{x}}; S)$$

b) The predicted residual error sum of squares (PRESS) is one of the techniques of crossvalidation used in linear regression analysis for detecting influential observations.

The PRESS statistic is

PRESS = 
$$\sum_{i=1}^{n} (y_i - \hat{y}_{(i)})^2$$
.

How is this obtained?

c) Use the formula (derived by Sherman-Morrison-Woodbury theorem, which you need not to recapitulate or invoke here)

$$\hat{y}_i = h_{ii}y_i + (1 - h_{ii})\hat{y}_{(i)},$$

where  $h_{ii}$  is the element on the main diagonal of H, to rewrite PRESS so that you only need to use predictors from full data linear regression. What important property of  $h_{ii}$ s is crucial here?

## Problem 5.5.

a) Consider

$$H_o: \beta_1 = \beta_2 = \ldots = \beta_k = 0$$

and the alternative hypothesis

$$H_1: \beta_i \neq 0$$
 for at least one  $\beta_i$ 

What null hypothesis about multiple linear regression are we testing here?

b) We assume normal multiple linear regression for the hypothesis testing. To test this we need  $H_c$  in the centered model. It has been shown in another problem in this exam generator that, if  $H_o$  is true,

$$\frac{\widehat{\varepsilon}^T \widehat{\varepsilon}}{\sigma^2} \sim \chi^2 \left( n - k - 1 \right). \tag{75}$$

This is valid even for the centered model. In the same manner we can show that for the model component of variance that

$$\frac{SS_{\rm R}}{\sigma^2} = \frac{\widehat{\mathbf{y}}^T C_{ce} \widehat{\mathbf{y}}}{\sigma^2} \sim \chi^2(k).$$
(76)

This is valid even for the centered model. It has been found that these two quadratic forms  $\hat{\sigma}^2$  and  $SS_{\rm R}$  are independent. It has been shown in another problem that

 $S_1 \sim \chi^2(d_1)$  and  $S_2 \sim \chi^2(d_2)$  are independent. Set

$$V := \frac{S_1/d_1}{S_2/d_2}.$$

What is the distribution of V?

c) Explain how the following ANOVA table for multiple linear regression is used to test the hypotheses above at a significance level  $\alpha$ . Remember to state explicitly the test statistic and its distribution, if  $H_o$  is true.

Source	df	Sum of Squares	MSS
Regression	k	$SS_{ m R}$	$SS_{ m R}/ m k$
Residual	n - k - 1	$SS_{ m Res}$	$SS_{\text{Res}}/(\text{n-k-1})$
Total	n-1	$SS_{\mathrm{T}}$	

Source = source of variation, df= degrees of freedom, SS= sum of squares, MSS= mean sum of squares.

## 6 Choice of Regression Models

### Problem 6.1.

Consider two models

$$\mathcal{M}_1$$
:  $E[\mathbf{Y}|X_1] = X_1 \boldsymbol{\beta}_1$  and  $\mathcal{M}_2$ :  $E[\mathbf{Y}|X_2] = X_2 \boldsymbol{\beta}_2$ ,

where  $X_1$  is  $n \times (k_1 + 1)$ ,  $X_2$  is  $n \times (k_2 + 1)$ ,  $\beta_1$  is  $(k_1 + 1) \times 1$ , and  $\beta_2$  is  $(k_2 + 1) \times 1$ .

We suppose that the models are nested so that  $\mathcal{M}_1 \subset \mathcal{M}_2$ : more explicitly,  $X_2 = [X_1 \ X_2]$ , where  $\tilde{X}_2$  is  $n \times (k_2 - k_1)$ . In other words, the first  $(k_1 + 1)$  columns of  $X_2$  are equal to those of  $X_1$ .

Denote by  $\widehat{\boldsymbol{\beta}}_1$  and  $\widehat{\boldsymbol{\beta}}_2$  the least squares estimates of  $\mathcal{M}_1$  and  $\mathcal{M}_2$ , respectively. Recall that the LSE errors are

$$Q_o^{(i)}(\widehat{\boldsymbol{\beta}}_i) = \|\mathbf{y} - X_i \widehat{\boldsymbol{\beta}}_i\|^2, \quad i = 1, 2.$$

The question is whether the larger, more flexible model  $\mathcal{M}_2$  is significantly better than the more restricted model  $\mathcal{M}_1$ . It has been shown that an *F*-statistic to decide this question is

$$F_{\mathcal{M}} = \frac{\left(Q_o^{(1)}\left(\widehat{\boldsymbol{\beta}}_1\right) - Q_o^{(2)}\left(\widehat{\boldsymbol{\beta}}_2\right)\right) / (k_2 - k_1)}{Q_o^{(2)}\left(\widehat{\boldsymbol{\beta}}_2\right) / (n - k_2)}.$$

In this problem we study the relationship between the *F*-test and likelihood ratio test of the the null hypothesis that the true model lies in  $\mathcal{M}_1$ . We have the likelihood functions

$$L_{\mathbf{y}}\left(\boldsymbol{\beta}_{i},\sigma^{2}\right) = \frac{1}{\sqrt{(\sigma^{2}2\pi)^{n}}} e^{-\frac{1}{2\sigma^{2}}\|\mathbf{y}-X_{i}\boldsymbol{\beta}_{i}\|^{2}}, \quad i = 1, 2.$$

The likelihood ratio (LR) statistic is

$$LR = \frac{\max_{\boldsymbol{\beta}_2, \sigma^2} L_{\mathbf{y}} \left(\boldsymbol{\beta}_2, \sigma^2\right)}{\max_{\boldsymbol{\beta}_1, \sigma^2} L_{\mathbf{y}} \left(\boldsymbol{\beta}_1, \sigma^2\right)}.$$
(77)

- a) Why does it always hold that  $LR \ge 1$ ?
- b) The MLEs of  $\sigma^2$  are (you do not need to show this here)

$$\widehat{\sigma}_{1}^{2} = \frac{\sum_{j=1}^{n} e_{1,j}^{2}}{n}, \widehat{\sigma}_{2}^{2} = \frac{\sum_{j=1}^{n} e_{2,j}^{2}}{n}$$

in the models in  $\mathcal{M}_1$  and  $\mathcal{M}_2$ , respectively. Here  $(e_{i,1}, \ldots, e_{i,n})$  are the least squares residuals w.r.t. model  $\mathcal{M}_i$ , i = 1, 2. Why does it always hold that

$$\widehat{\sigma}_1^2 \ge \widehat{\sigma}_2^2?$$

c) Show that

$$LR = \left(\frac{\widehat{\sigma}_1^2}{\widehat{\sigma}_2^2}\right)^{n/2}.$$

Aid: Use (77) and (58) in an approximate way.

- d) Show that  $F_{\mathcal{M}}$  is a monotone function of LR.
- e) The LR test rejects the null hypothesis that the true model lies in  $\mathcal{M}_1$  if  $\frac{\tilde{\sigma}_1^2}{\tilde{\sigma}_2^2} > c$  for a chosen constant c. The F-test rejects the null hypothesis that the true model lies in  $\mathcal{M}_1$ , if  $F_{\mathcal{M}} > F_{\alpha}(k_2 k_1, n k_2)$ , the  $\alpha$ -percentile. In view of the finding in d) above, what is the relationship between the F-test and the LR test for the choice of nested multiple regression models?

### Problem 6.2.

This problem deals with Mallows'  $C_p$  criterion for the choice of model dimension in multiple regression. The main statements below are recapitulated without proof in MPV on pp. 334-335. The symbol p in MPV is p = k + 1 in this document.

We are faced with multiple regression models  $\mathcal{M}_k$ , where k is the number of predictors in the model. The intercept is included in every model; hence  $\mathcal{M}_0$  is the model with only the intercept. M is the maximum number of possible predictors in a model. The true model for n observations of a response variable Y is OLS as in section 9.10 satisfying

$$\mathbf{Y} = X\boldsymbol{\beta}_* + \boldsymbol{\varepsilon},$$

where  $\boldsymbol{\beta}_*$  is a  $(k_*+1) \times 1$  vector and X is a  $n \times (k_*+1)$  matrix. The true number of predictors  $k_*$  is unknown. The purpose of developing Mallows'  $C_p$  is to estimate  $\hat{k}$  on the basis of n samples of Y so that overfitting is avoided.

Fix now any  $k \in \{1, \ldots, M\}$ . Then

$$\widehat{\mathbf{Y}} = \begin{pmatrix} \widehat{Y}_1 \\ \vdots \\ \widehat{Y}_n \end{pmatrix} = H_k \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \quad \widehat{\boldsymbol{\varepsilon}} = \begin{pmatrix} \widehat{\varepsilon}_1 \\ \vdots \\ \widehat{\varepsilon}_n \end{pmatrix} = \mathbf{Y} - H_k \mathbf{Y}$$

where  $H_k$  is the hat matrix for the model  $\mathcal{M}_k$ . Thus  $H_k$  is

$$H_k = X_k \left( X_k^T X_k \right)^{-1} X_k^T,$$

where  $X_k^T$  is a  $(k+1) \times n$  matrix. Hence  $X_k^T X$  is a well-defined  $(k+1) \times (k_*+1)$  matrix and  $H_k X \boldsymbol{\beta}_*$  is a well-defined  $n \times 1$  matrix.  $\hat{\boldsymbol{\varepsilon}}$  is the vector of LSE residuals computed in the model  $\mathcal{M}_k$ , i.e.  $\hat{\boldsymbol{\varepsilon}}(k)$ , but we use a simpler notation for ease of writing.

a) Verify that

$$E\left[\widehat{\mathbf{Y}}\right] = H_k X \boldsymbol{\beta}_*$$

Hence

$$E\left[\widehat{Y}_{i}\right] \neq E\left[Y_{i}\right].$$

and the decomposition (80) has a useful content.

b) Show that

$$C_{\hat{\mathbf{Y}}} = \sigma^2 H_k. \tag{78}$$

c) We have that

Show now that

$$E\left[\widehat{\boldsymbol{\varepsilon}}\right] = \left(\mathbb{I}_{n} - H_{k}\right) X \boldsymbol{\beta}_{*}.$$

$$C_{\widehat{\boldsymbol{\varepsilon}}} = \sigma^{2} \left(\mathbb{I}_{n} - H_{k}\right).$$
(79)

d) Check that

$$E\left[\left(\widehat{Y}_{i}-E\left[Y_{i}\right]\right)^{2}\right]=\left(E\left[\widehat{Y}_{i}\right]-E\left[Y_{i}\right]\right)^{2}+\operatorname{Var}\left[\widehat{Y}_{i}\right].$$
(80)

Aid: Start with  $E\left[\left(\widehat{Y}_i - E\left[Y_i\right]\right)^2\right] = E\left[\left(\widehat{Y}_i - E\left[\widehat{Y}_i\right] + E\left[\widehat{Y}_i\right] - E\left[Y_i\right]\right)^2\right]$ . Next, expand the sum inside the expectation, and then compute the expectation.

e) Show that

$$\sum_{i=1}^{n} \operatorname{Var}\left[\widehat{Y}_{i}\right] = (k+1)\sigma^{2}.$$

Aid: Note that the variances  $\operatorname{Var}\left[\widehat{Y}_{i}\right]$  are the diagonal elements of the covariance matrix of  $\widehat{\mathbf{Y}}$ , note (78), and recall Appendix C in the slides of Lecture 3.

f) Show that

$$E\left[\sum_{i=1}^{n}\widehat{\varepsilon}_{i}^{2}\right] = \sum_{i=1}^{n}\left(E\left[\widehat{Y}_{i}\right] - E\left[Y_{i}\right]\right)^{2} + (n - (k+1))\sigma^{2}.$$

Aid: Use a well known formula from the first course for  $E[\hat{\varepsilon}_i^2]$ . Then note (79) and use the results on slide 58/79 in Lecture 4 as in b) above.

g) Then show that

$$\Gamma_{k} := \frac{1}{\sigma^{2}} E\left[\left(\widehat{Y}_{i} - E\left[Y_{i}\right]\right)^{2}\right]$$
$$= \frac{1}{\sigma^{2}} \left[\sum_{i=1}^{n} \left(E\left[\widehat{Y}_{i}\right] - E\left[Y_{i}\right]\right)^{2} + \sum_{i=1}^{n} \operatorname{Var}\left[\widehat{Y}_{i}\right]\right]$$
$$= \frac{1}{\sigma^{2}} E\left[\sum_{i=1}^{n} \widehat{\varepsilon}_{i}^{2}\right] - n + 2(k+1).$$

Now Mallows replaces the mean with the sum of residual squares to define the (almost) final criterion of choice of model dimension as

$$C_k := \frac{1}{\sigma^2} \sum_{i=1}^n \widehat{\varepsilon}_i(k)^2 - n + 2(k+1)$$

where  $\widehat{\varepsilon}_i = \widehat{\varepsilon}_i(k)$  are the LSE residuals computed in the model  $\mathcal{M}_k, k \in \{1, \ldots, M\}$ . How should one estimate  $\sigma^2$ ?

## Problem 6.3.

a) Let X be  $n \times (k+1)$  and  $\boldsymbol{\beta}_*$  be  $(k+1) \times 1$  and

$$\mathbf{Y} = X\boldsymbol{\beta}_* + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N_n \left( \mathbf{0}, \sigma^2 \mathbb{I}_n, \right)$$
$$\mathbf{Z} = X\boldsymbol{\beta}_* + \boldsymbol{\nu}, \quad \boldsymbol{\nu} \sim N_n \left( \mathbf{0}, \sigma^2 \mathbb{I}_n \right).$$

In this  $\varepsilon$  and  $\nu$  are independent.

We use a training data set on the responses in Y

$$\mathcal{D}_{tr} = \left\{ \left( y_i, \mathbf{x}_i^T \right)_{i=1}^n \right\}$$

to find the LSE  $\hat{\boldsymbol{\beta}}$ . Next we wish to estimate **Z** by means of

$$\widehat{\mathbf{Z}} = X\widehat{\boldsymbol{\beta}} = H\mathbf{Y}.$$

Compute the mean square error per sample of  $\widehat{\mathbf{Z}}$  as

$$\frac{1}{n} \operatorname{E}\left[\left(\mathbf{Z} - \widehat{\mathbf{Z}}\right)^{T} \left(\mathbf{Z} - \widehat{\mathbf{Z}}\right)\right] = \sigma^{2} \left(1 + \frac{k+1}{n}\right).$$
(81)

Aid: Expand, use the properties of H, the rule (11) may be useful, too. Compute the involved traces.

b) The AIC (=Akaike Information Criterion) for model choice in a model family defined by the p.d.f.s  $f_{\mathbf{Y}}(y;\theta)$ , where  $p = \dim(\theta)$  is generally given by

$$AIC_p = -2\ln f_{\mathbf{Y}}\left(y;\widehat{\theta}_{MLE}\right) + 2 \cdot p.$$
(82)

The best model in the sense of AIC is found by

$$p_{\text{AIC}} = \operatorname{argmin}_{p} \left\{ -2 \ln f_{\mathbf{Y}} \left( y; \widehat{\theta}_{\text{MLE}} \right) + 2 \cdot p \right\}.$$

The assignment here is to apply AIC in multiple linear regression. Consider a family  $\mathbb{M}$  of K normal multiple regression models

$$\mathcal{M}_k: E\left[\mathbf{Y}|X_k\right] = X_k \boldsymbol{\beta}_k, \quad k = 1, 2, \dots, K,$$

where  $X_k$  is  $n \times (k+1)$ ,  $\boldsymbol{\beta}_i$  is  $(k+1) \times 1$ .

Show that AIC for model choice in  $\mathbb{M}$  given training sets  $\mathcal{D}_{tr} = \{(y_i)_{i=1}^n, X_k = \left(\mathbf{x}_{(1)}^T, \dots, \mathbf{x}_{(k)}^T\right)^T\},\$  $\mathbf{x}_{(l)}^T = \left(x_i^{(l)}\right)_{i=1}^n$  is

$$AIC_k = C_n + n \ln\left(\widehat{\sigma}_{MLE}^2\right) + 2 \cdot (k+1).$$
(83)

4.4 where  $C_n = n \ln(2\pi) + n$ . We do not include the variance  $\sigma^2$  in the number of parameters, as it is shared by all models.

*Aid:* You allowed to use, without working out or presenting the details, the relevant results in Problem 4.28 in section 4.4.

c) Some statisticians/ data analysts take (81) as a starting point and define an approximate AIC by

$$\operatorname{AIC}_{k}^{*} = \widehat{\sigma}_{\operatorname{MLE}}^{2} \left( 1 + \frac{k+1}{n} \right).$$
(84)

Describe concisely how  $AIC_k^*$  might incorporate the idea of parsimony, i.e., of optimal model fit and model complexity underlying the Akaike Information Criterion.

d) Give an approximate relation between  $\frac{\text{AIC}_k}{n}$  and a certain function of  $\text{AIC}_k^*$ . Aid: Approximation  $\ln(1+x) \approx x$ .

# 7 Generalized Linear Regression

In the problems of this section

$$\mathbf{x}^T \boldsymbol{eta} = eta_0 + \sum_{i=1}^p eta_i x_i$$

## Problem 7.1.

 $\epsilon \sim \text{Logistic}(0, 1)$ , iff its probability density function (pdf) is

$$\sigma'(x) = \frac{d}{dx}\sigma(x) = \frac{e^{-x}}{(1+e^{-x})^2}.$$

a) Verify the primitive function

$$\int \sigma'(t)dt = \frac{1}{1 + e^{-x}} + C$$
(85)

and determine C. Set

$$\sigma(x) := \frac{1}{1 + e^{-x}}$$

b) Show that

$$P(-\varepsilon \le x) = P(\varepsilon \le x).$$
(86)

You may use (85) even if you have failed to show it.

c) Define

$$Y^* = \mathbf{x}^T \boldsymbol{\beta} + \varepsilon \tag{87}$$

and construct

$$Y = \begin{cases} 1 & \text{if } Y^* > 0\\ 0 & \text{otherwise.} \end{cases}$$

Verify that

$$P(Y = 1 | \mathbf{x}) = \sigma (\mathbf{x}^T \boldsymbol{\beta}).$$

You may use (86) even if you have failed to show it.

d) The r.v. Y in this assignment is a special case of generalized linear regression (GLM). The way to construct a GLM goes via a link function. Present and justify the link function in this assignment.

### Problem 7.2.

Consider the normal multiple least squares model,

$$Y^* = \mathbf{x}^T \boldsymbol{\beta} + \varepsilon, \tag{88}$$

where  $\mathbf{x}^T = (1, x_1, x_2, \dots, x_p), \ \boldsymbol{\beta}^T = (\beta_0, \beta_1, \dots, \beta_p)$  and  $\varepsilon \sim N(0, 1)$ , and is independent of  $\mathbf{x}$ .

a) Explain why we have

$$P\left(-\varepsilon \le x\right) = P\left(\varepsilon \le x\right). \tag{89}$$

Aid: Let  $\Phi(x)$  be the cumulative distribution function of N(0,1) and  $\phi(x)$  be the corresponding pdf. Since  $\phi(x) = \phi(-x)$ , we have

$$\Phi(-x) = 1 - \Phi(x).$$

b) Construct

$$Y = \begin{cases} 1 & \text{if } Y^* > 0, \\ 0 & \text{otherwise.} \end{cases}$$

Verify that

$$P(Y=1 \mid \mathbf{x}) = \Phi\left(\mathbf{x}^T \boldsymbol{\beta}\right),$$

You may use (89) even if you have failed to show it.

c) The r.v. Y in this assignment is a special case of generalized linear regression (GLM). The way to construct a GLM goes via a link function. Present and justify the link function in this assignment.

## Problem 7.3.

Y is an exponentially distributed random variable with parameter  $\lambda > 0$ , we write  $Y \sim \text{Exp}(\lambda)$ , if its probability density function is

$$f(y;\lambda) = \begin{cases} \lambda e^{-\lambda y} & y \ge 0, \\ 0 & y < 0. \end{cases}$$

 $\lambda$  is sometimes called the rate parameter. We want to construct a generalized linear model (exponential regression). We know that  $E[Y] = 1/\lambda$ , this need not be derived here.

- a) Find the canonical link function relating  $\mathbf{x}^T \boldsymbol{\beta}$  to the mean. What does this require from  $\mathbf{x}^T \boldsymbol{\beta}$ ?
- b) Write down the probability density function

$$f(y; \mathbf{x}^T \boldsymbol{\beta})$$

c) What is  $E[Y|\mathbf{x}]$  in exponential regression?

## 8 The Woodbury Matrix Identity & Ridge Regression

## Problem 8.1.

Woodbury matrix identity is a generalization of the Sherman-Morrison-Woodbury theorem, and says that

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U \left(C^{-1} + VA^{-1}U\right)^{-1} VA^{-1},$$
(90)

where A, U, C and V are conformable matrices and the required inverses exist. You are not asked to prove this. By (37) we have the ridge regression estimate

$$\widehat{\boldsymbol{\beta}}_{\lambda} = \left(X^T X + \lambda \mathbb{I}_p\right)^{-1} X^T \mathbf{y}.$$
(91)

Now, apply the Woodbury matrix identity with  $A = X^T X$ , X has full column rank p = k + 1,  $\lambda \mathbb{I}_p = C$ ,  $V = U = \mathbb{I}_p$  to express  $\widehat{\beta}_{\lambda}$  as in the right hand side of (90). Or, we find that

$$\widehat{\boldsymbol{\beta}}_{\lambda} = \left(X^T X\right)^{-1} X^T \mathbf{y} - \text{CORR}$$
(92)

- a) What is the expression for CORR?
- b) Let next  $\lambda \to 0$ . By (91), if X has full column rank,

$$\widehat{\boldsymbol{\beta}}_{\lambda} \to \left( X^T X \right)^{-1} X^T \mathbf{y}.$$

Find the limit for  $\widehat{\boldsymbol{\beta}}_{\lambda}$  by means of (92), when  $\lambda \to 0$ . *Hint*: Use (90) once more in CORR.

# 9 Collections of Formulas and Auxiliary Results

## 9.1 Matrices and Matrix Rules

- A and B conformal,  $(AB)^T = B^T A^T$ . A and B invertible,  $(AB)^{-1} = B^{-1} A^{-1} \cdot (A^T)^T = A$ .
- A and B conformal,  $(A + B)^T = A^T + B^T$
- A is  $n \times n$  and invertible.

$$(A^T)^{-1} = (A^{-1})^T.$$
 (93)

Proof:  $A^T (A^{-1})^T = (A^{-1}A)^T = \mathbb{I}_n^T = \mathbb{I}_n$  and  $(A^{-1})^T A^T = (AA^{-1})^T = \mathbb{I}_n^T = \mathbb{I}_n$ . Hence, as  $X^T X$  is  $k \times k$ , and symmetric

$$((X^T X)^{-1})^T = (X^T X)^{-1}$$
(94)

The  $n \times n$  identity matrix is

$$\mathbb{I}_{n} = \begin{pmatrix}
1 & 0 & 0 & \dots & 0 \\
0 & 1 & 0 & \dots & 0 \\
0 & \ddots & \vdots & \dots & 0 \\
0 & 0 & 0 & \dots & 1
\end{pmatrix}.$$
(95)

The  $n \times 1$  vector of ones is denoted by

$$\mathbf{1}_{n} = \begin{pmatrix} 1\\1\\\vdots\\1 \end{pmatrix} \in \mathbb{R}^{n}.$$
(96)

The centering matrix is

$$C_{ce} := \mathbb{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T.$$
(97)

## 9.2 The Range Space (=Column Space) of a Matrix

Let A be any  $n \times p$  matrix. Then the range space of A is defined by

$$\mathcal{R}(A) = \{ \mathbf{x} \in \mathbb{R}^n | \text{there exists a } \mathbf{b} \in \mathbb{R}^p \text{ such that } \mathbf{x} = A\mathbf{b} \}.$$
(98)

Then we have the following statement.

**Proposition 9.1.** X is an  $n \times (k+1)$  matrix with full column rank = k + 1. Then

$$\mathcal{R}(X^T) = \mathcal{R}(X^T X).$$
(99)

*Proof:* We show only half of the assertion, namely that

$$\mathcal{R}\left(X^{T}X\right) \subset \mathcal{R}\left(X^{T}\right).$$
(100)

For that purpose take any  $\mathbf{w} \in \mathcal{R}(X^T X)$ . Thus there exists a  $\mathbf{b} \in \mathbb{R}^{k+1}$  such that

$$\mathbf{w} = X^T X \mathbf{b} = X^T \left( X \mathbf{b} \right) = X^T \mathbf{c},$$

where  $\mathbf{c} = X\mathbf{b} \in \mathbb{R}^n$ . Therefore  $\mathbf{w} \in \mathcal{R}(X^T)$  by definition of  $\mathcal{R}(X^T)$ . Hence we have established (100).

## 9.3 **Projection Matrix**

Any  $n \times n$  matrix P from  $\mathbb{R}^n$  to a subspace (=range space of P),  $\mathcal{R}(P) \subset \mathbb{R}^n$ , is a called a *projection matrix*, if it is idempotent and symmetric, i.e.,

$$P^2 = P, \quad P^T = P.$$

## 9.4 Trace of a Square Matrix

Let A be a square matrix. The **trace** Tr A of A is the sum of the entries in main diagonal:

$$\operatorname{Tr}\begin{pmatrix}a_{11} & \cdots & a_{1k}\\ \vdots & \ddots & \vdots\\ a_{k1} & \cdots & a_{kk}\end{pmatrix} = \sum_{1}^{k} a_{jj}$$

- 1. If A is a  $k \times n$ -matrix, and B an  $n \times k$ -matrix, then Tr(AB) = Tr(BA)
- 2. In particular, if a is a column-vector, then  $a^T a = \text{Tr}(aa^T)$ .
- 3. For any real numbers a and b, Tr(aC + bD) = a Tr C + b Tr D

**Exempel 9.2.** The hat matrix H is  $H = X(X^TX)^{-1}X^T$ . Set  $B = X^T$  and  $A = X(X^TX)^{-1}$ . Then

$$\operatorname{Tr} H = \operatorname{Tr} X (X^T X)^{-1} X^T = \operatorname{Tr} AB = \operatorname{Tr} BA,$$

where we used rule 2.. But  $BA = X^T X (X^T X)^{-1} = \mathbb{I}_{k+1}$  (How is k+1 there?). Thus

$$\operatorname{Tr} H = \operatorname{Tr} \mathbb{I}_{k+1} = k+1.$$

### 9.5 Factorization and Square Root of Covariance Matrices

If  $\Sigma$  is an  $n \times n$  symmetric matrix, then  $\Sigma$  can be written as

$$\Sigma = ADA^T,$$

where A is an orthogonal matrix  $(A^T A = A A^T = I_n)$  and D is an  $n \times n$  diagonal matrix, with the eigenvalues on the main diagonal. If  $\Sigma$  is a covariance matrix, its eigenvalues  $\lambda_i$  are non-negative. Then

$$\Sigma^{1/2} = A D^{1/2} A^T.$$

where  $D^{1/2}$  is an  $n \times n$  diagonal matrix, with  $\sqrt{\lambda_i}$  on the main diagonal. One checks now that  $\Sigma^{1/2}\Sigma^{1/2} = \Sigma$ .

When  $\Sigma$  is positive definite, its eigenvalues are positive and we define

$$\Sigma^{-1/2} = A D^{-1/2} A^T, \tag{101}$$

where  $D^{-1/2}$  is an  $n \times n$  diagonal matrix, with  $1/\sqrt{\lambda_i}$  on the main diagonal.  $\Sigma^{-1/2}$  is symmetric, since  $D^{-1/2}$  is symmetric. Clearly,  $D^{-1/2}D^{-1/2} = D^{-1}$ . Then

$$\Sigma^{-1} = \Sigma^{-1/2} \Sigma^{-1/2}.$$
 (102)

## 9.6 Linear Transformations of Covariance Matrices

$$E\left[\mathbf{X} + \mathbf{Y}\right] = \boldsymbol{\mu}_{\mathbf{X}} + \boldsymbol{\mu}_{\mathbf{Y}} \tag{103}$$

If  $\mathbf{Z} = A\mathbf{X} + \mathbf{b}$ , then

$$E\left[\mathbf{Z}\right] = A\boldsymbol{\mu}_{\mathbf{X}} + \mathbf{b},\tag{104}$$

and

$$C_{\mathbf{Z}} = A C_{\mathbf{X}} A^T. \tag{105}$$

$$C_{\mathbf{X}} = E\left[\mathbf{X}\mathbf{X}^{T}\right] - \boldsymbol{\mu}_{\mathbf{X}}\boldsymbol{\mu}_{\mathbf{X}}^{T}$$
(106)

$$\operatorname{Var}\left[\mathbf{a}^{T}\mathbf{X}\right] = \mathbf{a}^{T}C_{\mathbf{X}}\mathbf{a} \tag{107}$$

## 9.7 Generalized Inverse

If a generalized inverse G of A satisfies the four conditions below, then G is called the Bjerhammar -Moore-Penrose (BMP) inverse.

- MP1 AGA = A
- MP2 GAG = G
- MP3  $(AG)^T = AG$
- MP4  $(GA)^T = GA$

## 9.8 Rules of Computation with Finite Sums

(1) 
$$\sum_{i=1}^{n} x_i = x_1 + x_2 + \ldots + x_n$$
.

(2) 
$$\sum_{i=1}^{n} a \cdot x_i = a \sum_{i=1}^{n} x_i.$$

(3) 
$$\underline{\sum_{i=1}^{n} (x_i + y_i)} = \sum_{i=1}^{n} x_i + \sum_{i=1}^{n} y_i.$$

(4) 
$$\underline{\sum_{i=1}^{n} (ax_i + by_i)} = a \sum_{i=1}^{n} x_i + b \sum_{i=1}^{n} y_i$$

- (5)  $\underline{\sum_{i=1}^{n} (x_i + y_i)^2}_{\overline{x} := \frac{1}{n} \sum_{i=1}^{n} x_i^2} + 2 \sum_{i=1}^{n} x_i y_i + \sum_{i=1}^{n} y_i^2.$
- (6)  $\sum_{i=1}^{n} (x_i \overline{x}) = 0.$

(7) 
$$\underline{\sum_{i=1}^{n} (x_i - \overline{x}) \cdot (y_i - \overline{y})} = \underline{\sum_{i=1}^{n} (x_i - \overline{x}) y_i} = \underline{\sum_{i=1}^{n} x_i (y_i - \overline{y})}.$$

(8) 
$$\underline{\sum_{i=1}^{n} (x_i - \overline{x}) \cdot (y_i - \overline{y})} = \sum_{i=1}^{n} x_i y_i - n \overline{x} \overline{y}_i$$

(9) 
$$\underline{\sum_{i=1}^{n} (x_i - \overline{x})^2} = \underline{\sum_{i=1}^{n} x_i^2} - n\overline{x}^2.$$

## 9.9 Simple Linear Regression

## 9.9.1 The Model

Let

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots \\ 1 & x_n \end{pmatrix}, \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$$

be an  $(n \times 1)$ -vector,  $(n \times 2)$ -matrix, and  $(2 \times 1)$ -vector, respectively. The simple linear regression is now given by

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

where  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$  is an  $n \times 1$  random vector. We assume the following:

- 1)  $E[\boldsymbol{\varepsilon}] = \mathbf{0}_n$  (= the  $n \times 1$  zero vector), i.e.,  $\mathbf{0}_n \in \mathbb{R}^n$ .
- 2) The errors are uncorrelated: the covariance matrix  $C_{\boldsymbol{\varepsilon}}$  of  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$  is

$$C_{\boldsymbol{\varepsilon}} = \sigma^2 \mathbb{I}_n$$

3)  $\sigma^2$  does not depend on X.

9.9.2 LSE

$$Q(\beta_{0},\beta_{1}) = \sum_{i=1}^{n} (y_{i} - \beta_{0} - \beta_{1}x_{i})^{2}$$

$$\begin{cases} \frac{\partial}{\partial\beta_{0}}Q(\hat{\beta}_{0},\hat{\beta}_{1}) = 0\\ \frac{\partial}{\partial\beta_{1}}Q(\hat{\beta}_{0},\hat{\beta}_{1}) = 0. \end{cases}$$

$$\frac{\hat{\beta}_{0} = \overline{y} - \hat{\beta}_{1}\overline{x}.$$

$$(108)$$

$$\frac{\partial}{\partial\beta_{1}}Q(\hat{\beta}_{0},\hat{\beta}_{1}) = 0 \Leftrightarrow \sum_{i=1}^{n} \left(y_{i} - \hat{\beta}_{0} - \hat{\beta}_{1}x_{i}\right)x_{i} = 0.$$

$$\hat{\beta}_{1} = \frac{\sum_{i=1}^{n} \left(y_{i} - \overline{y}\right)x_{i}}{\sum_{i=1}^{n} \left(x_{i} - \overline{x}\right)x_{i}}.$$

$$(109)$$

Let  $\beta_0^*, \beta_1^*$  denote the true poulation values of the intercept and the regression coefficient. Assume that

$$\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T \sim N_n \left( \mathbf{0}_n, \sigma^2 \mathbb{I}_n \right)$$

Then

$$\hat{\beta}_0 \sim N\left(\beta_0^*, \sigma^2\left(\frac{1}{n} + \frac{\overline{x}^2}{\sum_{i=1}^n (x_i - \overline{x})^2}\right)\right)$$

$$\hat{\beta}_1 \sim N\left(\beta_1^*, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \overline{x})^2}\right).$$
(110)

## 9.10 Multiple Linear Regression

The Ordinary Least Squares Model (OLS)

We have the model

$$\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon},\tag{111}$$

where

$$X = \begin{pmatrix} \mathbf{x}_1^{\mathsf{T}} \\ \mathbf{x}_2^{\mathsf{T}} \\ \vdots \\ \mathbf{x}_n^{\mathsf{T}} \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

and we assume the following:

- 1)  $E[\boldsymbol{\varepsilon}] = \mathbf{0}_n$  (= the  $n \times 1$  zero vector), i.e.,  $\mathbf{0}_n \in \mathbb{R}^n$ .
- 2) The covariance matrix  $C_{\boldsymbol{\varepsilon}}$  of  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$  is

$$C_{\varepsilon} = \sigma^2 \mathbb{I}_n. \tag{112}$$

- 3)  $\sigma^2$  does not depend on X.
- 4) n > k + 1 and X has full rank.

## 9.11 Normal Linear Regression

OLS with  

$$\boldsymbol{\varepsilon} \sim N_n \left( \mathbf{0}_n, \sigma^2 \mathbb{I}_n \right)$$

$$\mathbf{Y} = X \boldsymbol{\beta}_* + \boldsymbol{\varepsilon}.$$
Let  $\mathbf{X} \sim N_n \left( \boldsymbol{\mu}, C \right)$  and  $\mathbf{Y} = B \mathbf{X} + \mathbf{b}, B$  is  $m \times n$ . Then  

$$\mathbf{Y} \sim N_m \left( B \boldsymbol{\mu} + \mathbf{b}, B C B^T \right).$$
(113)

## 9.12 Distribution of Quadratic Forms of Normal Vectors

1.  $\mathbf{X} \sim N_n(\mu, \Sigma)$  is an  $n \times 1$  Gaussian vector, where  $\Sigma$  is positive definite. The quadratic form is

$$\left(\mathbf{X}-\mu\right)^{T}\Sigma^{-1}\left(\mathbf{X}-\mu\right).$$

We use the factorization of  $\Sigma^{-1}$  in (102) to get

$$\left(\mathbf{X}-\boldsymbol{\mu}\right)^{T}\boldsymbol{\Sigma}^{-1}\left(\mathbf{X}-\boldsymbol{\mu}\right) = \left(\boldsymbol{\Sigma}^{-1/2}\left(\mathbf{X}-\boldsymbol{\mu}\right)\right)^{T}\boldsymbol{\Sigma}^{-1/2}\left(\mathbf{X}-\boldsymbol{\mu}\right).$$

Let  $\mathbf{Z} := \Sigma^{-1/2} (\mathbf{X} - \mu)$ . Then our rules of computation give that Z has the mean vector  $E[Z] = \mathbf{0}_n$  and Z has the covariance matrix

$$C_{\mathbf{Z}} = \Sigma^{-1/2} \Sigma \Sigma^{-1/2} = \Sigma^{-1/2} \Sigma^{1/2} \Sigma^{1/2} \Sigma^{-1/2} = \mathbf{I}_n.$$

 $\mathbf{Z} \sim N_n(\mathbf{0}_n, \mathbb{I}_n)$ . Hence

$$\left(\mathbf{X}-\boldsymbol{\mu}\right)^T \boldsymbol{\Sigma}^{-1} \left(\mathbf{X}-\boldsymbol{\mu}\right) = \mathbf{Z}^T \mathbf{Z} = \sum_{i=1}^n Z_i^2 \sim \chi^2(n).$$

This is Thm 9.1. in Gut, Allan: An Intermediate Course in Probability. Second Edition. Note that  $z_i \sim N(0, 1)$  are independent and that a finite sum of squares of independent standard normal r.v.'s has  $\chi^2(n)$  (chi-squared distribution with n degrees of freedom).

2.

**Proposition 9.3.** If  $\mathbf{Z} \sim N_n(\mathbf{0}_n, \mathbb{I}_n)$ , then

$$\mathbf{Z}^T A \mathbf{Z} \sim \chi^2(r) \tag{114}$$

if and only if A is an idempotent matrix with  $\operatorname{rank} A = r$ .

## 9.13 Sherman-Morrison-Woodbury Theorem

Suppose A is is an invertible square  $n \times n$  matrix and  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$  are column vectors. Then  $A + \mathbf{u}\mathbf{v}^{\mathsf{T}}$  is invertible iff  $1 + \mathbf{v}^{\mathsf{T}}A^{-1}\mathbf{u} \neq 0$ . In this case,

$$(A + \mathbf{u}\mathbf{v}^{\mathsf{T}})^{-1} = A^{-1} - \frac{A^{-1}\mathbf{u}\mathbf{v}^{\mathsf{T}}A^{-1}}{1 + \mathbf{v}^{\mathsf{T}}A^{-1}\mathbf{u}}.$$
(115)

## 9.14 Fundamental Analysis of Variance Identity

$$\sum_{i=1}^{n} (y_i - \bar{y})^2 = \sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$
(116)

is known as the Fundamental Analysis of Variance Identity

$$\sum_{i=1}^{n} (y_i - \bar{y})^2 \leftrightarrow SS_{\mathrm{T}}$$

 $\sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2 \leftrightarrow SS_{\rm R} \text{ regression or model sum of squares}$  $\sum_{i=1}^{n} (y_i - \hat{y}_i)^2 \leftrightarrow SS_{\rm Res} \text{ Residual Sum of Squares}$  $SS_{\rm T} = SS_{\rm R} + SS_{\rm Res}$ 

$$R^2 = \frac{SS_{\rm R}}{SS_{\rm T}} = 1 - \frac{SS_{\rm Res}}{SS_{\rm T}}.$$
(118)

(117)

## 9.15 Matrix Derivatives

Let **A** be a  $k \times k$  matrix of constants, **a** be a  $k \times 1$  vector of constants, and **y** be a  $k \times 1$  vector of variables.

- 1. If  $\mathbf{z} = \mathbf{a}^{\top} \mathbf{y}$ , then  $\frac{\partial \mathbf{z}}{\partial \mathbf{y}} = \frac{\partial \mathbf{a}^{\top} \mathbf{y}}{\partial \mathbf{y}} = \mathbf{a}.$ 2. If  $\mathbf{z} = \mathbf{y}^{\top} \mathbf{y}$ , then  $\frac{\partial \mathbf{z}}{\partial \mathbf{y}} = \frac{\partial \mathbf{y}^{\top} \mathbf{y}}{\partial \mathbf{y}} = 2\mathbf{y}.$
- 3. If  $\mathbf{z} = \mathbf{a}^{\top} \mathbf{A} \mathbf{y}$ , then

$$rac{\partial \mathbf{z}}{\partial \mathbf{y}} = rac{\partial \mathbf{a}^{\top} \mathbf{A} \mathbf{y}}{\partial \mathbf{y}} = \mathbf{A}^{\top} \mathbf{a}.$$

4. If  $\mathbf{z} = \mathbf{y}^{\top} \mathbf{A} \mathbf{y}$ , then

$$\frac{\partial \mathbf{z}}{\partial \mathbf{y}} = \frac{\partial \mathbf{y}^\top \mathbf{A} \mathbf{y}}{\partial \mathbf{y}} = \mathbf{A} \mathbf{y} + \mathbf{A}^\top \mathbf{y}.$$

If **A** is symmetric, then

$$\frac{\partial \mathbf{y}^{\top} \mathbf{A} \mathbf{y}}{\partial \mathbf{y}} = 2\mathbf{A} \mathbf{y}$$

# 9.16 Solution of (20) (for all interested)

We expand using rules of transpose and multiplication

$$(\mathbf{Y} - A\mathbf{X}) (\mathbf{Y} - A\mathbf{X})^T = (\mathbf{Y} - A\mathbf{X}) (\mathbf{Y}^T - (A\mathbf{X})^T)$$
$$= \mathbf{Y}\mathbf{Y}^T - \mathbf{Y}(A\mathbf{X})^T - A\mathbf{X}\mathbf{Y}^T + A\mathbf{X}(A\mathbf{X})^T.$$

Now we take the expectations

$$E \left[ \mathbf{Y}\mathbf{Y}^{T} - \mathbf{Y}(A\mathbf{X})^{T} - A\mathbf{X}\mathbf{Y}^{T} + A\mathbf{X}(A\mathbf{X})^{T} \right] = E \left[ \mathbf{Y}\mathbf{Y}^{T} \right] - E \left[ \mathbf{Y}(A\mathbf{X})^{T} \right]$$

$$(119)$$

$$-E \left[ A\mathbf{X}\mathbf{Y}^{T} \right] + E \left[ A\mathbf{X}(A\mathbf{X})^{T} \right].$$

Since mean vectors are here zero vectors, we use next (11), i.e.,  $Cov(A\mathbf{X}, B\mathbf{Y}) = ACov(\mathbf{X}, \mathbf{Y})B^T$  with

$$E\left[\mathbf{Y}(A\mathbf{X})^{T}\right] = \operatorname{Cov}(\mathbf{Y}, A\mathbf{X}) = \operatorname{Cov}(\mathbf{Y}, \mathbf{X})A^{T} = E\left[(\mathbf{Y}\mathbf{X}^{T}\right]A^{T},$$
(120)

and

$$E\left[A\mathbf{X}\mathbf{Y}^{T}\right] = \operatorname{Cov}(A\mathbf{X}, \mathbf{Y}) = A\operatorname{Cov}(\mathbf{X}, \mathbf{Y}) = AE\left[\mathbf{X}\mathbf{Y}^{T}\right], \qquad (121)$$

and

$$E\left[A\mathbf{X}(A\mathbf{X})^{T}\right] = \operatorname{Cov}(A\mathbf{X}, A\mathbf{X}) = A\operatorname{Cov}(\mathbf{X}, \mathbf{X})A^{T} = A\operatorname{C}_{\mathbf{X}}A^{T}.$$
(122)

By definition

$$E\left[\mathbf{Y}\mathbf{Y}^{T}\right] = \mathbf{C}_{\mathbf{Y}}.$$
(123)

When we insert from (120) - (123) in the right hand side of (119) we get

$$E\left[ (\mathbf{Y} - A\mathbf{X}) (\mathbf{Y} - A\mathbf{X})^T \right] = C_{\mathbf{Y}}$$
$$-E\left[ \mathbf{Y}\mathbf{X}^T \right] A^T - AE\left[ \mathbf{X}\mathbf{Y}^T \right]$$
$$+AC_{\mathbf{X}}A^T.$$
(124)

Here the expression (19) of the optimal matrix A gives

$$E \begin{bmatrix} \mathbf{Y}\mathbf{X}^T \end{bmatrix} A^T = \operatorname{Cov}(\mathbf{Y}, \mathbf{X}) A^T = \operatorname{Cov}(\mathbf{Y}, \mathbf{X}) (\operatorname{C}_{\mathbf{X}}^{-1})^T \operatorname{Cov}(\mathbf{Y}, \mathbf{X})^T$$
  
=  $\operatorname{Cov}(\mathbf{Y}, \mathbf{X}) \operatorname{C}_{\mathbf{X}}^{-1} \operatorname{Cov}(\mathbf{X}, \mathbf{Y}).$  (125)

Here we first invoked (9) by  $\operatorname{Cov}(\mathbf{Y}, \mathbf{X})^T = (\operatorname{Cov}(\mathbf{X}, \mathbf{Y})^T)^T = \operatorname{Cov}(\mathbf{X}, \mathbf{Y})$ . Second, we used (93), as well as the fact that the covariance matrix  $C_{\mathbf{X}}$  is symmetric. Again by (19),

$$AE\left[\mathbf{X}\mathbf{Y}^{T}\right] = \operatorname{Cov}(\mathbf{Y}, \mathbf{X}) \operatorname{C}_{\mathbf{X}}^{-1} \operatorname{Cov}(\mathbf{X}, \mathbf{Y}).$$
(126)

Finally by (19)

$$AC_{\mathbf{X}}A^{T} = Cov(\mathbf{Y}, \mathbf{X}) \underbrace{C_{\mathbf{X}}^{-1}C_{\mathbf{X}}}_{=\mathbb{I}_{n}} \left( Cov(\mathbf{Y}, \mathbf{X})C_{\mathbf{X}}^{-1} \right)^{T}$$
$$= Cov(\mathbf{Y}, \mathbf{X}) \left(C_{\mathbf{X}}^{-1}\right)^{T} Cov(\mathbf{Y}, \mathbf{X})^{T}$$

By (9),  $\operatorname{Cov}(\mathbf{Y}, \mathbf{X})^T = (\operatorname{Cov}(\mathbf{X}, \mathbf{Y})^T)^T = \operatorname{Cov}(\mathbf{X}, \mathbf{Y})$ . We note also again (93) and that the covariance matrix  $C_{\mathbf{X}}$  is symmetric. Hence we have found that

$$AC_{\mathbf{X}}A^{T} = Cov(\mathbf{Y}, \mathbf{X})C_{\mathbf{X}}^{-1}Cov(\mathbf{X}, \mathbf{Y}).$$
(127)

Next we substitute (125), (126) and (127) in the right hand side of (124) we obtain

$$E\left[\left(\mathbf{Y} - A\mathbf{X}\right)\left(\mathbf{Y} - A\mathbf{X}\right)^{T}\right] = C_{\mathbf{Y}} -Cov(\mathbf{Y}, \mathbf{X})C_{\mathbf{X}}^{-1}Cov(\mathbf{X}, \mathbf{Y}) - Cov(\mathbf{Y}, \mathbf{X})C_{\mathbf{X}}^{-1}Cov(\mathbf{X}, \mathbf{Y}) +Cov(\mathbf{Y}, \mathbf{X})C_{\mathbf{X}}^{-1}Cov(\mathbf{X}, \mathbf{Y}).$$
  
$$= C_{\mathbf{Y}} - Cov(\mathbf{Y}, \mathbf{X})C_{\mathbf{X}}^{-1}Cov(\mathbf{X}, \mathbf{Y}).$$
(128)

Clearly, the rightmost expression in (128) verifies (20).