Semester 1 Assessment, 2016

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 8 pages (including this page)

**Authorised materials**:

- Scientific calculators are premitted, but not graphical calculators.

- One A4 double-sided handwritten sheet 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 (9 marks)**

(a) Let $A$ be a square matrix and suppose that $A^k = A^{k+1}$ for some $k \geq 1$. Show that $A$ is idempotent.

(b) Let $X$ be an $n \times p$ matrix of full rank, where $n > p$. Show that $H = X(X^TX)^{-1}X^T$ is idempotent, and find its rank. (You may assume that $H$ is symmetric.)

(c) Show that if a square matrix $A$ is positive semidefinite, then its eigenvalues are non-negative.

**Question 2 (10 marks)** Let

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \sim MVN \left( \begin{bmatrix} a \\ -a \\ 0 \end{bmatrix}, \begin{bmatrix} 2 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix} \right),$$

where $a$ is a constant.

(a) What is the distribution of $y_1 + y_2$?

(b) What is the distribution of $\frac{1}{2} \left( y_1^2 - 2y_1y_2 + y_2^2 + y_3^2 \right)$?

(c) Suppose $a = 0$. For what values of $c$ does

$$c \frac{y_1^2 - 2y_1y_2 + y_2^2 + y_3^2}{y_1^2 + 2y_1y_2 + y_2^2}$$

have an $F$ distribution?

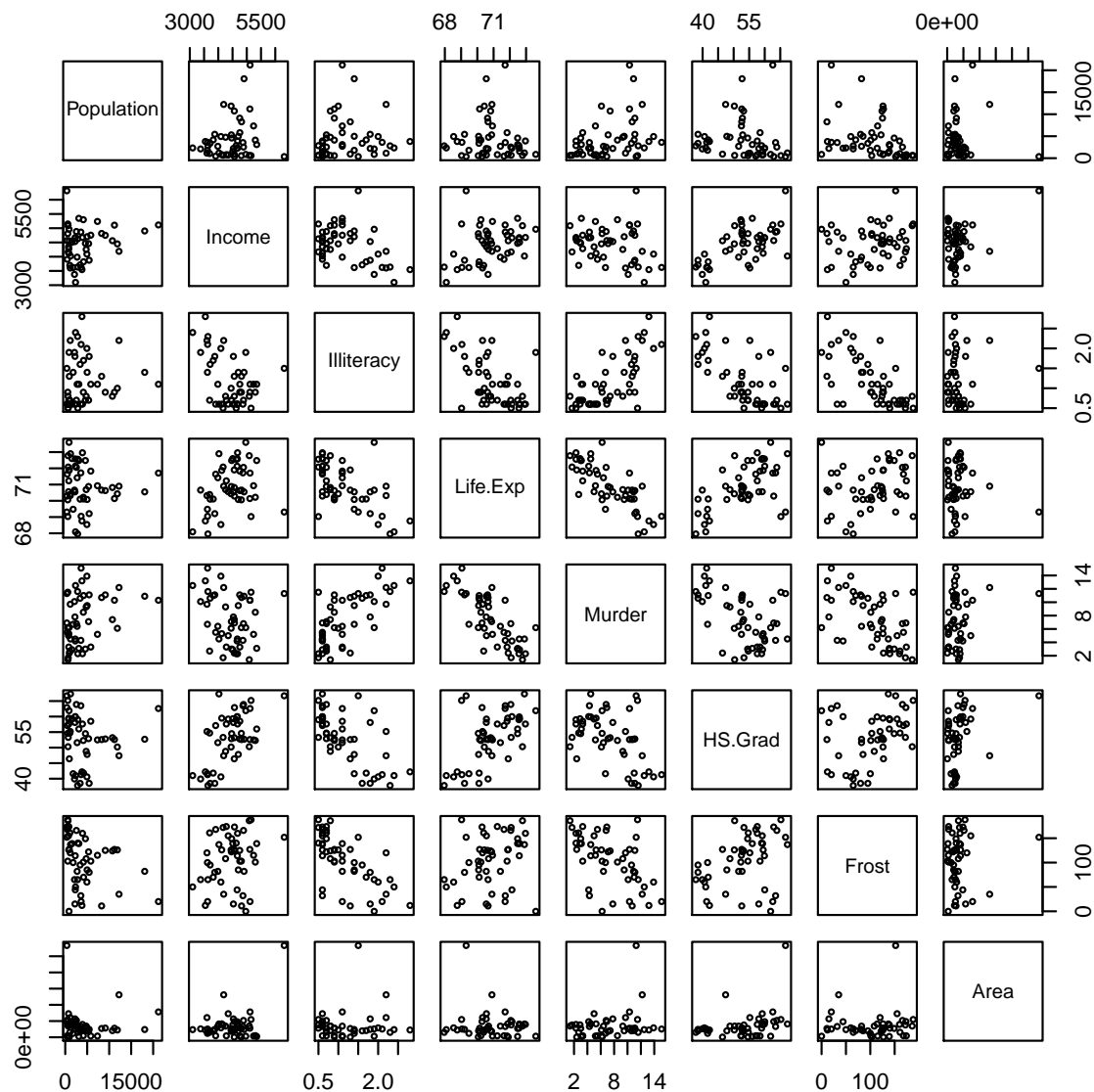**Question 3 (14 marks)** Consider the full rank linear model, $\mathbf{y} = X\beta + \varepsilon$.

(a) State the assumptions involved in fitting this model.

(b) Define the term BLUE (best linear unbiased estimator).

(c) Is it better to fit this model using the method of least squares or maximum likelihood estimation? Justify your answer.

(d) Define and explain the purpose of the leverage of a point.

(e) Explain the difference between a model relevance test and a model relevance test using a corrected sum of squares.

(f) When is a model with fewer explanatory variables more desirable or less desirable than a model with more explanatory variables?

(g) Explain why the residual sum of squares $SS_{Res}$ is not an appropriate goodness-of-fit measure for model selection.

**Question 4 (17 marks)** In this question, we study a dataset of 50 US states. This dataset contains the variables:

- `Population`: population estimate as of July 1, 1975

- `Income`: per capita income (1974)

- `Illiteracy`: illiteracy (1970, percent of population)

- `Life.Exp`: life expectancy in years (1969–71)

- `Murder`: murder and non-negligent manslaughter rate per 100,000 population (1976)

- `HS.Grad`: percentage of high-school graduates (1970)

- `Frost`: mean number of days with minimum temperature below freezing (1931–1960) in capital or large city

- `Area`: land area in square miles

We use linear models to model life expectancy in terms of the other variables. The following R output is produced.

```
> data(state)
> statedata <- data.frame(state.x77, row.names=state.abb, check.names=TRUE)
> pairs(statedata, cex=0.5)
```

```
> statedata$Population <- log(statedata$Population)
> statedata$Area <- log(statedata$Area)
> nullmodel <- lm(Life.Exp ~ 1, data = statedata)
> fullmodel <- lm(Life.Exp ~ ., data = statedata)
> model <- step(fullmodel, scope = ~ .)

Start:  AIC=-23.6
Life.Exp ~ Population + Income + Illiteracy + Murder + HS.Grad +
    Frost + Area

              Df Sum of Sq     RSS       AIC
- Income       1    0.0018  22.650  -25.5934
- Illiteracy   1    0.0556  22.704  -25.4746
```

```
- Area        1    0.2106 22.859 -25.1344
<none>                     22.648 -23.5973
- Frost       1    1.2374 23.886 -22.9374
- Population  1    1.8854 24.533 -21.5992
- HS.Grad     1    2.4375 25.086 -20.4864
- Murder      1   23.2760 45.924   9.7483

Step:  AIC=-25.59
Life.Exp ~ Population + Illiteracy + Murder + HS.Grad + Frost +
    Area


              Df Sum of Sq    RSS      AIC
- Illiteracy  1    0.0556 22.705 -27.4708
- Area        1    0.2197 22.870 -27.1107
<none>                     22.650 -25.5934
- Frost       1    1.2602 23.910 -24.8862
+ Income      1    0.0018 22.648 -23.5973
- Population  1    2.1909 24.841 -22.9768
- HS.Grad     1    4.0374 26.687 -19.3918
- Murder      1   24.2130 46.863   8.7601

Step:  AIC=-27.47
Life.Exp ~ Population + Murder + HS.Grad + Frost + Area


              Df Sum of Sq    RSS      AIC
- Area        1    0.2157 22.921 -28.998
<none>                     22.705 -27.471
+ Illiteracy  1    0.0556 22.650 -25.593
+ Income      1    0.0017 22.704 -25.475
- Population  1    2.2792 24.985 -24.688
- Frost       1    2.3760 25.082 -24.495
- HS.Grad     1    4.9491 27.655 -19.612
- Murder      1   29.2296 51.935  11.899

Step:  AIC=-29
Life.Exp ~ Population + Murder + HS.Grad + Frost


              Df Sum of Sq    RSS      AIC
<none>                     22.921 -28.998
+ Area        1    0.216 22.705 -27.471
+ Illiteracy  1    0.052 22.870 -27.111
+ Income      1    0.011 22.911 -27.021
- Frost       1    2.214 25.135 -26.387
- Population  1    2.450 25.372 -25.920
- HS.Grad     1    6.959 29.881 -17.741
- Murder      1   34.109 57.031  14.578
```

```
> summary(model)

Call:
lm(formula = Life.Exp ~ Population + Murder + HS.Grad + Frost,
    data = statedata)

Residuals:
     Min      1Q   Median      3Q      Max
-1.41760 -0.43880  0.02539  0.52066  1.63048

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 68.720810   1.416828  48.503  < 2e-16 ***
Population   0.246836   0.112539   2.193 0.033491 *
Murder      -0.290016   0.035440  -8.183 1.87e-10 ***
HS.Grad      0.054550   0.014758   3.696 0.000591 ***
Frost       -0.005174   0.002482  -2.085 0.042779 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.7137 on 45 degrees of freedom
Multiple R-squared:  0.7404,        Adjusted R-squared:  0.7173
F-statistic: 32.09 on 4 and 45 DF,  p-value: 1.17e-12

> anova(nullmodel, model, fullmodel)

Analysis of Variance Table

Model 1: Life.Exp ~ 1
Model 2: Life.Exp ~ Population + Murder + HS.Grad + Frost
Model 3: Life.Exp ~ Population + Income + Illiteracy + Murder + HS.Grad +
    Frost + Area
  Res.Df    RSS Df Sum of Sq       F    Pr(>F)
1     49 88.299
2     45 22.921  4    65.378 30.3101 6.901e-12 ***
3     42 22.648  3     0.273  0.1688    0.9168
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

> signif(vcov(model), 6)

             (Intercept)    Population        Murder       HS.Grad         Frost
(Intercept)   2.00740000 -1.18811e-01 -1.98357e-02 -1.44506e-02 -1.42795e-03
Population    -0.11881100  1.26650e-02 -3.56651e-04  2.36109e-04  8.91432e-05
Murder        -0.01983570 -3.56651e-04  1.25601e-03  1.84375e-04  3.42863e-05
HS.Grad       -0.01445060  2.36109e-04  1.84375e-04  2.17798e-04 -3.18945e-06
Frost         -0.00142795  8.91432e-05  3.42863e-05 -3.18945e-06  6.15931e-06
```

(a) Why do we take a logarithmic transformation of population and area?

(b) Find the Akaike's Information Criterion for the model with variables `Population`, `Murder`, `Frost`, and `Area`.

(c) Write down the final fitted model (including any variable transforms used).

(d) Calculate the sample variance $s^2$ for the final model.

(e) Calculate a 95% confidence interval for $\beta_{Population} - \beta_{Murder}$. (The 97.5% critical value for a $t$ distribution with 45 d.f. is 2.014.)

(f) What conclusions do you draw from the tests in the ANOVA table?

(g) If you were to perform an $F$ test of $H_0 : \beta_{Frost} = 0$ in the final model, what would your $F$ statistic and $p$-value be?

(h) Explain the $F$-statistic for the final model (last line of the `summary` call). Why is it different to the $F$-value in line 2 of the ANOVA table?

**Question 5 (14 marks)** Consider the general linear model $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, which may be of full or less than full rank.

(a) Define the term estimable.

(b) Show that if $\mathbf{t}^T = \mathbf{t}^T (X^T X)^c X^T X$, then $\mathbf{t}^T \boldsymbol{\beta}$ is estimable.

(c) Show that in a one-factor model, all treatment contrasts are estimable.

(d) If $\mathbf{t}^T \boldsymbol{\beta}$ is estimable, derive the distribution of $\mathbf{t}^T \mathbf{b}$, where $\mathbf{b}$ is the least squares estimator of $\boldsymbol{\beta}$.

(e) If $\mathbf{t}^T \boldsymbol{\beta}$ is estimable, show that $\mathbf{t}^T \mathbf{b}$ is independent of the sample variance $s^2$.

**Question 6 (12 marks)** The nursing director at a private hospital wishes to compare the weekly number of complaints received against the nursing staff during the 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 | |
|---|---|---|---|
| | number of observations | mean | sample variance |
| shift 1 | 5 | 10 | 2 |
| shift 2 | 6 | 9 | 4.8 |
| shift 3 | 6 | 12 | 4.4 |

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

(a) What kind of experimental design is this?

(b) Calculate the sample variance $s^2$ for the linear model.

(c) Calculate a 95% prediction interval for the total number of complaints received in a day. (The 97.5% critical value of a $t$ distribution with 14 d.f. is 2.145.) (*Hint:* You will need to modify the formula for a prediction interval.)

(d) Test the hypothesis that shift has no effect on the number of complaints. (The 95% critical value of an $F$ distribution with 2 and 14 d.f. is 3.739.)

**Question 7 (14 marks)**

(a) Discuss when it is best to use a completely randomised design, complete block design, or Latin square design.

(b) For a complete block design, why do we fit an additive model and not an interaction model?

(c) Write down a design matrix and parameter vector for a balanced incomplete block design for a model with 3 treatments and 3 blocks, each of size 2.

(d) Calculate the reduced design matrix $X_{2|1}$ for this model.

(e) Do you expect the reduced normal equations for this model to have the same solution as the normal equations for a completely randomised design of 6 experimental units over 3 treatments?

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

**Author/s:**

Mathematics and Statistics

**Title:**

Linear Statistical Models, 2016 Semester 1, MAST30025

**Date:**

2016

**Persistent Link:**

http://hdl.handle.net/11343/127833