



Semester 1 Assessment, 2017

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 permitted, but not graphical calculators.
- 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 (9 marks)

- (a) Show that if A is a symmetric and idempotent matrix, then $r(A) = \text{tr}(A)$.
- (b) Show, without using the above result, that if X is an $n \times p$ matrix with $n > p$, then $r(X(X^T X)^c X^T) = r(X)$.
- (c) Let $B\mathbf{x} = \mathbf{g}$ be a consistent linear system. Show that $\mathbf{x} = B^c\mathbf{g}$ is a solution to this system.

Question 2 (13 marks) Let

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \sim MVN \left(\begin{bmatrix} 0 \\ -1 \end{bmatrix}, \begin{bmatrix} 2 & 1 \\ 1 & 3 \end{bmatrix} \right).$$

You are given the following R calculations.

```
> V <- matrix(c(2,1,1,3),2,2)
> P <- eigen(V)$vectors
> P%*%diag(sqrt(eigen(V)$values))%*%t(P)

      [,1]      [,2]
[1,] 1.3763819 0.3249197
[2,] 0.3249197 1.7013016

> P%*%diag(1/sqrt(eigen(V)$values))%*%t(P)

      [,1]      [,2]
[1,] 0.7608452 -0.1453085
[2,] -0.1453085  0.6155367
```

- (a) Find two independent standard normal random variables which are linear combinations of \mathbf{y} elements (and constants).
- (b) Calculate $E[y_1^2 - 2y_1y_2]$.
- (c) Find all quadratic forms in \mathbf{y} which have a non-central χ^2 distribution with 2 degrees of freedom.
- (d) Show that $4y_1^2 - 4y_1y_2 + y_2^2$ is independent of $y_1^2 + 6y_1y_2 + 9y_2^2$.

Question 3 (14 marks) The following data is a sample of 5 random countries. For each country we measure the following variables:

- **logPPgdp**: The logarithm of the 2001 gross domestic product per person in US dollars;
- **logFertility**: The logarithm of the birth rate per 1000 females in the year 2000;
- **Purban**: The percentage of the population which lives in an urban area.

Name	logFertility	logPPgdp	Purban
Moldova	0.23	5.8	41
Netherlands	0.38	10.1	90
Estonia	0.14	8.3	69
Uganda	1.36	5.5	15
Hungary	0.13	8.6	65

We wish to predict **logFertility** using **logPPgdp** and **Purban**, using a linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ with no constant (intercept) term.

You are given the following R calculations.

```
> qt(0.975,1:5)
[1] 12.706205 4.302653 3.182446 2.776445 2.570582

> qf(0.95,1,1:5)
[1] 161.447639 18.512821 10.127964 7.708647 6.607891

> qf(0.95,2,1:5)
[1] 199.500000 19.000000 9.552094 6.944272 5.786135
```

- Calculate the least squares estimates of $\boldsymbol{\beta}$.
- Calculate and interpret a 95% confidence interval for the parameter associated with **Purban**. You are given the sample variance $s^2 = 0.0887$.
- Test for model relevance at a 5% significance level.
- Croatia has a gross domestic product of \$4500 per person, and 58% of its population lives in an urban area. Calculate a 95% prediction interval for its birth rate.

Question 4 (13 marks) Consider a full rank linear model $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$. We wish to derive the formula for a prediction interval for the sum of the responses of two independent future observations with predictors \mathbf{x}_1 and \mathbf{x}_2 respectively. This uses the unbiased point estimator $(\mathbf{x}_1 + \mathbf{x}_2)^T \mathbf{b}$, where \mathbf{b} is the least squares estimator of $\boldsymbol{\beta}$.

- (a) Calculate the variance of the prediction error of this estimator.
- (b) Show that this estimator is normally distributed.
- (c) Show that this estimator is independent of SS_{Res} .
- (d) Derive a t -distributed quantity based on this estimator and state its degrees of freedom.
- (e) Thus write down a formula for a $100(1 - \alpha)\%$ prediction interval for the sum of two responses.

Question 5 (14 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) State two methods that can be used to fit this model to data and compare them.
- (b) Explain the difference between a confidence interval and a prediction interval.
- (c) Define Akaike's information criterion and explain why it is useful as a goodness-of-fit measure for model selection.
- (d) Give an advantage and a disadvantage of stepwise selection over forward selection.
- (e) Define interaction between two categorical predictors.
- (f) List two principles of experimental design.
- (g) Explain what a Latin square is and its use in experimental design.

Question 6 (15 marks) Data was collected on the world record times for the one-mile run. For males, the records are from the period 1861–2003, and for females, from the period 1967–2003. This data is analysed below.

```
> mile <- read.csv('mile.csv', header=T)
> mile$Gender <- factor(mile$Gender)
> plot(Time ~ Year, data = mile, pch=as.character(Gender))
> imodel <- lm(Time ~ (Year + Gender)^2, data = mile)
> summary(imodel)
```

Call:

```
lm(formula = Time ~ (Year + Gender)^2, data = mile)
```

Residuals:

Min	1Q	Median	3Q	Max
-5.4512	-1.6160	-0.1137	1.1784	13.7265

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2309.4247	202.0583	11.429	< 2e-16 ***
Year	-1.0337	0.1021	-10.126	1.95e-14 ***
GenderMale	-1355.6778	203.1441	-6.673	1.03e-08 ***
Year:GenderMale	0.6675	0.1027	6.502	2.00e-08 ***

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 2.989 on 58 degrees of freedom

Multiple R-squared: 0.9663, Adjusted R-squared: 0.9645

F-statistic: 553.8 on 3 and 58 DF, p-value: < 2.2e-16

```
> amodel <- lm(Time ~ Year + Gender, data = mile)
> summary(amodel)
```

Call:

```
lm(formula = Time ~ Year + Gender, data = mile)
```

Residuals:

Min	1Q	Median	3Q	Max
-5.9071	-2.0988	-0.1141	1.2002	13.1863

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1003.00334	27.84691	36.02	<2e-16 ***
Year	-0.37364	0.01406	-26.57	<2e-16 ***
GenderMale	-34.85078	1.30099	-26.79	<2e-16 ***

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 3.896 on 59 degrees of freedom

```
Multiple R-squared:  0.9417,           Adjusted R-squared:  0.9397
F-statistic: 476.3 on 2 and 59 DF,  p-value: < 2.2e-16
```

```
> anova(amodel, imodel)
```

Analysis of Variance Table

```
Model 1: Time ~ Year + Gender
Model 2: Time ~ (Year + Gender)^2
Res.Df   RSS Df Sum of Sq    F    Pr(>F)
1      59 895.62
2      58 518.03  1     377.59 42.276 2.0001e-08 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
> linearHypothesis(imodel, c(0,1,0,1), -0.3)
```

Linear hypothesis test

Hypothesis:
Year + Year:GenderMale = - 0.3

```
Model 1: restricted model
Model 2: Time ~ (Year + Gender)^2
```

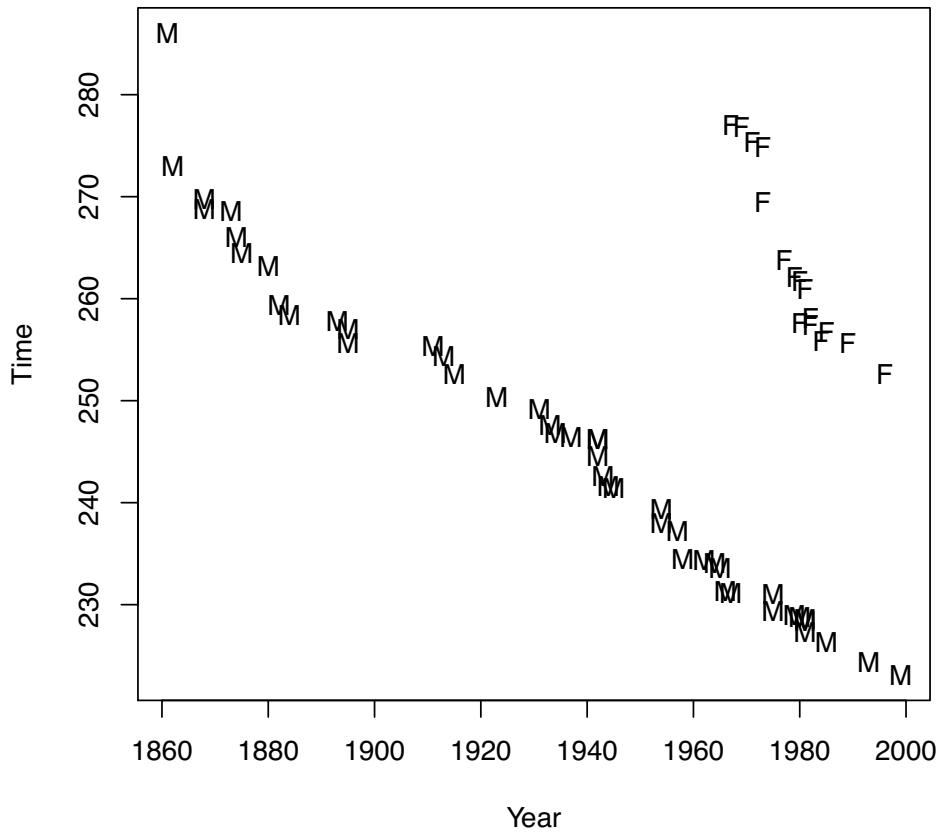
```
Res.Df   RSS Df Sum of Sq    F    Pr(>F)
1      59 850.63
2      58 518.03  1     332.6 37.238 9.236e-08 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
> qt(c(0.95,0.975,0.99,0.995), df=58)
```

```
[1] 1.671553 2.001717 2.392377 2.663287
```

```
> qt(c(0.95,0.975,0.99,0.995), df=59)
```

```
[1] 1.671093 2.000995 2.391229 2.661759
```



- Do you think that this data satisfies the linear model assumptions? Explain.
- What are the covariates used in the `imodel` object (in the context of the study)?
- What is being tested in the `anova` function call? What do you conclude from the results?
- Write down the final fitted models for the male and female records.
- Calculate a point estimate for the year when the female world record will equal the male world record. Do you expect this estimate to be accurate? Why or why not?
- Calculate a 95% confidence interval for the amount by which the gap between the male and female world records narrow every year.
- What is the hypothesis being tested in the `linearHypothesis` function call? What do you conclude from the output?

Question 7 (12 marks)

- (a) Randomisation eliminates confounding from both known and unknown factors. Explain why it is additionally advantageous to use blocking for some confounding factors.
- (b) A study is to be conducted to evaluate the effect of a drug on brain function. The evaluation consists of measuring the response of a particular part of the brain using an MRI scan. The drug is prescribed in doses of 1, 2 and 5 milligrams. Funding allows only 24 observations to be taken in the current study.

Explain how control might be used in a design of this experiment.

- (c) For the scenario above, explain how replication might be used in a design of this experiment.
- (d) For a complete block design with b blocks and k treatments, the reduced design matrix $X_{2|1}$ satisfies

$$X_{2|1}^T X_{2|1} = b \left[I_k - \frac{1}{k} J_k \right],$$

where I_k is the $k \times k$ identity matrix and J_k is the $k \times k$ matrix with all elements equal to 1. Show that

$$(X_{2|1}^T X_{2|1})^c = \frac{1}{b} I_k.$$

- (e) For the above complete block design, show that if a quantity $\mathbf{t}^T \boldsymbol{\tau}$ involving only the treatment parameters is estimable, then it must be a treatment contrast. (That is, $\mathbf{t}^T \mathbf{1} = 0$.)



THE UNIVERSITY OF MELBOURNE

Library Course Work Collections

Author/s:

Mathematics and Statistics

Title:

Linear Statistical Models, 2017 Semester1, MAST30025

Date:

2017

Persistent Link:

<http://hdl.handle.net/11343/190930>

Question 1 (9 marks)

- (a) Show that if A is a symmetric and idempotent matrix, then $r(A) = \text{tr}(A)$.
- (b) Show, without using the above result, that if X is an $n \times p$ matrix with $n > p$, then $r(X(X^T X)^c X^T) = r(X)$.
- (c) Let $Bx = g$ be a consistent linear system. Show that $x = B^c g$ is a solution to this system.

a) Diagonalise A as $A = PDP^T$ where D is diagonal matrix with diagonal entries as eigenvalues of A .

Since A is idempotent, eigenvalues of A are either 0 or 1.

\Rightarrow Diagonals of D are either 0 or 1.

$$\Rightarrow r(D) = \text{tr}(D)$$

$r(A) = r(PDP^T) = r(D)$ since P and P^T are orthogonal, thus invertible.

$$\text{tr}(A) = \text{tr}(PDP^T) = \text{tr}(DP^T P) = \text{tr}(D) \text{ since } P^T P = I.$$

$$\Rightarrow r(A) = \text{tr}(A).$$

b) $r(x) \geq r(x(X^T X)^c X^T) \geq r(X(X^T X)^c X^T) \geq r(x)$

$$\Rightarrow r(x) = r(X(X^T X)^c X^T).$$

c) $B B^c g = B B^c B x = B x = g$

Since $B(B^c g) = g$, $B^c g$ is a solution of $Bx = g$.

Question 2 (13 marks) Let

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \sim MVN \left(\begin{bmatrix} 0 \\ -1 \end{bmatrix}, \begin{bmatrix} 2 & 1 \\ 1 & 3 \end{bmatrix} \right).$$

You are given the following R calculations.

```
> V <- matrix(c(2,1,1,3),2,2)
> P <- eigen(V)$vectors
> P%*%diag(sqrt(eigen(V)$values))%*%t(P)  $\sqrt{V}$ 
[,1]      [,2]
[1,] 1.3763819 0.3249197
[2,] 0.3249197 1.7013016

> P%*%diag(1/sqrt(eigen(V)$values))%*%t(P)  $V^{-1/2}$ 
[,1]      [,2]
[1,] 0.7608452 -0.1453085
[2,] -0.1453085  0.6155367
```

- (a) Find two independent standard normal random variables which are linear combinations of \mathbf{y} elements (and constants).
- (b) Calculate $E[y_1^2 - 2y_1y_2]$.
- (c) Find all quadratic forms in \mathbf{y} which have a non-central χ^2 distribution with 2 degrees of freedom.
- (d) Show that $4y_1^2 - 4y_1y_2 + y_2^2$ is independent of $y_1^2 + 6y_1y_2 + 9y_2^2$.

a) We need $A\mathbf{y} + b \sim MVN(\mathbf{0}, \mathbf{I})$.

Let $\mathbf{y} \sim MVN(\mu, \Sigma)$, $\mu = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$, $\Sigma = \begin{bmatrix} 2 & 1 \\ 1 & 3 \end{bmatrix}$.

Σ is obviously invertible.

$\Rightarrow \mathbf{y} - \mu \sim MVN(\mathbf{0}, \Sigma)$

$\Sigma^{-1/2}(\mathbf{y} - \mu) \sim MVN(\mathbf{0}, \Sigma^{-1/2}\Sigma(\Sigma^{-1/2})^\top) \sim MVN(\mathbf{0}, \mathbf{I})$.

$\Sigma = PDP^\top$, where $D = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$, λ_1, λ_2 are eigenvalues of Σ .

$\Sigma^{-1/2} = PD^{-1/2}P^\top$, $D^{-1/2} = \begin{bmatrix} \sqrt{\lambda_1} & 0 \\ 0 & \sqrt{\lambda_2} \end{bmatrix}$

$\Rightarrow \Sigma^{-1/2} = \begin{bmatrix} 0.7608452 & -0.1453085 \\ -0.1453085 & 0.6155367 \end{bmatrix}$

$\Sigma^{-1/2}\mu = \begin{bmatrix} 0.1453085 \\ -0.6155367 \end{bmatrix}$

\Rightarrow Two elements of $\Sigma^{-1/2}\mathbf{y} - \Sigma^{-1/2}\mu$ are two independent $N(0, 1)$.

$$b) E[y_1^2 - 2y_1 y_2]$$

$$\begin{aligned} &= E\left[\begin{bmatrix} y_1 & y_2 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}\right] \\ &= \text{tr}\left(\begin{bmatrix} 1 & -1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 1 & 3 \end{bmatrix}\right) + \begin{bmatrix} 0 & -1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ -1 \end{bmatrix} \\ &= \text{tr}\left(\begin{bmatrix} 1 & -2 \\ -2 & 1 \end{bmatrix}\right) + 0 \\ &= 0 \end{aligned}$$

c) We need a symmetric 2×2 matrix A such that
 AV is idempotent and $r(AV)=2$.

An idempotent 2×2 matrix with rank 2, i.e. full rank, can only be I_2 .

$$\Rightarrow AV = I, A = V^{-1} = \frac{1}{6-1} \begin{bmatrix} 3 & -1 \\ -1 & 2 \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 3 & -1 \\ -1 & 2 \end{bmatrix}$$

Only satisfying quadratic form of y is $y^T V^{-1} y$, where $V^{-1} = \frac{1}{5} \begin{bmatrix} 3 & -1 \\ -1 & 2 \end{bmatrix}$.

$$d) 4y_1^2 - 4y_1 y_2 + y_2^2 = \begin{bmatrix} y_1 & y_2 \end{bmatrix} \begin{bmatrix} 4 & -2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = y^T A y.$$

$$y_1^2 + 6y_1 y_2 + 9y_2^2 = \begin{bmatrix} y_1 & y_2 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 3 & 9 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = y^T B y.$$

$$AVB = \begin{bmatrix} 4 & -2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 3 & 9 \end{bmatrix}$$

$$= \begin{bmatrix} 6 & -2 \\ -3 & 1 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 3 & 9 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

\Rightarrow independent.

Question 3 (14 marks) The following data is a sample of 5 random countries. For each country we measure the following variables:

- **logPPgdp**: The logarithm of the 2001 gross domestic product per person in US dollars;
- **logFertility**: The logarithm of the birth rate per 1000 females in the year 2000;
- **Purban**: The percentage of the population which lives in an urban area.

Name	logFertility	logPPgdp	Purban
Moldova	0.23	5.8	41
Netherlands	0.38	10.1	90
Estonia	0.14	8.3	69
Uganda	1.36	5.5	15
Hungary	0.13	8.6	65

We wish to predict **logFertility** using **logPPgdp** and **Purban**, using a linear model $\mathbf{y} = \mathbf{X}\beta + \varepsilon$ with no constant (intercept) term.

You are given the following R calculations.

```
> qt(0.975,1:5)
[1] 12.706205 4.302653 3.182446 2.776445 2.570582

> qf(0.95,1,1:5)
[1] 161.447639 18.512821 10.127964 7.708647 6.607891

> qf(0.95,2,1:5)
[1] 199.500000 19.000000 9.552094 6.944272 5.786135
```

- Calculate the least squares estimates of β .
- Calculate and interpret a 95% confidence interval for the parameter associated with **Purban**. You are given the sample variance $s^2 = 0.0887$.
- Test for model relevance at a 5% significance level.
- Croatia has a gross domestic product of \$4500 per person, and 58% of its population lives in an urban area. Calculate a 95% prediction interval for its birth rate.

$$a) \quad Y = \begin{bmatrix} 0.23 \\ 0.38 \\ 0.14 \\ 1.36 \\ 0.13 \end{bmatrix} \quad X = \begin{bmatrix} 5.8 & 41 \\ 10.1 & 90 \\ 8.3 & 69 \\ 5.5 & 15 \\ 8.6 & 65 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 308.75 & 2361 \\ 2361 & 18992 \end{bmatrix} \quad (X^T X)^{-1} = \begin{bmatrix} 0.066 & -0.008 \\ -0.008 & 0.001 \end{bmatrix}$$

$$b = (X^T X)^{-1} X^T Y = \begin{bmatrix} 0.310 \\ -0.034 \end{bmatrix}$$

$$b) \quad 95\% \text{ CI} \text{ is } b_2 \pm t_{0.975} \cdot \sqrt{s^2 \cdot (X^T X)^{-1}_{22}}$$

$$= -0.034 \pm 3.182446 \cdot \sqrt{0.0887} \times 0.001$$

$$= (-0.0640, 0.00403).$$

$$c) \quad SSR_{\text{Reg}} = Y^T X b = 1.817.$$

$$\text{F-stat} = \frac{SS_{\text{Reg}}/2}{s^2} = 10.242 > 9.552$$

\Rightarrow model is relevant at 5% level.

$$d) \quad t^T = [\log(4500) \ 58] = [8.412 \ 58]$$

$$\log(Y^*) = t^T \beta.$$

95% PI for $\log(Y^*)$.

$$t^T b \pm t_{0.975} \cdot \sqrt{s^2 \cdot (1 + t^T (X^T X)^{-1} t)}$$

$$= [8.412, 58] \begin{bmatrix} 0.310 \\ -0.034 \end{bmatrix} + 3.182446 \cdot \sqrt{0.0887} \times (1 + [8.412 \ 58] \begin{bmatrix} 1 & 1 \end{bmatrix})$$

$$= (a, b) \Rightarrow 95\% \text{ PI of } Y^* \text{ is } (e^a, e^b).$$

Question 4 (13 marks) Consider a full rank linear model $\mathbf{y} = \mathbf{X}\beta + \boldsymbol{\varepsilon}$. We wish to derive the formula for a prediction interval for the sum of the responses of two independent future observations with predictors \mathbf{x}_1 and \mathbf{x}_2 respectively. This uses the unbiased point estimator $(\mathbf{x}_1 + \mathbf{x}_2)^T \mathbf{b}$, where \mathbf{b} is the least squares estimator of β .

- Calculate the variance of the prediction error of this estimator.
- Show that this estimator is normally distributed.
- Show that this estimator is independent of SS_{Res} .
- Derive a t -distributed quantity based on this estimator and state its degrees of freedom.
- Thus write down a formula for a $100(1 - \alpha)\%$ prediction interval for the sum of two responses.

$$\begin{aligned} a) \quad & \text{var}(\mathbf{y}_1 + \mathbf{y}_2 - (\mathbf{x}_1 + \mathbf{x}_2)^T \mathbf{b}) \\ &= \text{var}(\mathbf{y}_1) + \text{var}(\mathbf{y}_2) + \text{var}((\mathbf{x}_1 + \mathbf{x}_2)^T \mathbf{b}) \\ &= \sigma^2 + \sigma^2 + (\mathbf{x}_1 + \mathbf{x}_2)^T \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{x}_1 + \mathbf{x}_2) \\ &= \sigma^2 \left[2 + (\mathbf{x}_1 + \mathbf{x}_2)^T (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{x}_1 + \mathbf{x}_2) \right] \end{aligned}$$

$$\begin{aligned} b) \quad & (\mathbf{x}_1 + \mathbf{x}_2)^T \mathbf{b} = (\mathbf{x}_1 + \mathbf{x}_2)^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \text{ which is in the form } \mathbf{A}\mathbf{y} \\ & \Rightarrow (\mathbf{x}_1 + \mathbf{x}_2)^T \mathbf{b} \text{ is normally distributed.} \end{aligned}$$

$$\begin{aligned} c) \quad & \text{Since } (\mathbf{x}_1 + \mathbf{x}_2)^T \mathbf{b} \text{ only depends on } \mathbf{b}, \quad SS_{Res} = (n-p)\sigma^2 \\ & \mathbf{b} \text{ and } \sigma^2 \text{ are independent} \Rightarrow (\mathbf{x}_1 + \mathbf{x}_2)^T \mathbf{b} \text{ is independent of } SS_{Res}. \end{aligned}$$

OR

$$SS_{Res} = \mathbf{y}^T (\mathbf{I} - \mathbf{H}) \mathbf{y}, \text{ where } \mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

$$(\mathbf{x}_1 + \mathbf{x}_2)^T \mathbf{b} = (\mathbf{x}_1 + \mathbf{x}_2)^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\begin{aligned} \mathbf{B}^T \mathbf{A} &= (\mathbf{x}_1 + \mathbf{x}_2)^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I} (\mathbf{I} - \mathbf{H}) \\ &= \sigma^2 \left[(\mathbf{x}_1 + \mathbf{x}_2)^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - (\mathbf{x}_1 + \mathbf{x}_2)^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{H} \right] \\ &= \sigma^2 \left[(\mathbf{x}_1 + \mathbf{x}_2)^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - (\mathbf{x}_1 + \mathbf{x}_2)^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right], \text{ since } \mathbf{X}^T \mathbf{H} = \mathbf{X}^T \\ &= 0 \\ \Rightarrow & (\mathbf{x}_1 + \mathbf{x}_2)^T \mathbf{b} \text{ is independent of } SS_{Res}. \end{aligned}$$

$$d) T = \frac{\bar{z}}{\sqrt{s^2/r}} \sim t_r$$

$$z = \frac{(y_1 + y_2) - (x_1 + x_2)^T b}{\sqrt{s^2(2 + (x_1 + x_2)^T (x^T x)^{-1} (x_1 + x_2))}}, \quad \chi^2 = \frac{SSRes}{s^2}, \quad r = np.$$

$$\Rightarrow T = \frac{(y_1 + y_2) - (x_1 + x_2)^T b}{\sqrt{s^2(2 + (x_1 + x_2)^T (x^T x)^{-1} (x_1 + x_2))}} / \sqrt{\frac{SSRes}{(n-p)s^2}}$$

$$= \frac{(y_1 + y_2) - (x_1 + x_2)^T b}{\sqrt{s^2(2 + (x_1 + x_2)^T (x^T x)^{-1} (x_1 + x_2))}} / \sqrt{\frac{s^2}{(n-p)s^2}}$$

$$= \frac{(y_1 + y_2) - (x_1 + x_2)^T b}{\sqrt{s^2(2 + (x_1 + x_2)^T (x^T x)^{-1} (x_1 + x_2))}} \sim t_{n-p}$$

Q) 100(1-\alpha)% PI for $y_1 + y_2$ is

$$(x_1 + x_2)^T b \pm t_{\alpha/2} \cdot \sqrt{s^2(2 + (x_1 + x_2)^T (x^T x)^{-1} (x_1 + x_2))}$$

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

- (a) State two methods that can be used to fit this model to data and compare them.
- (b) Explain the difference between a confidence interval and a prediction interval.
- (c) Define Akaike's information criterion and explain why it is useful as a goodness-of-fit measure for model selection.
- (d) Give an advantage and a disadvantage of stepwise selection over forward selection.
- (e) Define interaction between two categorical predictors.
- (f) List two principles of experimental design.
- (g) Explain what a Latin square is and its use in experimental design.

a) LSE or MLE.

Both give the same estimator of $\boldsymbol{\beta}$ which is BLUE.

LSE gives unbiased estimator of σ^2 but MLE of σ^2 is biased.

b) CI is interval of expected response

PI is interval of a single prediction, thus more variant. As a result,
PI is wider than CI.

c) $AIC = -2 \log(\text{likelihood}) + 2p$, p is number of parameters.

In general linear model, $AIC = \log\left(\frac{SS_{\text{Res}}}{n}\right) + 2p + \text{const.}$

Therefore, AIC has a balance on how good the model fits (by considering SS_{Res} , or $\log(\text{likelihood})$) as well as how complex the model is (by considering p), thus discouraging overfit.

d) Advantage: at each step, possible to add or remove variables \Rightarrow more flexible.

Disadvantage: may include insignificant variables since it's based on goodness-of-fit rather than F test.

e) There is interaction when effect of one factor differs for different levels of another factor.

f) Control and comparison;

Blocking;

Randomisation;

Blind and double blind experiments.

g) For a $n \times n$ Latin Square, every number from 1 to n appears exactly once at each row and column.

Latin Square is used in experimental designs where there are two or more confounding factors and CRD is impractical.

Question 6 (15 marks) Data was collected on the world record times for the one-mile run. For males, the records are from the period 1861–2003, and for females, from the period 1967–2003. This data is analysed below.

```
> mile <- read.csv('mile.csv', header=T)
> mile$Gender <- factor(mile$Gender)
> plot(Time ~ Year, data = mile, pch=as.character(Gender))
> imodel <- lm(Time ~ (Year + Gender)^2, data = mile)
> summary(imodel)
```

Call:

```
lm(formula = Time ~ (Year + Gender)^2, data = mile)
```

Residuals:

Min	1Q	Median	3Q	Max
-5.4512	-1.6160	-0.1137	1.1784	13.7265

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2309.4247	202.0583	11.429	< 2e-16 ***
Year	-1.0337	0.1021	-10.126	1.95e-14 ***
GenderMale	-1355.6778	203.1441	-6.673	1.03e-08 ***
Year:GenderMale	0.6675	0.1027	6.502	2.00e-08 ***

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 2.989 on 58 degrees of freedom

Multiple R-squared: 0.9663, Adjusted R-squared: 0.9645

F-statistic: 553.8 on 3 and 58 DF, p-value: < 2.2e-16

```
> amodel <- lm(Time ~ Year + Gender, data = mile)
> summary(amodel)
```

Call:

```
lm(formula = Time ~ Year + Gender, data = mile)
```

Residuals:

Min	1Q	Median	3Q	Max
-5.9071	-2.0988	-0.1141	1.2002	13.1863

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1003.00334	27.84691	36.02	<2e-16 ***
Year	-0.37364	0.01406	-26.57	<2e-16 ***
GenderMale	-34.85078	1.30099	-26.79	<2e-16 ***

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 3.896 on 59 degrees of freedom

```
Multiple R-squared:  0.9417,           Adjusted R-squared:  0.9397
F-statistic: 476.3 on 2 and 59 DF,  p-value: < 2.2e-16
```

```
> anova(amodel, imodel)
```

Analysis of Variance Table

```
Model 1: Time ~ Year + Gender
Model 2: Time ~ (Year + Gender)^2
Res.Df   RSS Df Sum of Sq    F    Pr(>F)
1      59 895.62
2      58 518.03  1     377.59 42.276 2.0001e-08 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
> linearHypothesis(imodel, c(0,1,0,1), -0.3)
```

Linear hypothesis test

Hypothesis:

Year + Year:GenderMale = - 0.3

```
Model 1: restricted model
Model 2: Time ~ (Year + Gender)^2
```

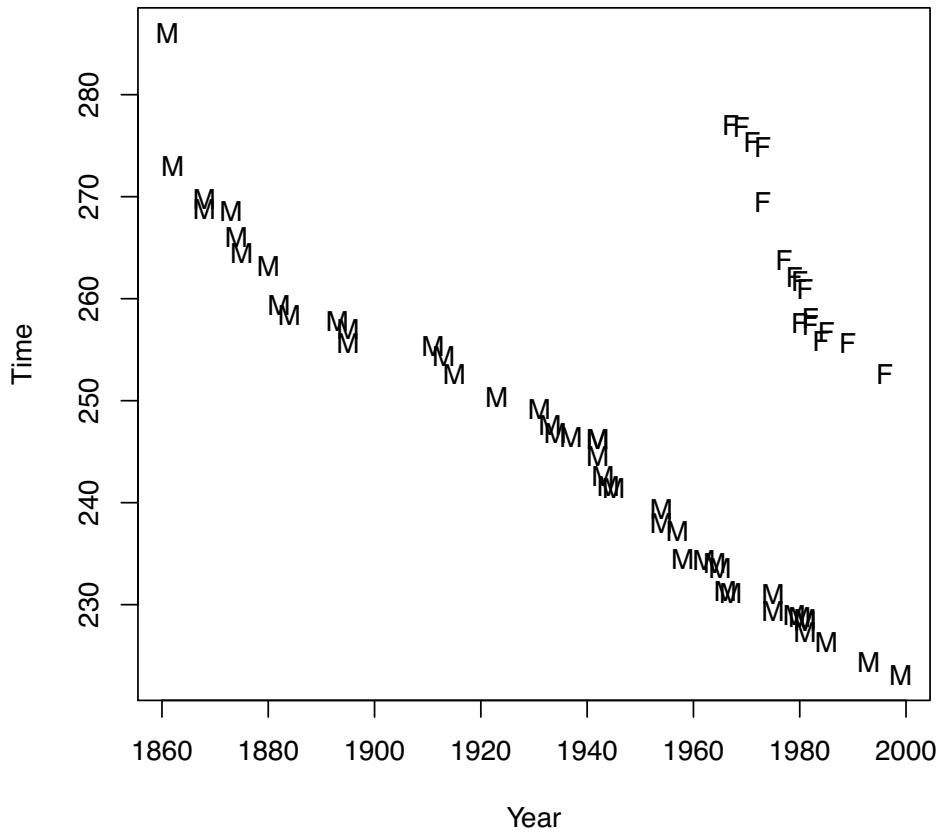
```
Res.Df   RSS Df Sum of Sq    F    Pr(>F)
1      59 850.63
2      58 518.03  1     332.6 37.238 9.236e-08 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
> qt(c(0.95,0.975,0.99,0.995), df=58)
```

[1] 1.671553 2.001717 2.392377 2.663287

```
> qt(c(0.95,0.975,0.99,0.995), df=59)
```

[1] 1.671093 2.000995 2.391229 2.661759



- Do you think that this data satisfies the linear model assumptions? Explain.
- What are the covariates used in the `imodel` object (in the context of the study)?
- What is being tested in the `anova` function call? What do you conclude from the results?
- Write down the final fitted models for the male and female records.
- Calculate a point estimate for the year when the female world record will equal the male world record. Do you expect this estimate to be accurate? Why or why not?
- Calculate a 95% confidence interval for the amount by which the gap between the male and female world records narrow every year.
- What is the hypothesis being tested in the `linearHypothesis` function call? What do you conclude from the output?

a) No. In this context, records can only decrease.

That means responses are not independent.

b) Year, which is continuous.

Gender, which is categorical and is either "male" or "female"

Interaction between Year and Gender.

c) ANOVA tests for relevance of interaction.

From the result, p-value is well below 0.05, implying there is significant interaction.

d) We use model with interaction.

Coefficients:

	Estimate	
(Intercept)	2309.4247	$\mu + \tau_1$
Year	-1.0337	$\beta + \gamma_1$
GenderMale	-1355.6778	$\tau_2 - \tau_1$
Year:GenderMale	0.6675	$\gamma_2 - \gamma_1$

For Female: Time = $2309.4247 - 1.0337 \text{ Year}$. i.e. $y_i = \mu + \tau_1 + (\beta + \gamma_1)x$

$$\begin{aligned} \text{For Male: } \text{Intercept} &= \mu + \tau_2 \\ &= (\mu + \tau_1) + (\tau_2 - \tau_1) \\ &= 2309.4247 - 1355.6778 \\ &= 953.7469 \end{aligned}$$

$$\begin{aligned} \text{gradient: } \beta + \gamma_2 &= (\beta + \gamma_1) + (\gamma_2 - \gamma_1) \\ &= -1.0337 + 0.6675 \\ &= -0.3662. \end{aligned}$$

$\Rightarrow \text{Time} = 953.7469 - 0.3662 \text{ Year}$. $y_i = \mu + \tau_2 + (\beta + \gamma_2)x$

$$2) 2309.4247 - 1.0337x = 953.7469 - 0.3662x$$

$$\Rightarrow x = 2030.95$$

OR $x = \frac{1355.6778}{0.6675} \rightarrow$ how much difference between M and F at year 0.
 \rightarrow how much F catch up the difference each year.
 $= 2030.95$

Not accurate, since we are extrapolating beyond the relevant range.

f) Want CI for $\bar{z}_2 - \bar{z}_1$,

$$95\% \text{ CI is } \bar{z}_2 - \bar{z}_1 \pm se(\bar{z}_2 - \bar{z}_1) \cdot t_{0.975} (\text{df}=58)$$

$$= 0.6675 \pm 0.1027 \times 2.001717$$

$$= (0.462, 0.873).$$

g) H_0 : time for male records decreases at a rate of 0.3 per year.

Since p-value is very small, it implies time records do not decrease at 0.3 per year.

Question 7 (12 marks)

- (a) Randomisation eliminates confounding from both known and unknown factors. Explain why it is additionally advantageous to use blocking for some confounding factors.
- (b) A study is to be conducted to evaluate the effect of a drug on brain function. The evaluation consists of measuring the response of a particular part of the brain using an MRI scan. The drug is prescribed in doses of 1, 2 and 5 milligrams. Funding allows only 24 observations to be taken in the current study.
- Explain how control might be used in a design of this experiment.

- (c) For the scenario above, explain how replication might be used in a design of this experiment.
- (d) For a complete block design with b blocks and k treatments, the reduced design matrix $X_{2|1}$ satisfies

$$X_{2|1}^T X_{2|1} = b \left[I_k - \frac{1}{k} J_k \right],$$

where I_k is the $k \times k$ identity matrix and J_k is the $k \times k$ matrix with all elements equal to 1.
Show that

$$(X_{2|1}^T X_{2|1})^C = \frac{1}{b} I_k.$$

- (e) For the above complete block design, show that if a quantity $\mathbf{t}^T \boldsymbol{\tau}$ involving only the treatment parameters is estimable, then it must be a treatment contrast. (That is, $\mathbf{t}^T \mathbf{1} = 0$.)

a) Blocking reduces variance within each block.

b) A control group taking placebo having 0 milligrams of drug.

Effectively 4 treatments with $\{0, 1, 2, 5\}$ milligrams

$$c) X_{2|1}^T X_{2|1} \left(\frac{1}{b} I_k \right) X_{2|1}^T X_{2|1}$$

$$= b [I_k - \frac{1}{k} J_k] \left(\frac{1}{b} I_k \right) \cdot b [I_k - \frac{1}{k} J_k]$$

$$= b [I_k - \frac{1}{k} J_k] [I_k - \frac{1}{k} J_k]$$

$$= b [I_k - \frac{1}{k} J_k - \frac{1}{k} J_k + \frac{1}{k^2} J_k J_k]$$

$$= b [I_k - \frac{2}{k} J_k + \frac{1}{k^2} (k J_k)]$$

$$= b [I_k - \frac{2}{k} J_k + \frac{1}{k} J_k]$$

$$= b [I_k - \frac{1}{k} J_k] = X_{2|1}^T X_{2|1}$$

$\Rightarrow \frac{1}{b} I_k$ is a conditional inverse of $X_{2|1}^T X_{2|1}$.

$$\text{i.e. } (X_{2|1}^T X_{2|1})^C = \frac{1}{b} I_k.$$

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 3 & 3 \\ 3 & 3 & 3 \\ 3 & 3 & 3 \end{bmatrix} = 3 \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

e) Since $t^T \tau$ is estimable.

$$t^T (x_{21}^T x_{21})^c x_{21}^T x_{21} = t^T$$

$$\Rightarrow t^T \left(\frac{1}{b} I_k \right) b [I_k - \frac{1}{k} J_k] = t^T$$

$$= t^T [I_k - \frac{1}{k} J_k] = t^T$$

$$t^T 1 = t^T [I_k - \frac{1}{k} J_k] 1$$

$$= t^T 1 - \frac{1}{k} t^T (J_k 1)$$

$$= t^T 1 - \frac{1}{k} t^T (k 1)$$

$$= t^T 1 - t^T 1$$

$$= 0$$

$$\begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 3 \\ 3 \end{bmatrix} = 3 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$