



--

Semester 1 Assessment, 2019

School of Mathematics and Statistics

**MAST30025 Linear Statistical Models**

Writing time: 3 hours

Reading time: 15 minutes

This is NOT an open book exam

This paper consists of 9 pages (including this page)

**Authorised materials:**

- The only permitted scientific calculator is the Casio FX82.
- Two A4 double-sided handwritten sheets of notes.

**Instructions to Students**

- You must NOT remove this question paper at the conclusion of the examination.
- You should attempt all questions. Marks for individual questions are shown.
- The total number of marks available is 90.

**Instructions to Invigilators**

- Students must NOT remove this question paper at the conclusion of the examination.

This paper may be held in the Baillieu Library

**This paper must not be removed from the examination room**

**Question 1 (10 marks)**

- (a) [3 marks] Let  $A_1, \dots, A_m$  be a set of symmetric idempotent matrices whose sum is also idempotent. Show directly that

$$r\left(\sum_{i=1}^m A_i\right) = \sum_{i=1}^m r(A_i).$$

- (b) [3 marks] Let  $\mathbf{y}$  be a random vector with  $\text{var } \mathbf{y} = V$ . Show that  $V$  is positive semidefinite.
- (c) [4 marks] Show directly that for any matrix  $A$ , we have  $A = A(A^T A)^c A^T A$ . (*Hint: If  $M^T M = 0$ , then  $M = 0$ .*)

**Question 2 (12 marks)** Let

$$\mathbf{y} \sim MVN\left(\begin{bmatrix} 4 \\ -1 \end{bmatrix}, \begin{bmatrix} 3 & 1 \\ 1 & 3 \end{bmatrix}\right), \quad A = \frac{1}{2} \begin{bmatrix} c & 1-3c \\ 1-3c & c \end{bmatrix}.$$

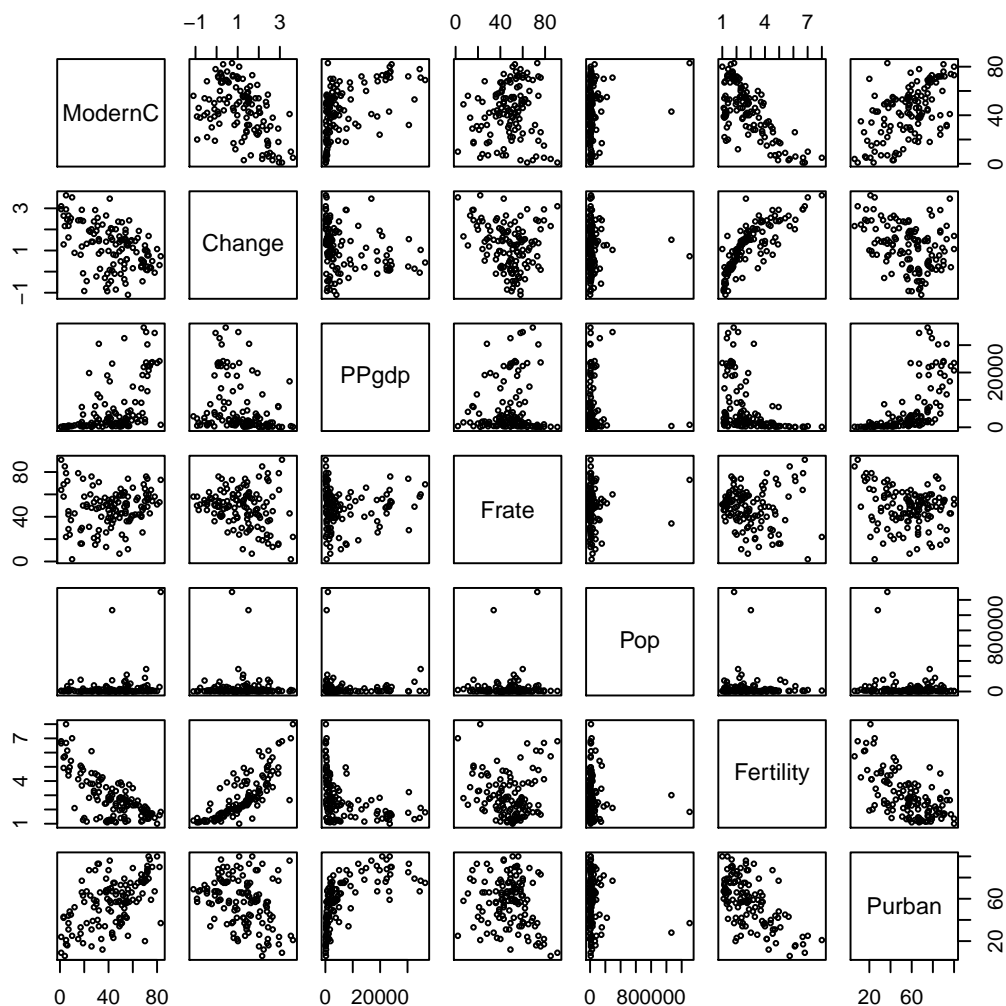
- (a) [3 marks] Let  $c = 1$ . Calculate  $E[\mathbf{y}^T A \mathbf{y}]$ .
- (b) [3 marks] Let  $c = -1$ . Describe the distribution of  $A \mathbf{y}$ .
- (c) [4 marks] Find all values of  $c$  for which  $\mathbf{y}^T A \mathbf{y}$  has a non-central  $\chi^2$  distribution.
- (d) [2 marks] Let  $c = \frac{1}{2}$ . Determine if  $\mathbf{y}^T A \mathbf{y}$  is independent of  $y_1 + y_2$ .

**Question 3 (18 marks)** In this question, we study a dataset of 125 countries, collected by the United Nations. This dataset contains the variables:

- **ModernC**: Percent of unmarried women using a modern method of contraception
- **Change**: Annual population growth rate, percent
- **PPgdp**: Per capita gross national product, US dollars
- **Frate**: Percent of females over age 15 economically active
- **Pop**: Total 2001 population, 1000s
- **Fertility**: Expected number of live births per female, 2000
- **Purban**: Percent of population that is urban, 2001

We wish to model the birth rate (**Fertility**) in terms of the other variables. The following R calculations are produced (with some removed):

```
> UN <- read.csv('UN3.csv', header=T)
> pairs(UN, cex=0.5)
```



```
> UN$Fertility <- log(UN$Fertility)
> UN$PPgdp <- log(UN$PPgdp)
> UN$Pop <- log(UN$Pop)
> fullmodel <- lm(Fertility ~ ., data=UN)
> deviance(fullmodel)
```

```
[1] 4.929815
```

```
> # Input removed
```

```
Start: AIC=-390.13
```

```
Fertility ~ ModernC + Change + PPgdp + Frate + Pop + Purban
```

	Df	Sum of Sq	RSS	AIC
- Frate	1	0.0003	4.9301	-392.12
- PPgdp	1	0.0277	4.9575	-391.43
<none>			4.9298	-390.13
- Pop	1	0.3084	5.2382	-384.54
- ModernC	1	0.3461	5.2759	-383.65
- Purban	1	0.5826	5.5124	-378.16
- Change	1	10.2407	15.1705	-251.62

```
Step: AIC=-392.12
```

```
Fertility ~ ModernC + Change + PPgdp + Pop + Purban
```

	Df	Sum of Sq	RSS	AIC
- PPgdp	1	0.0285	4.9586	-393.40
<none>			4.9301	-392.12
+ Frate	1	0.0003	4.9298	-390.13
- Pop	1	0.3102	5.2403	-386.49
- ModernC	1	0.3559	5.2860	-385.41
- Purban	1	0.6135	5.5436	-379.46
- Change	1	10.9186	15.8487	-248.15

```
Step: AIC=-393.4
```

```
Fertility ~ ModernC + Change + Pop + Purban
```

```
> summary(model)
```

```
Call:
```

```
lm(formula = Fertility ~ ModernC + Change + Pop + Purban, data = UN)
```

```
Residuals:
```

Min	1Q	Median	3Q	Max
-0.4866	-0.1282	0.0084	0.1321	0.5862

```
Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	1.285452	0.102858	12.497	< 2e-16 ***

```

ModernC      -0.003971    0.001121   -3.541 0.000567 ***
Change       0.323302    0.019679   16.429 < 2e-16 ***
Pop          -0.024369    0.009324   -2.613 0.010110 *
Purban       -0.006111    0.000985   -6.205 8.09e-09 ***

```

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

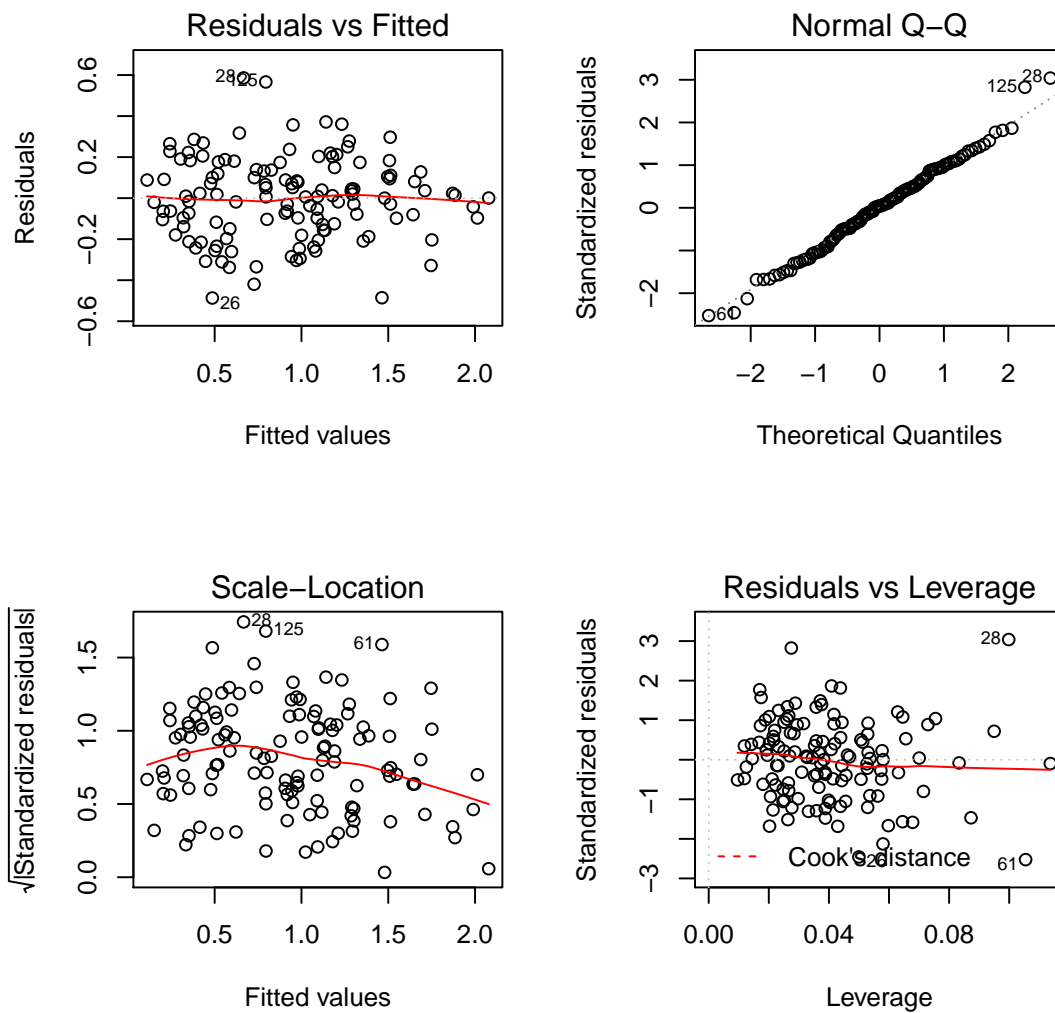
```
Residual standard error: 0.2033 on 120 degrees of freedom
```

```
Multiple R-squared:  0.8462,      Adjusted R-squared:  0.841
```

```
F-statistic: 165 on 4 and 120 DF,  p-value: < 2.2e-16
```

```
> par(mfrow=c(2,2))
```

```
> plot(model)
```



```
> library(car)
> linearHypothesis(model, c(0,0,0,log(2),0), log(0.98))
```

Linear hypothesis test

Hypothesis:

0.693147180559945 Pop = - 0.0202027073175195

Model 1: restricted model

Model 2: Fertility ~ ModernC + Change + Pop + Purban

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	121	4.9694				
2	120	4.9586	1	0.010847	0.2625	0.6093

```
> qt(0.975,118:120)
```

```
[1] 1.980272 1.980100 1.979930
```

```
> qf(0.975,1,118:120)
```

```
[1] 5.154550 5.153431 5.152331
```

```
> qf(0.975,2,118:120)
```

```
[1] 3.806642 3.805631 3.804638
```

- [2 marks]** Give two reasons in favour of the logarithmic transformation of the **PPgdp** variable.
- [4 marks]** Test the hypothesis  $H_0 : \beta_{\text{Pop}} = 0$  for the full model (with all variables included) at the 5% significance level. Clearly state your  $F$ -statistic and critical value, and interpret your conclusion in the context of the study.
- [2 marks]** Identify the variable selection procedure that has been used.
- [2 marks]** Perform one step of backwards elimination on **model**.
- [2 marks]** Write down the final fitted model, including any variable transformations used.
- [2 marks]** Calculate a 95% confidence interval for the parameter corresponding to the **Change** variable.
- [2 marks]** Interpret the results of the **linearHypothesis** function in the context of the study.
- [2 marks]** From the diagnostic plots, comment on the suitability of the linear model for this data.

**Question 4 (13 marks)** Consider the full rank linear model  $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  with  $p$  parameters.

- (a) [3 marks] Show that  $\mathbf{b}$ , the least squares estimator of  $\boldsymbol{\beta}$ , and  $s^2$ , the sample variance, are independent.
- (b) [3 marks] Calculate the expected value of the maximum likelihood estimator of  $\sigma^2$  and thereby show that it is biased.
- (c) [4 marks] Given a (symmetric) variance matrix  $\text{var } \mathbf{y} = V$ , the method of generalised least squares estimates the parameters  $\boldsymbol{\beta}$  by minimising  $\mathbf{e}^T V^{-1} \mathbf{e}$ . Derive an expression for the resulting estimators.
- (d) [3 marks] Calculate the variance of the estimators derived in question 4c.

**Question 5 (12 marks)** Consider the general linear model,  $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ . This model may be of full or less than full rank.

- (a) [2 marks] State two methods that can be used to fit this model to data, and compare them.
- (b) [2 marks] Is a high leverage or low leverage more desirable? Explain.
- (c) [2 marks] What are two possible reasons to transform a response variable?
- (d) [2 marks] Define overfitting and explain how it may be avoided.
- (e) [2 marks] Define interaction between two continuous predictors and explain how to model it.
- (f) [2 marks] Under what circumstances should one prefer a complete block design to a completely randomised design, and vice versa?

**Question 6 (14 marks)** The nursing director at a private hospital wishes to compare the weekly number of complaints received against the nursing staff during three daily shifts: first (7am–3pm), second (3pm–11pm), and third (11pm–7am). Her plan is to sample 17 weeks and select a shift at random from each week sampled, recording the number of complaints received during the selected shift.

The following data is collected:

	number of complaints		
	observations	mean	sample variance
shift 1	9, 9, 11, 9, 12	10	2
shift 2	8, 11, 6, 8, 9, 12	9	4.8
shift 3	15, 14, 10, 11, 10, 12	12	4.4

The data is analysed using a one-way classification model:

$$y_{ij} = \mu + \tau_i + \varepsilon_{ij}, \quad i = 1, 2, 3, j = 1, \dots, n_i.$$

You are also given the following R calculations:

```
> qt(0.975, 14:17)
```

```
[1] 2.144787 2.131450 2.119905 2.109816
```

```
> qf(0.975, 1, 14:17)
```

```
[1] 6.297939 6.199501 6.115127 6.042013
```

```
> qf(0.975, 2, 14:17)
```

```
[1] 4.856698 4.765048 4.686665 4.618874
```

- [2 marks]** Write down the normal equations for this model.
- [2 marks]** Write down two distinct estimates for the parameters of the model. One of these should estimate  $\hat{\mu} = 10$ .
- [3 marks]** The sample variance  $s^2$  is calculated to be 3.857. Calculate a 95% prediction interval for the weekly number of complaints generated in the first shift.
- [3 marks]** The nursing director wishes to test if the third shift generates as many complaints as the average of the other two. Express this as a formal hypothesis and show that it is testable.
- [4 marks]** Test this hypothesis at the 5% significance level. Clearly state your  $F$ -statistic and critical value.



**Question 7 (11 marks)**

- (a) **[4 marks]** Sleepzeze is a new drug that is claimed to reduce insomnia. You plan to conduct an experiment to test this claim. You are given resources to monitor a total of 120 person-nights (120 people for 1 night each, or 60 people for two nights each, etc.), and wish to perform an experiment with a complete block design, with gender as a blocking factor.
- Briefly describe how control, blocking, randomisation, and blinding may be used in a design of this experiment.
- (b) **[2 marks]** Write down a design matrix and parameter vector for a model which has a complete block design with 2 treatments and 3 blocks, and one sample per treatment/block combination.
- (c) **[2 marks]** Write down estimators for the model from question 7b for the treatment parameters only, as a linear function of the response vector  $\mathbf{y}$ .
- (d) **[3 marks]** For the model from question 7b, directly calculate the variance of the estimator of the difference between the two treatment parameters. You may write it in terms of the error variance  $\sigma^2$ .

**End of Exam—Total Available Marks = 90.**