



The
University
Of
Sheffield.

MAS6003

SCHOOL OF MATHEMATICS AND STATISTICS

Spring Semester 2016–2017

Linear Models

3 hours

Restricted Open Book Examination.

Candidates may bring to the examination lecture notes and associated lecture material (but no textbooks) plus a calculator which conforms to University regulations.

*Marks will be awarded for your best **five** answers.*

Total marks 100.

Please leave this exam paper on your desk
Do not remove it from the hall

Registration number from U-Card (9 digits)
to be completed by student

--	--	--	--	--	--	--	--	--

Blank

1 The level of light output of bulbs covered with two types of coating (A and B) is studied. Data are collected and given in the following table.

Length of operation (hours)	Coat	Drop in light output (percent of original output)
0	A	0
400	A	6
800	A	22
1200	A	27
1600	A	32
2000	A	36
2400	A	38
0	B	0
400	B	4
800	B	6
1200	B	9
1600	B	10
2000	B	11
2400	B	12

(i) Write down, assess and interpret the model for which R output is given in Table 1 overleaf. **(6 marks)**

(ii) Assuming that the rate of drop in output per hour of duration for both types of bulbs is identical, set up a general hypothesis and a test procedure for checking whether the relationship of output to length of operation is the same for both types of bulbs. Do not carry out the test, but provide a detailed description of how this would be done. **(7 marks)**

(iii) A new variable is created, each element of which is the result of averaging the two values of the drop in light output for coats A and B (see table below).

Length of operation (hours)	Averaged Drop in light output (percent of original output)
0	0
400	5
800	13
1200	18
1600	21
2000	23.5
2400	25

Then, a new linear model is set up by regressing the new Averaged Drop in light output on the length of operation. The R output of this model is given in Table 2. Comment on the performance of this model and state its limitations, if any. **(5 marks)**

(iv) Compare the models obtained in (i) and (iii) and suggest which of the two is better. **(2 marks)**

1 (continued)

Table 1

Call: `lm(formula = Drop ~ Length + factor(Coat) - 1)`

Residuals:

Min	1Q	Median	3Q	Max
-10.25	-4.116	2.536	4.375	5.321

Coefficients:

	Value	Std. Error	t value	Pr(> t)
Length	0.0106	0.0020	5.2080	0.0003
Coat A	10.2500	3.3646	3.0464	0.0111
Coat B	-5.3214	3.3646	-1.5816	0.1421

Residual standard error: 6.107 on 11 degrees of freedom

Multiple R-Squared: 0.9256

F-statistic: 45.59 on 3 and 11 degrees of freedom,
the p-value is 1.703e -006

Table 2

Call: `lm(formula = Drop ~ Length)`

Residuals:

1	2	3	4	5	6	7
-2.214	-1.5	2.214	2.929	1.643	-0.1429	-2.929

Coefficients:

	Value	Std. Error	t value	Pr(> t)
(Intercept)	2.2143	1.7219	1.2860	0.2548
Length	0.0107	0.0012	8.9742	0.0003

Residual standard error: 2.527 on 5 degrees of freedom

Multiple R-Squared: 0.9415

F-statistic: 80.54 on 1 and 5 degrees of freedom,
the p-value is 0.00028 66

2 Consider the exponential family of distributions for the random variable Y , with probability function

$$f(y; \theta, \phi) = \exp \left[\frac{y\theta - b(\theta)}{\phi/w} + c(y, \phi) \right], \quad (1)$$

where θ is the natural parameter, ϕ is the dispersion parameter, w is the weight and the function $b(\theta)$ is assumed to be twice differentiable.

(i) Using the result that the variance of the score statistic is equal to minus the expectation of the second partial derivative of the log-likelihood of θ , with respect to θ , show that the variance of Y is

$$\text{Var}(Y) = \frac{\phi}{w} b''(\theta),$$

where $b''(\theta)$ is the second derivative of $b(\cdot)$ with respect to θ . **(5 marks)**

(ii) Show that the mode $\hat{Y} = \hat{y}$ of the distribution of Y , evaluated at observation y , satisfies the differential equation

$$\frac{\partial c(\hat{y}, \phi)}{\partial y} = -\frac{w\theta}{\phi},$$

where Y is continuously distributed and $c(y, \phi)$ is assumed to be differentiable with respect to y . **(4 marks)**

(iii) Consider the two-parameter inverse Gaussian distribution, with p.d.f.

$$f(y; \mu, \lambda) = \sqrt{\frac{\lambda}{2\pi y^3}} \exp \left[-\frac{\lambda(y - \mu)^2}{2\mu^2 y} \right],$$

where $\mu, \lambda > 0$ are the two parameters and $f(y; \mu, \lambda)$ is positive for $y > 0$.

(a) Write $f(y; \mu, \lambda)$ in the exponential form (1), hence determine θ , ϕ , w , $b(\theta)$ and $c(y, \phi)$ in terms of μ and λ . **(6 marks)**

(b) Using the result of part (i) find the variance of Y . **(2 marks)**

(c) Write down the canonical link of the distribution $f(y; \mu, \lambda)$. Explain why the canonical link may not be a good choice if we wish to fit a generalised linear model with that link. Hence, suggest another link function which is a better choice. **(3 marks)**

3 The following table summarizes the response to chemotherapy of 30 patients suffering from lymphoma, a form of cancer.

	Gender of patient			
	Male ($j = 1$)		Female ($j = 2$)	
	Tumour type		Tumour type	
	Nodular	Diffuse	Nodular	Diffuse
	$k = 1$	$k = 2$	$k = 1$	$k = 2$
No response ($i = 1$)	1	12	2	3
Response ($i = 2$)	4	1	6	1

Gender and tumour type are controlled variables. Generalized linear models with Poisson errors and a log link function fitted to these data gave the following results.

Model fitted	Deviance	df	
G*T+R	14.83	3	R = no response/response
G*T+R*G	12.02	2	G = gender
G*T+R*T	0.81	?	T = type of tumour
G*T+R*T+R*G	0.65	†	

(i) Identify the degrees of freedom entries represented by ? and † in these results. (3 marks)

(ii) Let π_{ijk} be the conditional probability $\pi_{ijk} = P(R = i | G = j, T = k)$, for any $i, j, k = 1, 2$. Verify that

$$\sum_i \hat{\pi}_{ijk} = 1, \quad \text{for all } j, k,$$

where $\hat{\pi}_{ijk}$ is the estimate of the probability π_{ijk} in the model G*T+R*G. (2 marks)

(iii) Show (without using the R output) that the fitted value $\hat{\mu}_{112}$ for the G*T+R*G model is 9.4. (2 marks)

(iv) By first finding the fitted value using $\mu_{ijk} = n_{jk}\pi_{ijk}$, calculate the Pearson and deviance residuals for observation y_{221} for the G*T+R*T model. You must not use the R output in this question. (4 marks)

(v) Using **only** the R output on the next page, verify the fitted value in (iii) above. (2 marks)

(vi) Using changes in scaled deviance in the table above, what conclusions about the dependence of response on gender and type of tumour would you draw? Is the model you decide on a good fit? (5 marks)

(vii) What model-checking would you perform? (2 marks)

3 (continued)

Call:

```
glm(formula = count ~ factor(gender) * factor(type) + factor(response) *
     factor(gender), family = poisson(log), data = data)
```

Deviance Residuals:

1	2	3	4	5	6	7	8
-1.6292	0.8166	-0.7895	0.9274	1.8000	-1.6292	0.5909	-0.9859

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	1.28402	0.47050	2.729	0.00635 **
factor(gender)2	-0.08004	0.68046	-0.118	0.90636
factor(type)2	0.95551	0.52623	1.816	0.06941 .
factor(response)2	-0.95551	0.52623	-1.816	0.06941 .
factor(gender)2:factor(type)2	-1.64866	0.80742	-2.042	0.04116 *
factor(gender)2:factor(response)2	1.29198	0.78726	1.641	0.10077

The following R output is additionally provided

```
> qchisq(0.95, 2)
[1] 5.991465
```

4 The Junior School Project collected longitudinal data on the performance of students from 49 primary schools in inner London. Over a period of 6 years, for each pupil they recorded

`school` - the school they attend (a factor with levels 1 to 49)

`class` - their class within the school (a factor with levels 1 to 4)

`gender` - the pupils gender (a factor with levels boy and girl)

`english` - their score on the English test that year (a number between 1 and 40)

`year` - the school year in which the test is taken (a number between 1 and 6)

`id` - a identification number unique to each student.

Interest lies in exploring the differences in English test performance between boys and girls. The data are stored in R as the data frame `jsp`.

(i) Initially, only the test results for year 2 are used. A model is fitted in R using the command

```
> fit1 <- lmer(english ~ gender-1+(1|school/class), data=jsp, subset = year==2)
```

If Y_{ijk} denotes the test score of student k in class j in school i , write down the equation of the model that has been used, defining your notation carefully.

(3 marks)

(ii) Explain, given the aim of the study, why the designation of explanatory variables as fixed and random is reasonable, and comment upon the structure of the random effects.

(3 marks)

4 (continued)

(iii) Use the R output below to give estimates for all of the model parameters.

```
> summary(fit1)
Linear mixed model fit by REML ['lmerMod']
Formula: english ~ gender - 1 + (1 | school/class)
Data: jsp
Subset: jsp$year == 2

REML criterion at convergence: 8462.7

Scaled residuals:
    Min       1Q   Median       3Q      Max
-2.41335 -0.75200 -0.09772  0.68060  2.90022

Random effects:
 Groups          Name          Variance Std.Dev.
class:school (Intercept)  34.15    5.843
school      (Intercept)  53.93    7.344
Residual                                382.99   19.570
Number of obs: 953, groups: class:school, 90; school, 48

Fixed effects:
              Estimate Std. Error t value
genderboy    39.751      1.581    25.14
gendergirl   44.588      1.564    28.51

Correlation of Fixed Effects:
              gndrby
gendergirl  0.652
```

(2 marks)

(iv) What is the estimated correlation between the test results for two students in the same school but different classes?

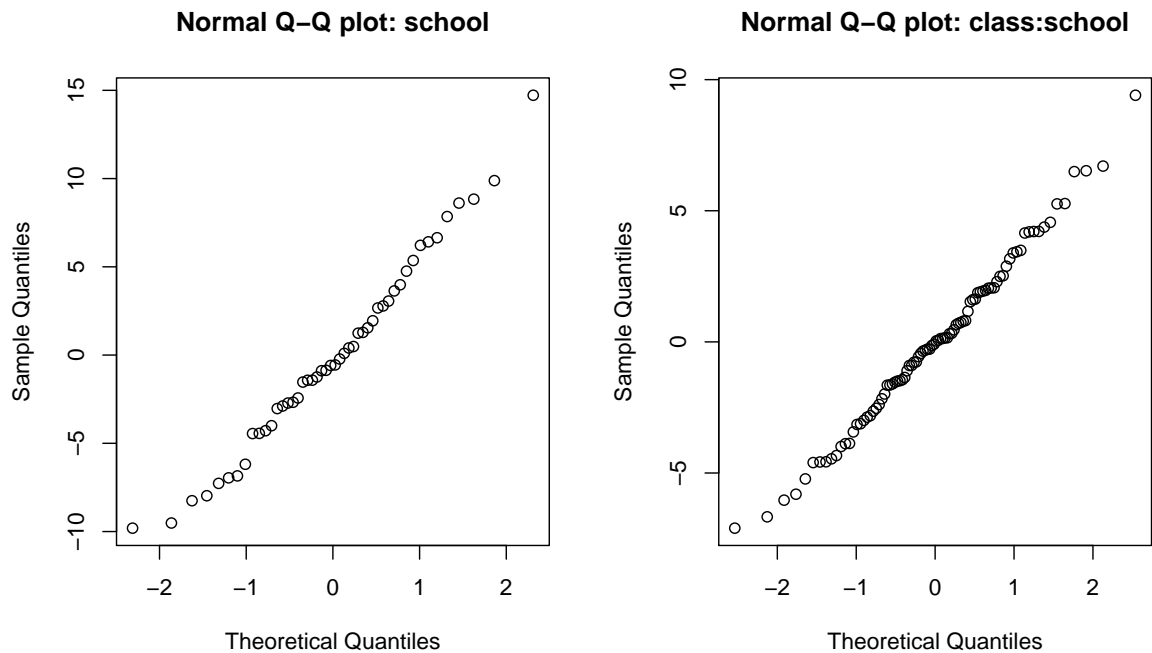
(4 marks)

(v) The R commands

```
> par(mfrow=c(1,2))
> qqnorm(unlist(ranef(fit1)$'school'))
> qqnorm(unlist(ranef(fit1)$'class:school'))
```

were used to produce the following plots.

4 (continued)



Explain the reason for producing these plots. What do you conclude from these plots?

(2 marks)

(vi) The R session is continued below.

```
> jsp3=dplyr::filter(jsp, year==2)
> attach(jsp3)
> fit1 <- lmer(english ~ gender-1+(1|school/class), REML=F)
> fit2 <- lmer(english ~ 1+(1|school/class), REML=F)
>
> Lambda <- replicate(10^4,{
+   english.new<-unlist(simulate(fit2))
+   fit1.new <- lmer(english.new ~ gender-1+(1|school/class), REML=F)
+   fit2.new <- lmer(english.new ~ 1+(1|school/class), REML=F)
+   -2*(logLik(fit2.new) - logLik(fit1.new))
+ })
>
> Lambda.obs <- -2*(logLik(fit2) - logLik(fit1))
> sum(Lambda>Lambda.obs)
[1] 5
```

Give the name of the procedure that has been used, state what is being tested, and interpret the output.

(4 marks)

4 (continued)

(vii) The full dataset contains the test results for the pupils over their 6 years at school (longitudinal data). Write down the R command you would use to fit an appropriate mixed-effects model to investigate the trend in test performance for pupils over the 6 years they spend at primary school.

(2 marks)

5 The sequence of random variables X_1, \dots, X_n are iid with an exponential distribution with mean $\frac{1}{\lambda}$, i.e., they have pdf

$$f(x; \lambda) = \begin{cases} \lambda \exp(-\lambda x) & \text{for } x \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

The X_i are not observed directly. Instead, we observe

$$Y_i = \mathbb{I}_{X_i \leq z_i} = \begin{cases} 1 & \text{if } X_i \leq z_i \\ 0 & \text{otherwise} \end{cases}$$

for some known sequence of thresholds z_1, \dots, z_n . In other words, X_i is left-censored at z_i if $Y_i = 1$ and right-censored at z_i if $Y_i = 0$. The data can thus be summarized as $(z_1, Y_1), \dots, (z_n, Y_n)$.

(i) Derive the likelihood function for λ given Y_1, \dots, Y_n and explain why this might necessitate the use of the EM algorithm if we want to fit the model.

(4 marks)

(ii) Derive the value for

$$R_i := \mathbb{E}(X_i | X_i > z_i).$$

Hint: recall the memoryless property of the exponential distribution.

(3 marks)

(iii) Show that

$$L_i := \mathbb{E}(X_i | X_i \leq z_i) = \frac{1 - (1 + \lambda z_i)e^{-\lambda z_i}}{\lambda(1 - e^{-\lambda z_i})}.$$

Hint: you may find your answer from part (i) and the law of total expectation useful

$$\mathbb{E}(X) = \mathbb{E}(\mathbb{E}(X|Y)).$$

(6 marks)

(iv) Hence derive an EM-algorithm for finding the maximum likelihood estimator of λ using $(X_1, y_1, z_1), \dots, (X_n, y_n, z_n)$ as the completed dataset. You may leave your answer in terms of L_i and R_i .

(7 marks)

6 (i) Consider a bivariate sample on (Y_1, Y_2) from a model that has parameters μ and ψ . For each of the following missing-data mechanisms, state whether the data are missing completely at random (MCAR), missing at random (MAR), or not missing at random (NMAR).

$$(a) \quad \mathbb{P}(y_2 \text{ missing} | y_1, y_2, \mu, \sigma, \psi) = \frac{e^{\mu + \psi y_1}}{1 + e^{\mu + \psi y_1}}$$

$$(b) \quad \mathbb{P}(y_2 \text{ missing} | y_1, y_2, \mu, \sigma, \psi) = \frac{e^{\mu + \psi}}{1 + e^{\mu + \psi}}$$

$$(c) \quad \mathbb{P}(y_2 \text{ missing} | y_1, y_2, \mu, \sigma, \psi) = \frac{e^{\mu + \psi y_1 y_2}}{1 + e^{\mu + \psi y_1 y_2}}$$

Which of these mechanisms are ignorable for likelihood-based inference?

(4 marks)

(ii) Consider the following dataset, where NA is used to denote a missing observation.

x	y
-6.3	-9.5
1.8	6.4
NA	-13.7
16.0	NA

Interest lies in the expected value of y . Calculate the mean of y using

- (a) Complete-case analysis
- (b) Available-case analysis
- (c) Mean imputation
- (d) (Mean) Regression imputation

You may find the following R output useful.

```
> coef(lm(y~x))
(Intercept)          x
  2.866667    1.962963
> coef(lm(x~y))
(Intercept)          y
-1.460377    0.509434
```

(5 marks)

6 (continued)

(iii) Consider a dataset from a large postal survey on the psychology of debt. It contains responses from $n = 464$ individuals. The covariates are

- AttDebt - score on a scale of attitudes to debt (high values=favourable to debt)
- CCard - how often did s/he use credit cards (1=never... 3=regularly)
- BuildSoc - does the respondent have a building society account?
- LOC - the respondents' *locus of control*, which is the degree to which the respondent believes that they have control over the outcome of events in their life (high values = belief that one's life can be controlled).

The output below shows the structure of the dataset.

```
> str(Debt)
'data.frame': 464 obs. of 4 variables:
 $ AttDebt : num  2.71 3.88 3.06 4.29 3.82 ...
 $ CCard   : Factor w/ 3 levels "1","2","3": 2 3 2 2 3 2 2 1 NA 2 ...
 $ BuildSoc: Factor w/ 2 levels "0","1": NA NA NA 1 1 1 1 1 1 NA ...
 $ LOC     : num  2.83 4.83 3.83 4.83 3.17 ...
> head(Debt)
  AttDebt CCard BuildSoc LOC
1    2.71     2    <NA> 2.83
2    3.88     3    <NA> 4.83
3    3.06     2    <NA> 3.83
4    4.29     2     0 4.83
5    3.82     3     0 3.17
6    3.06     2     0 3.83
```

(a) The following R command is used.

```
> Debt.mice <- mice(Debt, m=5, method = c('norm', 'polyreg', 'logreg', 'norm'))
```

Describe in detail the statistical procedure that is used to fill in the missing values.

(6 marks)

6 (continued)

(b) Interest lies in the coefficient of LOC (β_3) in the linear regression model

$$\text{AttDebt} = \beta_0 + \beta_1 \text{CCard} + \beta_2 \text{BuildSoc} + \beta_3 \text{LOC}$$

Use the R output below to calculate an expected value of β_3 and its standard error.

```
> fit.mice <- with(Debt.mice, lm(AttDebt ~ CCard+BuildSoc+LOC))
> (coefs = sapply(fit.mice$analyses, coef))
      [,1] [,2] [,3] [,4] [,5]
(Intercept) 3.61 3.74 3.73 3.68 3.75
CCard2      0.11 0.15 0.17 0.19 0.12
CCard3      0.45 0.45 0.41 0.39 0.45
BuildSoc2   -0.21 -0.19 -0.16 -0.13 -0.16
LOC         -1.12 -2.11 -0.53 -1.32 -1.01
>
> apply(coefs, 1, var)
(Intercept)      CCard2      CCard3      BuildSoc2      LOC
 0.00337      0.00112      0.00080      0.00095      0.33307
>
> sapply(fit.mice$analyses, function(x) vcov(x)[5,5])
[1] 0.67 0.78 0.46 0.90 0.46
```

(5 marks)

End of Question Paper