



The  
University  
Of  
Sheffield.

**MAS473**

**SCHOOL OF MATHEMATICS AND STATISTICS**

**Spring Semester  
2013–2014**

**MAS473 Extended linear models**

**2 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.*

*Answer all questions. Total marks 60.*

- 1 An experiment has been conducted to investigate the time taken for a clot to form in a sample of blood, when treated with a drug. Three drugs are compared. There are ten volunteers in the study. Each volunteer donates six blood samples, and each drug is used on two of these samples, so that there are 60 observations in total. In an R dataset, the clot formation time (in seconds) is stored as a variable `clot.time` and the drug and volunteer labels are stored as factor variables `drug` and `volunteer` respectively. Below is some output from an R session.

```
> fm1<-lmer(clot.time~drug-1+(1|volunteer/drug))
> summary(fm1)
Linear mixed model fit by REML ['lmerMod']
Formula: clot.time ~ drug - 1 + (1 | volunteer/drug)
```

REML criterion at convergence: -44.3055

Random effects:

Groups	Name	Variance	Std.Dev.
drug:volunteer	(Intercept)	0.030827	0.17558
volunteer	(Intercept)	0.001150	0.03391
Residual		0.008216	0.09064

Number of obs: 60, groups: drug:volunteer, 30; volunteer, 10

Fixed effects:

	Estimate	Std. Error	t value
drug1	10.07687	0.06007	167.8
drug2	10.99666	0.06007	183.1
drug3	12.04935	0.06007	200.6

Correlation of Fixed Effects:

	drug1	drug2
drug2	0.032	
drug3	0.032	0.032

- (i) Write down the equation of the model that has been fitted and assigned to the name `fm1`, defining your notation carefully. *(3 marks)*
- (ii) Give the estimated parameter values for each of the variance parameters in your model in (i). *(1 mark)*
- (iii) Calculate the estimated variance for any observation. *(1 mark)*
- (iv) Calculate the estimated correlation between any two different observations involving the same volunteer and the same drug. *(2 marks)*

1 (continued)

- (v) The estimator for the fixed effect for drug  $i$  is the sample mean of all the observations involving drug  $i$ .
- (a) Derive an expression for the variance of this estimator, and hence verify that the estimated standard error is 0.06 (to 2 d.p.)  
*(3 marks)*
- (b) Derive an expression for the correlation between the fixed effect estimators of two different drugs, and hence verify that the estimated correlation is 0.03 (to 2 d.p.)  
*(4 marks)*
- (c) A fixed effects model is also fitted to the data using the command `lm(clot.time~drug*volunteer-1,contrasts=list(volunteer=contr.sum))`. Edited output corresponding to drug1 is given below.

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
drug1	10.051984	0.020269	495.938	< 2e-16 ***

Residual standard error: 0.09064 on 30 degrees of freedom

The drug1 estimate is again given by the mean of all the observations involving drug 1. Explain why the estimated standard error is smaller compared to model fm1. Explain the difference in interpretation of the drug 1 term between the fixed and mixed effects models.  
*(3 marks)*

- (vi) The session is continued below.

```
> fm2<-lmer(clot.time~drug-1+(1|volunteer))
> logLik(fm1)
'log Lik.' 22.15275 (df=6)
> logLik(fm2)
'log Lik.' 9.289661 (df=5)
> qchisq(0.99,1)
[1] 6.634897
```

Compare the models defined as fm1 and fm2 using a generalised likelihood ratio test. State clearly what the null hypothesis is, in terms of the parameters of model fm1, and interpret the result.  
*(3 marks)*

2 A car tyre manufacturing company wants to assess the relationship between car tyre thickness (labelled as `thickness` in R) and the probability of a tyre splitting after 20,000 miles of use on two different road surfaces A and B (labeled as `surface` in R). For various tyre thicknesses and for each road surface they determine what proportion of the tyres split. The number of tyres tested at each tyre thickness is labelled as `tested` and the number splitting at each tyre thickness is labeled as `split`. The information is given in the table below.

(i) The following command is used to fit a model to the data:

```
lm1<-glm(split/tested~surface*thickness,weights=tested,
family=binomial(logit))
```

Write down, in terms of the linear predictor  $\eta_i$ , the statistical model for  $E(y_i)$  that is fitted to the data. Specify  $\eta_i$  in terms of the variables and parameters in the model. (3 marks)

(ii) The deviance in a glm with a binomial response is given by

$$D(y, \hat{\mu}) = -2 \sum_i n_i \left\{ y_i \log \left( \frac{\hat{\mu}_i}{y_i} \right) + (1 - y_i) \log \left( \frac{1 - \hat{\mu}_i}{1 - y_i} \right) \right\}.$$

For surface A, given the information in the table below, write down the numerical values of  $n_i$  and  $y_i$  for all values of  $i$ . (2 marks)

Surface A			Surface B		
thickness (mm)	number tested	number split	thickness (mm)	number tested	number split
2.3	100	75	2.3	100	88
2.9	50	25	2.9	50	26
3.4	40	11	3.4	40	14
3.9	50	10	3.9	50	12
4.3	50	3	4.3	50	13

2 (continued)

(iii) Some further edited R commands and output are provided below

```
> summary(lm1)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	5.09125	0.64787	7.858	3.89e-15
surfaceB	0.04737	0.90763	0.052	0.958
thickness	-1.74354	0.21124	-8.254	< 2e-16
surfaceB:thickness	0.16636	0.28716	0.579	0.562

```
> anova(lm1)
```

Analysis of Deviance Table

Model: binomial, link: logit

Response: split/tested

Terms added sequentially (first to last)

	Df	Deviance	Resid. Df	Resid. Dev
NULL			9	194.483
surface	1	5.822	8	188.662
thickness	1	176.614	7	12.048
surface:thickness	1	0.337	6	11.711

```
> print(vcov(lm1),digits=3)
```

	(Intercept)	surfaceB	thickness	surfaceB:thickness
(Intercept)	0.420	-0.420	-0.1336	0.1336
surfaceB	-0.420	0.824	0.1336	-0.2543
thickness	-0.134	0.134	0.0446	-0.0446
surfaceB:thickness	0.134	-0.254	-0.0446	0.0825

```
> qchisq(0.95,1)
```

```
[1] 3.841459
```

- By performing a series of hypothesis tests, assess the effects of road surface and tyre thickness on the response. *(4 marks)*
- Using the output above, calculate the Pearson residual for 3.4 mm thick tyres tested on surface A. *(4 marks)*
- Using the output above, calculate the odds ratio of a tyre with a thickness of 5mm splitting on surface A compared to a tyre with a thickness of 3.4mm splitting on surface A. Give a 95% confidence interval for this odds ratio and interpret the result. *(5 marks)*

**2** (continued)

- (d) Give an example, in words, of an odds ratio of a tyre splitting involving the variables in the output above that requires the covariance between the parameter estimate for **surface** and the parameter estimate for the interaction between **surface** and **thickness** but does not require any of the covariances involving the parameter estimate for **thickness**. *(2 marks)*

- 3 Data were collected relating to collisions involving cyclists in a particular city over a given period of time. The interest was in the factors affecting the outcome for cyclists (indicated by the variable `outcome`). The two main factors of interest were whether the cyclist was wearing a helmet (indicated by the variable `helmet`) and what the cyclist collided with (indicated by the variable `collision`). `Collision` is either `lorry`, `car` or `pedestrian`. `Outcome` is either `serious` or `minor`. The data is stored in the data frame `cycling` in R. The first few lines of the data frame are shown below.

```
> head(cycling)
  count collision helmet outcome
1    25    lorry     1         1
2    23    lorry     1         0
3    18    lorry     0         1
4    12    lorry     0         0
5    10     car      1         1
6    44     car      1         0
```

The full data are shown below.

helmet			no helmet		
collision	serious	minor	collision	serious	minor
lorry	25	23	lorry	18	12
car	10	44	car	17	20
pedestrian	1	8	pedestrian	1	6

- (i) Explain why `outcome` is a response whilst `helmet` and `collision` are controlled factors. What is the minimal model when fitting linear models to these data? *(2 marks)*
  
- (ii) By looking at the proportions of cyclists suffering serious accidents in each category, make some brief initial observations about the relationship between the probability of having a serious cycling accident and helmet use and type of collision in this data set. *(3 marks)*

3 (continued)

- (iii) Various models are fitted to the data. A summary of the residual deviances for the models fitted and some quantiles are given in the table below.

Model	Residual Deviance	Df
helmet*collision+outcome	25.643	5
helmet*collision+outcome*helmet	20.751	4
helmet*collision+outcome*collision	8.371	3
helmet*collision+outcome*helmet+outcome*collision	2.318	2

```
> qchisq(0.95,1)
[1] 3.841459
> qchisq(0.95,2)
[1] 5.991465
```

Specify the terms in the model with linear predictor given by `helmet*collision+outcome` and hence explain why the degrees of freedom is 5. *(2 marks)*

- (iv) Specify an algebraic form for the linear predictor `helmet*collision+outcome*helmet` *(2 marks)*
- (v) With reference to the residual deviances in the table above, describe the most suitable model for the data. Discuss whether the model you select based on the residual deviances is consistent with your observations from part (ii). *(7 marks)*
- (vi) Calculate the estimated expected number of serious injuries for cyclists wearing helmets in collision with a car for the model with linear predictor `helmet*collision+outcome*collision`. *(4 marks)*

**End of Question Paper**