

**BSc, MSci and MSc EXAMINATIONS (MATHEMATICS)**

**May-June 2018**

This paper is also taken for the relevant examination for the Associateship of the Royal College of Science

**Statistical Modelling II**

Date: Friday, 25 May 2018

Time: 2:00 PM - 4:30 PM

Time Allowed: 2.5 hours

**This paper has 5 questions.**

Candidates should use ONE main answer book.

Supplementary books may only be used after the relevant main book(s) are full.

All required additional material will be provided.

- DO NOT OPEN THIS PAPER UNTIL THE INVIGILATOR TELLS YOU TO.
- Affix one of the labels provided to each answer book that you use, but DO NOT USE THE LABEL WITH YOUR NAME ON IT.
- Each question carries equal weight.
- Calculators may not be used.

1. For the analysis of a random sample of data points  $\mathbf{Y} = (y_1, \dots, y_n)^T$ , a Normal linear model is proposed

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(0, \sigma^2 \mathbf{I}_n),$$

where  $\mathbf{X}$  is an  $n \times p$  matrix ( $p < n$ ) of covariates, including an intercept, and  $\mathbf{X}$  has full rank.

- Write down a system of linear equations satisfied by the maximum likelihood estimator  $\hat{\boldsymbol{\beta}}$ .
- Give an expression for the matrix  $\mathbf{P}$  such that the fitted values  $\hat{\mathbf{y}}$  can be written as  $\mathbf{P}\mathbf{y}$ . Show that  $\mathbf{P}$  is a projection matrix.
- Show that  $\mathbf{I} - \mathbf{P}$  is also a projection matrix.
- Show that the vector of residuals,  $\mathbf{e}$ , can be written as  $(\mathbf{I}_n - \mathbf{P})\mathbf{y}$ , and find the expectation and variance-covariance matrix of  $\mathbf{e}$ .
- Stating clearly any properties of projection matrices that you use, explain why

$$\frac{\mathbf{e}^T \mathbf{e}}{n - p}$$

is an unbiased estimator of  $\sigma^2$ .

In an experimental study, the model above is fitted to some data. Figure 1 shows leverage against standardized residuals for this model.

- State, with brief justification, which of the two labelled points A and B in Figure 1:
  - has a residual with a higher variance.
  - would result in a greater change in the estimates of the model coefficients if removed.
- State one property of the unlabelled residuals in Figure 1 that indicates a poor model fit. Suggest another plot that could be used to confirm this.

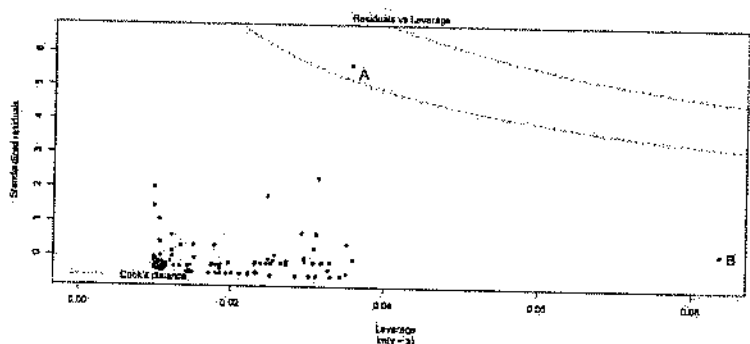


Figure 1: Standardized residuals plotted against leverage for the Normal linear model.

2. (a) Explain what is meant by

- (i) A generalized linear model.
- (ii) The canonical link function for a generalized linear model.
- (iii) A saturated model.

(b) Explain why a saturated model would not usually be fit to data.

The random variable  $Y$  has a density function written in exponential family form as follows

$$f(y; \theta, \phi) = \exp \left( \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \right).$$

(c) Show that  $E(Y) = b'(\theta)$ .

(d) Determine an expression for  $\text{Var}(Y)$  in terms of  $\theta$  and  $\phi$ .

(e) The inverse Gaussian distribution has density function

$$f(y; \lambda, \mu) = \left[ \frac{\lambda}{2\pi y^3} \right]^{\frac{1}{2}} \exp \left\{ -\frac{\lambda(y - \mu)^2}{2\mu^2 y} \right\}, \quad \lambda, \mu, y > 0.$$

Show that the inverse Gaussian density can be written in exponential family form, stating clearly the forms of  $\theta$ ,  $b(\theta)$  and  $\phi$  in terms of  $\lambda$  and  $\mu$ .

(f) State the canonical link for the inverse Gaussian distribution.

(g) Derive the form of the deviance for this model, in terms of  $\hat{\mu}$ .

(h) Suppose that an inverse Gaussian GLM is fitted to a large random sample of size  $n$ , with linear predictor  $\eta = X\beta$ , where  $\beta$  is of length  $p$ . Suggest a test statistic that can be used to carry out a test of the hypothesis that

$$\beta_{q+1} = \dots = \beta_p = 0,$$

and state the approximate distribution it follows under the null hypothesis. (Relevant asymptotic results may be assumed to hold.)

3. A Normal linear mixed model is written as

$$Y = X\beta + Z\nu + \epsilon,$$

where  $\nu \sim N(0, \sigma_\nu^2 I_m)$  is a random vector of length  $m$ ,  $\epsilon \sim N(0, \sigma_\epsilon^2 I_n)$  is a vector of length  $n$ , independent of  $\nu$ ,  $X$  is the design matrix for the fixed effects, and  $Z$  is the model matrix for the random effects. Assume that the design matrix has full rank,  $p$ .

- Determine the distribution of the random vector  $Y$ .
- Calculate the intra-class correlation coefficient for this model.
- Show that left multiplying  $Y$  by the matrix  $L^T = I_n - X(X^T X)^{-1} X^T$  gives a quantity that does not depend on the fixed effect coefficients  $\beta$ , and explain the relevance of this construction to the estimation of  $\sigma_\nu^2$ .
- If a different matrix  $P$ , such that  $\text{rank}(P) = n - p$  and  $P^T X = 0$ , were instead applied to  $Y$ , what would be the effect on the estimation?

The code below is used to test the null hypothesis that the random effects variance  $\sigma_\nu^2 = 0$ .

```
> mylm <- lm(response~1,data=dat)
> mylme2 <- lmer(response~1+(1|group),data=dat,REML=FALSE)
> d<-as.numeric(2*(logLik(mylme2)-logLik(mylm)))
> d
[1] 2.54379

> ds <- numeric(1000)
> for(i in 1:1000){
+   y <- unlist(simulate(mylm))
+   nullmod <- lm(y~1)
+   altmod <- lmer(y~1+(1|group),data=dat,REML=FALSE)
+   ds[i] <- as.numeric(2*(logLik(altmod)-logLik(nullmod)))
+ }
> mean(ds < 0.00001)
[1] 0.81
mean(ds>d)
[1] 0.01
```

- Explain briefly what the code does. (Your response should address the statistical purpose of the code, rather than simply describing it line-by-line.)
- Use the output given to explain how the usual asymptotic result for the generalized likelihood ratio test can be seen to be invalid here. State the assumption that is violated.
- State with reasons which of the models is to be preferred.
- Describe how the sampling uncertainty in the p-value above could be estimated.

4. The R output below concerns a binomial GLM,  $Y_i \sim \text{BINOMIAL}(n_i, \pi_i)$  for  $i = 1, \dots, n$ , where  $\pi$  is related to covariates  $X$  and  $Z$ .  $Y_i$  is the number of people who suffer from a disease. Throughout,  $n_i$  may be assumed to be large.

Call:

```
glm(formula = y ~ x + z, family = binomial)
```

Deviance Residuals:

|  | Min     | 1Q      | Median | 3Q     | Max    |
|--|---------|---------|--------|--------|--------|
|  | -5.2999 | -1.6093 | 0.0299 | 1.8933 | 4.8947 |

Coefficients:

|             | Estimate | Std. Error | z value | Pr(> z )     |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 0.23923  | 0.06404    | 3.736   | 0.000187 *** |
| x           | 0.38472  | 0.08242    | 4.668   | 3.05e-06 *** |
| z           | 0.44893  | 0.05360    | 8.376   | < 2e-16 ***  |

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 643.64 on 99 degrees of freedom  
Residual deviance: 541.87 on 97 degrees of freedom  
AIC: 1016.4

Number of Fisher Scoring iterations: 4

- State the link function used for this model, and one desirable property that results from this choice of link.
- Explain the effect of a one unit change in  $x$  on the odds of having the disease.
- Stating any distributional assumptions necessary, and taking the fitted model to be adequate, give an approximate 95% confidence interval for the coefficient of  $x$  in the model. You need not simplify your answer.
- Use appropriate values from the output to comment on the goodness of fit. Justify briefly any approximations you use.
- Explain what is meant by the AIC, and how it could be used to distinguish between the model here and a model that excluded the variable  $Z$ , with AIC 1015.2.
- State whether (and if so, how) the estimates of the parameters, and their standard errors, would change if a plug-in estimator of the dispersion were used. State the distribution that would be used to compute the probabilities in the fourth column in this case.

[CONTINUED]

- (g) Clustering is a common source of overdispersion. Suppose the sample is modelled as consisting of  $k$  clusters, each of size  $\frac{m}{k}$ , with

$$Y_i = \sum_{j=1}^{\frac{m}{k}} Z_{ij},$$

where the variables  $Z_{ij}|p_{ij} \sim \text{BINOMIAL}(k, p_{ij})$ , independently, and the probabilities  $p_{ij}$  are such that  $E(p_{ij}) = \pi_i$  and  $\text{Var}(p_{ij}) = \tau^2 \pi_i(1 - \pi_i)$ .

Show that in this case,  $E(Y_i) = m\pi_i$  and  $\text{Var}(Y_i) = \phi V(\pi_i)$ , where  $V$  is the variance function for the generalized linear model and  $\phi$  depends on  $\tau$  and  $k$  in a way that you should determine.

5. (a) Briefly explain why feature selection is desirable when fitting a linear model with a large number of predictors.
- (b) Explain why  $R^2$  is not a suitable goodness of fit measure for selecting amongst models with different numbers of predictors.
- (c) Explain why ridge regression is only applied after standardizing predictor variables.
- (d) For the model

$$Y = X\beta + \epsilon,$$

where  $\epsilon \sim N(0, \sigma^2 I_n)$ , write down the log likelihood function and show that if independent Normal priors with zero mean and variance  $\tau^2$  are assumed for the parameters, the estimates that maximize the resulting posterior distribution are identical to the parameter estimates from ridge regression.

- (e) Show that the ridge regression estimators  $\hat{\beta}_r$  satisfy  $\hat{\beta}_r = (X^T X + \lambda I_n)^{-1} X^T y$ , where  $\lambda$  is the tuning parameter.
- (f) Consider a lab that is trying to design a clinical screen for diabetes. A large number of blood markers are included as predictors in a Normal linear model in which the response variable is blood glucose level. The clinicians' understanding is that the level of glucose is strongly determined by a small number of markers, but they have no knowledge of which markers are likely to be involved. Explain whether ridge regression or the lasso is likely to perform better in this context.
- (g) Consider the problem of selecting from amongst two candidate models. Model  $A$  has 3 parameters and model  $B$  has 4 parameters. Two statisticians propose different approaches. The first statistician chooses the simpler model unless the test statistic

$$2(l(B) - l(A))$$

is greater than the 95th percentage point of the distribution of a  $\chi^2(1)$  variable. The second statistician always chooses the model with the smaller AIC. Determine which statistician has the lower type I error rate.

Note that  $\Pr(\chi^2(1) > 3.84) \approx 0.05$ .

Imperial College  
London

Course: M3S2/M4S2  
Setter:  
Checker:  
Editor:  
External:  
Date: April 24, 2018

MSc EXAMINATIONS (MATHEMATICS)  
May 2018

M3S2/M4S2

Statistical Modelling II [SOLUTIONS]

|                    |                     |                    |
|--------------------|---------------------|--------------------|
| Setter's signature | Checker's signature | Editor's signature |
| .....              | .....               | .....              |



1. (a) [1 mark]  $\hat{\beta}$  satisfies

$$X^T X \hat{\beta} = X^T y,$$

or, since  $X$  has full rank,

$$\hat{\beta} = (X^T X)^{-1} X^T y.$$

[Seen]

- (b) [4 marks] The fitted values are given by  $\hat{y} = X\hat{\beta} = X(X^T X)^{-1} X^T y$ , so the matrix  $P = X(X^T X)^{-1} X^T$ .

We check this is a projection matrix by checking it is both symmetric and idempotent.

It is symmetric, since  $P^T = (X(X^T X)^{-1} X^T)^T = (X^T)^T ((X^T X)^{-1})^T X^T = X(X^T X)^{-1} X^T$ , using the facts that  $(AB)^T = B^T A^T$  and  $(A^{-1})^T = (A^T)^{-1}$ .

It is idempotent since

$$P^2 = X(X^T X)^{-1} X^T X(X^T X)^{-1} X^T = X(X^T X)^{-1} I X^T = X(X^T X)^{-1} X^T = P.$$

[Seen]

- (c) [2 marks] Let  $Q = I - P$ . We show it is symmetric and idempotent.  $(Q^T)_{ii} = 1 - P_{ii} = Q_{ii}$  and  $(Q^T)_{ij} = Q_{ji} = -P_{ji} = -P_{ij} = Q_{ij}$  as  $P = P^T$ . Thus  $Q = Q^T$ . For idempotence,  $Q^2 = QQ = Q^T Q$  as  $Q$  is symmetric by above. Then  $Q^T Q = (I_n - P)^T (I_n - P) = I_n - 2P + \underbrace{P^T P}_{=P} = I_n - P = Q$ . [Seen]

- (d) [3 marks] The residuals are given by  $e = y - \hat{y} = y - Py = (I_n - P)y$ .

To compute the mean of the residuals,

$$E(e) = E((I_n - P)y) = (I_n - P)E(y) = X\beta - PX\beta = X\beta - X(X^T X)^{-1} X^T X\beta = X\beta - X I_p \beta = 0.$$

To compute the variance of the residuals,

$$\begin{aligned} \text{Cov}(e) &= \text{Cov}((I_n - P)y) = (I_n - P)\text{Cov}(y)(I_n - P)^T = (I_n - P)\sigma^2 I_n (I_n - P)^T \\ &= \sigma^2 (I_n - P)(I_n - P)^T = \sigma^2 (I_n - P), \end{aligned}$$

using that  $I_n - P$  is both symmetric and idempotent. [Seen]

(e) [4 marks]

$$\begin{aligned}
 E(e^T e) &= E(Y^T (I_n - P)^T (I_n - P) Y) = E(Y^T \underbrace{(I_n - P)}_{=Q} Y) \\
 &= E\left(\sum_i \sum_j Y_i Y_j Q_{ij}\right) = E(\text{trace}(Q Y Y^T)) \\
 &= \text{trace}(Q E(Y Y^T)) = \text{trace}(Q \{ \text{Cov}(Y) + E(Y) E(Y)^T \}) \\
 &= \text{trace}(Q \text{Cov}(Y)) + \text{trace}(Q E(Y) E(Y)^T) \\
 &= \text{trace}(Q \sigma^2) + \underbrace{\text{trace}(Q X \beta (X \beta)^T)}_{=0} \\
 &= \sigma^2 \text{trace}(I_n - P) \\
 &= \sigma^2(n - p), \text{ as } \text{trace}(P) = \text{rank}(P) = \text{rank}(X) = p.
 \end{aligned}$$

[Seen]

- (f) [2 marks] Standardized residuals are asymmetrically distributed about zero, but they should be approximately standard normal, and so symmetric. A normal Q-Q plot would confirm this. [Unseen]
- (g) (i) [2 marks]  $\text{Var}(e_i) = (1 - P_{ii})\sigma^2$ , where  $P_{ii}$  is the leverage. High leverage corresponds to small variance, so point *B* has smallest variance. [Unseen]
- (ii) [2 marks] Cook's distance gives change in the fit as a result of removing a particular observation. Since observation *A* has a higher Cook's distance, removing it will cause the larger difference in fit. [Seen, Similar]

2. (a) (i) [3 marks] A generalized linear model is defined by

- The **random component** specifies the probability distribution of the response variables. Specifically, the components of  $y$  have pdf or pmf from an exponential family of distributions, with  $E(Y) = \mu$ .
- The **systematic component** specifies a linear predictor  $\eta = X\beta$  as a function of the covariates and the unknown parameters.
- The **link function**  $g$  may be any monotonic differentiable function. The link function provides a functional relationship between the systematic component and the expectation of the response in the random component; namely  $\eta = g(\mu)$ . [Seen]

(ii) [1 mark] If the density of  $Y$  is written in exponential family form,

$$f(y; \theta, \phi) = \exp \left( \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \right),$$

then the canonical link is defined by the relation

$$\theta = \eta. \text{ [Seen]}$$

(iii) [1 mark] A model is said to be saturated if there as many independent parameters as observations. [Seen]

(b) [1 mark] A saturated model involves estimating  $n$  parameters from  $n$  data points. Typically, this means that the model can fit the data perfectly. The act of fitting a statistical model presupposes that the process under consideration is subject to random fluctuations, which we cannot expect to model. A saturated model is therefore fit to noise as well as systematic variation of interest, and inferences and predictions made from it will, on average, be subject to larger errors than a model with fewer parameters. [Seen]

(c) [3 marks]

Note first, that since  $f$  is a density function,

$$\int f(y; \theta, \phi) dy = 1.$$

Hence,

$$\frac{\partial}{\partial \theta} \int f(y; \theta, \phi) dy = \frac{\partial}{\partial \theta} 1 = 0.$$

Differentiating the density function, we see that

$$\frac{\partial}{\partial \theta} f(y; \theta, \phi) = \frac{y - b'(\theta)}{a(\phi)} f(y; \theta, \phi).$$

Assuming sufficient regularity that integration and partial differentiation with respect to  $\theta$  can be interchanged, this gives

$$\int \left( \frac{y - b'(\theta)}{a(\phi)} \right) f(y; \theta, \phi) dy = 0,$$

which simplifies to

$$\frac{1}{a(\phi)} (E(Y) - b'(\theta)) = 0,$$

by definition of the expected value. Thus

$$\mu \equiv E(Y) = b'(\theta). \quad (1)$$

[Seen]

(d) [2 marks] Differentiating again with respect to  $\theta$ ,

$$\frac{\partial^2}{\partial \theta^2} f(y; \theta, \phi) = -\frac{b''(\theta)}{a(\phi)} f(y; \theta, \phi) + \left( \frac{y - b'(\theta)}{a(\phi)} \right)^2 f(y; \theta, \phi).$$

Integrating over  $y$  then gives

$$\begin{aligned} 0 &= \int \frac{\partial^2}{\partial \theta^2} f(y; \theta, \phi) dy = -\frac{b''(\theta)}{a(\phi)} + \int \left( \frac{y - b'(\theta)}{a(\phi)} \right)^2 f(y; \theta, \phi) dy \\ &= -\frac{b''(\theta)}{a(\phi)} + \frac{\text{Var}(Y)}{a(\phi)^2}. \end{aligned}$$

Rearranging yields

$$\text{Var}(Y) = b''(\theta) a(\phi).$$

[Seen]

- (e) [4 marks] The density function is

$$f(y; \lambda, \mu) = \left[ \frac{\lambda}{2\pi y^3} \right]^{\frac{1}{2}} \exp \left\{ \frac{-\lambda(y - \mu)^2}{2\mu^2 y} \right\}, \quad \lambda, \mu, y > 0.$$

Putting this all inside an exponential gives

$$\begin{aligned} f(y; \lambda, \mu) &= \exp \left\{ \frac{1}{2}(\log \lambda - \log 2\pi) - \frac{3}{2} \log y - \frac{\lambda(y - \mu)^2}{2\mu^2 y} \right\} \\ &= \exp \left\{ \frac{1}{2}(\log \lambda - \log 2\pi) - \frac{3}{2} \log y - \frac{\lambda(y^2 - 2\mu y + \mu^2)}{2\mu^2 y} \right\} \\ &= \exp \left\{ \frac{1}{2}(\log \lambda - \log 2\pi) - \frac{3}{2} \log y - \frac{\lambda y}{2\mu^2} + \frac{\lambda}{\mu} - \frac{\lambda}{2y} \right\} \end{aligned}$$

By considering the term in  $\frac{1}{y}$ , this is clearly an exponential family with  $\phi \propto \frac{2}{\lambda}$ . Taking the constant of proportionality as 1 (we only have uniqueness up to a constant), using the form of the exponential family gives:

$$\begin{aligned} \theta &= -\frac{1}{\mu^2} \\ b(\theta) &= -\frac{2}{\mu} = -2\sqrt{-\theta}. \end{aligned}$$

[Unseen]

- (f) [1 mark] For the canonical link, we need  $\eta = \theta = -\frac{1}{\mu^2}$ . [Seen Similar]  
 (g) [2 marks]

$$D = 2\phi(l(\hat{y}, y) - l(\hat{\mu}, y)).$$

The log likelihood for a single observation is

$$\frac{1}{2}(\log \lambda - \log 2\pi) - \frac{3}{2} \log y - \frac{\lambda(y - \mu)^2}{2\mu^2 y}$$

Note first that terms in the log likelihood involving  $y$  and  $\lambda$  alone cancel. Moreover, the remaining term for the saturated model is clearly zero, since it is a function of  $(y_i - \mu_i)$ , and for the saturated model  $\hat{\mu}_i = y_i$ . Hence, since  $\phi = \frac{\lambda}{2}$  here, the deviance is

$$D = \sum_{i=1}^n \frac{(y_i - \hat{\mu}_i)^2}{\hat{\mu}_i^2 y_i}.$$

[Seen Similar]

- (h) [2 marks] Under the null hypothesis we have (approximately, for large samples) that the scaled deviances  $D^* = D/\phi$  are  $\chi^2$ :

$$D_1^* \sim \chi_{n-p}^2 \quad \text{and} \quad D_0^* - D_1^* \sim \chi_{p-q}^2.$$

If we consider  $D_1^*$  and  $D_0^* - D_1^*$  as asymptotically independent, then

$$\frac{(D_0^* - D_1^*)/(p-q)}{D_1^*/(n-p)} \sim F_{(p-q), (n-p)}.$$

We can multiply top and bottom by  $\phi$  to get a test statistic based on the deviance:

$$\frac{(D_0 - D_1)/(p-q)}{D_1/(n-p)} \sim F_{(p-q), (n-p)}.$$

[Seen Similar]

3. (a) [3 marks] As a linear combination of multivariate normal variables,  $Y$  must be multivariate normal. It suffices to compute its mean vector and variance-covariance matrix.  
For the mean

$$\begin{aligned} E(Y) &= E(X\beta + Z\nu + \epsilon) \\ &= X\beta + Z \underbrace{E(\nu)}_{=0} + \underbrace{E(\epsilon)}_{=0} \\ &= X\beta \end{aligned}$$

And for the variance-covariance matrix,

$$\begin{aligned} \text{Cov}(Y) &= \text{Cov}(X\beta + Z\nu + \epsilon) \\ &= Z\text{Cov}(\nu)Z^T + \text{Cov}(\epsilon) \quad (\text{since } \epsilon, \nu \text{ indep.}) \\ &= ZI_m\sigma_\nu^2Z^T + I_n\sigma_\epsilon^2 \\ &= \sigma_\epsilon^2 \left( ZI_m\frac{\sigma_\nu^2}{\sigma_\epsilon^2}Z^T + I_n \right) \\ &= \sigma_\epsilon^2 (I_n + Z\Psi Z^T + I_m), \end{aligned}$$

where  $\Psi = \frac{\sigma_\nu^2}{\sigma_\epsilon^2}I_m$ . [Seen]

- (b) [3 marks]

The correlation for observations in the same group is:

$$\text{corr}(Y_{1,j}, Y_{2,j}) = \frac{E[(Y_{1,j} - E(Y_{1,j}))(Y_{2,j} - E(Y_{2,j}))]}{\sqrt{\text{Var}(Y_{1,j})\text{Var}(Y_{2,j})}}$$

We have that  $\text{Var}(Y_{1,j}) = \text{Var}(Y_{2,j}) = \sigma_\epsilon^2 + \sigma_\nu^2$ , and thus the denominator is  $\sigma_\epsilon^2 + \sigma_\nu^2$ . Next, we have  $Y_{i,j} - E(Y_{i,j}) = \nu_j + \epsilon_{ij}$ , by noting that  $E(Y_{i,j}) = \mu$  and rearranging the model equation. Therefore the numerator is

$$\begin{aligned} E[(Y_{1,j} - E(Y_{1,j}))(Y_{2,j} - E(Y_{2,j}))] &= E[(\nu_j + \epsilon_{1,j})(\nu_j + \epsilon_{2,j})] \\ &= E(\nu_j^2) + E(\nu_j(\epsilon_{1,j} + \epsilon_{2,j})) + E(\epsilon_{1,j}\epsilon_{2,j}) \\ &= E(\nu_j^2) = \sigma_\nu^2, \end{aligned}$$

where we used the independence between the  $\nu_j$  and  $\epsilon_{ij}$ . Plugging in these results gives the intraclass correlation coefficient as

$$\rho = \frac{\sigma_\nu^2}{\sigma_\nu^2 + \sigma_\epsilon^2}.$$

[Seen]

- (c) [3 marks] Let  $L^T = I_n - X(X^T X)^{-1} X^T$ , then

$$\begin{aligned} L^T Y &= L^T X\beta + L^T Z\nu + L^T \epsilon \\ (I_n - X(X^T X)^{-1} X^T) Y &= X\beta - X(X^T X)^{-1} (X^T X)\beta + L^T Z\nu + L^T \epsilon \\ Y - X\hat{\beta} &= L^T Z\nu + L^T \epsilon, \end{aligned}$$

so we see that the term in  $X\beta$  cancels. Now let  $B = L^T Y$ , then

$$B \sim N(0, \sigma_\epsilon^2 L^T V L),$$

where  $V = I_n + \frac{\sigma_\nu^2}{\sigma_\epsilon^2} Z Z^T$ .

Then the joint pdf of  $B$  (note that strictly this is an improper density on the  $n - p$  dimensional subspace orthogonal to the columns of  $X$ ) is

$$\frac{1}{(2\pi)^{(n-p)/2} |\sigma_\epsilon^2 L^T V L|^{1/2}} \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} b^T (L^T V L)^{-1} b \right\}$$

so that the log-likelihood of the variance components vector  $\tau = (\sigma_\epsilon^2, \sigma_\nu^2)$  is

$$\ell_R(\tau; b) = -\frac{(n-p)}{2} \log(2\pi) - \frac{1}{2} \log |\sigma_\epsilon^2 L^T V L| - \frac{1}{2\sigma_\epsilon^2} b^T (L^T V L)^{-1} b.$$

The REML (restricted maximum likelihood) estimator of  $\tau$  is defined as the maximiser of  $\ell_R$ . This maximiser can now be found by solving the REML equation

$$\frac{\partial \ell_R}{\partial \tau} = 0.$$

[Seen]

- (d) [2 marks] To show that the REML estimators are invariant to the choice of  $L$ , we show that  $\ell(\tau; y)$  is invariant to the choice of  $L$ . To show this, suppose that we have another  $L_2$  such that  $\text{rank}(L_2) = n - p$  and  $L_2^T X = 0$ . Then, there exists a non-singular matrix  $M$  such that  $L^T = M L_2^T$ . Further

$$|L^T V L| = |M L_2^T V L_2 M^T| = |M|^2 |L_2^T V L_2|$$

and

$$\begin{aligned} b^T (L^T V L)^{-1} b &= y^T L (L^T V L)^{-1} L^T y \\ &= y^T L_2 M^T (M L_2^T V L_2 M^T)^{-1} M L_2^T y \\ &= y^T L_2 M^T (M^T)^{-1} (L_2^T V L_2)^{-1} M^{-1} M L_2^T y \\ &= y^T L_2 (L_2^T V L_2)^{-1} L_2^T y. \end{aligned}$$

Substituting these into the restricted log-likelihood, we get

$$\ell(\tau; y) = -\frac{(n-p)}{2} \log(2\pi) - \log |M| - \frac{1}{2} \log |\sigma_\epsilon^2 L_2^T V L_2| - \frac{1}{2\sigma_\epsilon^2} y^T L_2 (L_2^T V L_2)^{-1} L_2^T y.$$

[Seen Similar]

- (e) [2 marks] The code is a parametric bootstrap routine to estimate the probability, under the null hypothesis that  $\sigma_\nu^2 = 0$ , that the log likelihood ratio test statistic is as large as the one observed. It achieves this by generating 1000 independent samples from the null model, and computing the observed value of the test statistic. [Seen Similar]



- (f) [3 marks] The output shows the distribution of the test statistic is concentrated near zero. This is not consistent with the usual asymptotic result that the GLRT statistic has an approximate  $\chi^2$  distribution. To derive the asymptotic  $\chi^2$  distribution, we assume the parameter lies in the interior of the parameter space. But since a variance cannot be negative,  $\sigma_v^2 = 0$  is on the boundary, and so the assumption does not hold. [Seen Similar]
- (g) [2 marks] The improvement in likelihood is larger than would be expected due to chance, so prefer the model with  $\sigma_v^2 > 0$ . [Seen Similar]
- (h) [2 marks] Standard error for the binomial proportion given by

$$\sqrt{\frac{p(1-p)}{n}} = \sqrt{0.01 \times 0.99/1000}$$

(negligible here). [Unseen]

4. (a) [3 marks] Since the link function is not explicitly specified in the output, it must be the canonical link. For the binomial family this is the logit link,  $\eta = \log \frac{\pi}{1-\pi}$ .

With this choice of link function, the fitted values are normalized to lie in an appropriate range for probabilities. Also, the effects of predictors are interpretable on the odds scale. Moreover, since the link is canonical, the observed information is the same as the expected information. [Seen]

- (b) [2 marks] With the logit link,

$$\log \left( \frac{\pi_i}{1 - \pi_i} \right) = \sum_{j=1}^p x_{ij} \beta_j,$$

so on the odds scale, the predictors act multiplicatively. Hence a one unit increase in  $x$  leads to the odds of having the disease being multiplied by  $\exp(\hat{\beta}_x) = \exp(0.38)$ . [Seen Similar]

- (c) [2 marks] Assuming that the estimator for the parameter is normally distributed, a 95% confidence interval is given by  $0.38 \pm 1.96 \times 0.08$ . [Seen Similar]
- (d) [3 marks] When  $m$  (number of binomial samples per observation) is large, the scaled deviance should be roughly  $\chi^2(n-p)$ , where  $n$  is the number of observations and  $p$  is the number of parameters. But this means that  $\frac{D}{n-p}$  should estimate  $\phi$ , the dispersion parameter. But this is just 1 for the binomial model, which does not match up well with the value 542/97 from the table. [Seen Similar]
- (e) [2 marks] AIC is given by  $-2l(\hat{\beta}) + 2p$ . It is an attempt to penalize models that contain too many parameters - we select the model with the lower AIC. Therefore, we prefer the model without  $Z$ . [Seen Similar]
- (f) [3 marks] The estimates themselves would not change - these are independent of  $\phi$ . However, the covariance matrix of  $\hat{\beta}$  is  $\phi(X^T W X)^{-1}$ , where  $W$  is the matrix of weights from iterated weighted least squares. The standard errors are the square roots of the diagonal entries of the covariance matrix, so the standard errors would be multiplied by a factor of  $\sqrt{\phi}$ . The  $t$  distribution would be used, since the dispersion parameter is estimated from data. [Seen Similar]

(g) [5 marks] First note that

$$E(Y_i) = E_{\pi_i} E(Y_i | \pi_i) = E \left( \sum_{j=1}^{\frac{m}{k}} k p_{ij} \right) = \frac{m}{k} \times k \pi_i = m \pi_i.$$

Use law of total variance

$$\begin{aligned} \text{Var}(Y_i) &= E \left( \sum_{j=1}^{\frac{m}{k}} \text{Var}(Z_{ij} | p_{ij}) \right) + \text{Var} \left( \sum_{j=1}^{\frac{m}{k}} E(Z_{ij} | p_{ij}) \right) \\ &= E \left( \sum_{j=1}^{\frac{m}{k}} k p_{ij} (1 - p_{ij}) | p_{ij} \right) + \text{Var} \left( \sum_{j=1}^{\frac{m}{k}} E(Z_{ij} | p_{ij}) \right) \end{aligned}$$

The first term is

$$\begin{aligned} E \left( \sum_{j=1}^{\frac{m}{k}} \text{Var}(Z_{ij} | p_{ij}) \right) &= E \left( \sum_{j=1}^{\frac{m}{k}} k p_{ij} (1 - p_{ij}) \right) = E \left( \sum_{j=1}^{\frac{m}{k}} k p_{ij} \right) - E \left( \sum_{j=1}^{\frac{m}{k}} k p_{ij}^2 \right) \\ &= m \pi_i - \sum_{j=1}^{\frac{m}{k}} k \text{Var}(p_{ij}) - k E(p_{ij})^2 \\ &= m \pi_i - m \tau \pi_i (1 - \pi_i) - m \pi_i^2 = m \pi_i (1 - \pi_i) - m \tau \pi_i (1 - \pi_i). \end{aligned}$$

The second term is

$$\text{Var} \left( \sum_{j=1}^{\frac{m}{k}} E(Z_{ij} | p_{ij}) \right) = \sum_{j=1}^{\frac{m}{k}} \text{Var}(k p_{ij}) = \frac{m}{k} \times k^2 \tau \pi_i (1 - \pi_i) = m k \tau \pi_i (1 - \pi_i).$$

Combining these two gives

$$\text{Var}(Y_i) = m \pi_i (1 - \pi_i) (1 + (k - 1) \tau^2).$$

so define  $\phi = 1 + (k - 1) \tau^2$ . [Unseen]

5. (a) [2 marks] By choosing a subset of predictors, we reduce the tendency of least squares to overfitting. For a model with  $p \approx n$ , i.e. number of predictors comparable to number of observations, the parameter estimates will have high variance. *in text*
- (b) [3 marks] When a predictor is added, the residual sum of squares must decrease - even an irrelevant predictor explains some variability. Since  $R^2 = 1 - \frac{RSS}{TSS}$ ,  $R^2$  necessarily increases when a new predictor is added. *[Seen]*
- (c) [3 marks] As stated in the text, ridge regression is not invariant under scale transformation - changing the scale of the predictors would lead to very different parameter estimates, as it is the total  $l^2$  norm of the parameter estimates that it penalized. *in text*
- (d) [4 marks] The log likelihood function is given by

$$-\frac{1}{2\sigma^2}(Y - X\beta)^T(Y - X\beta).$$

Independent normal priors on  $\beta_1, \dots, \beta_p$  give rise to a prior on the log scale

$$-\frac{1}{2\tau^2}\beta^T\beta.$$

Bayes' theorem then gives the log posterior as

$$-\frac{1}{2\sigma^2}(Y - X\beta)^T(Y - X\beta) - \frac{1}{2\tau^2}\beta^T\beta + K,$$

where  $K$  is the log of a normalizing constant that does not depend on  $\beta$ .

Maximizing this function is clearly equivalent to minimizing the ridge regression criterion

$$(Y - X\beta)^T(Y - X\beta) + \lambda\beta^T\beta,$$

where  $\lambda = \frac{\sigma^2}{\tau^2}$ .

- (e) [3 marks] Let  $S(\beta) = (Y - X\beta)^T(Y - X\beta) + \lambda\beta^T\beta$ . Differentiating with respect to  $\beta$  gives

$$-2X^TY + 2X^TX\beta + 2\lambda\beta.$$

Setting this expression equal to zero for a critical point gives

$$(X^TX + \lambda I_n)\beta = X^TY.$$

*[Unseen] an exercise later in the text, not in the extract given.*

- (f) [2 marks] The lasso is equivalent to a prior under which most predictors have essentially zero effect, but a small number may have large effects. Ridge regression on the other hand supports a large number of variables contributing in moderate amounts. Lasso therefore better respects the clinicians' intuitions about this data. *[Unseen]*
- (g) [3 marks]  $AIC = -2l(\hat{\beta}) + 2p$ , so selecting the model with the smaller AIC in this case (where B has one more parameter than A) corresponds to rejecting the null hypothesis that model A is correct whenever  $2(l(B) - l(A)) > 2$ . This means that the first statistician will have the lower type I error rate, since they are using a more extreme point of the  $\chi^2(1)$  distribution. *[Unseen]*