

i Cover Page

Department of Mathematical Sciences

Examination paper for TMA4315 - Generalized Linear Models

Examination date: 5th December 2023

Examination time (from-to): 15:00 - 19:00

Permitted examination support material: C

- Tabeller og formler i statistikk (Tapir forlag, Fagbokforlaget),
- one yellow A4 sheet with your own handwritten notes (stamped by the Department of Mathematical Sciences),
- specified calculator

Academic contact during examination: Bob O'Hara
Phone: 91554416

Academic contact present at the exam location: NO

OTHER INFORMATION

Get an overview of the question set before you start answering the questions.

Read the questions carefully and make your own assumptions. If a question is unclear/vague, make your own assumptions and specify them in your answer. The academic person is only contacted in case of errors or insufficiencies in the question set. Address an invigilator if you suspect errors or insufficiencies. Write down the question in advance.

Hand drawings: For question 12 you are meant to answer on handwritten sheets. Other questions should be answered directly in Inspera. At the bottom of the question you will find a seven-digit code. Fill in this code in the top left corner of the sheets you wish to submit. We recommend that you do this during the exam. If you require access to the codes after the examination time ends, click "Show submission".

Weighting: The maximum achievable score for each question is given with the question.

Notifications: If there is a need to send a message to the candidates during the exam (e.g. if there is an error in the question set), this will be done by sending a notification in Inspera. A dialogue box will appear. You can re-read the notification by clicking the bell icon in the top right-hand corner of the screen.

Withdrawing from the exam: If you become ill or wish to submit a blank test/withdraw from the exam for another reason, go to the menu in the top right-hand corner and click "Submit blank". This cannot be undone, even if the test is still open.

Access to your answers: After the exam, you can find your answers in the archive in Inspera. Be aware that it may take a working day until any hand-written material is available in the archive.

A few years ago a group of researchers investigated the numbers of deaths caused by hurricanes, and concluded that hurricanes with female sounding names were deadlier.

Of course, one has to control for the strength of the hurricane, so a variable called NDAM (=“normalised damage”) was used.

We can assume that the number of deaths follows a Poisson distribution, and fit a GLM.

1 Parts of a GLM

List and explain the 3 parts of a Generalised Linear Model (1-2 sentences for each element)

Fill in your answer here

Maximum marks: 6

The negative binomial distribution can be written in this form

$$f(y|r, p) = \binom{y+r-1}{y} (1-p)^r p^y$$

with $y=0,1,2,\dots$, $r>0$, $0<p<1$. We are interested in modelling p , and we can treat r as a constant.

We should write this as a member of the exponential family:

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

2 b()

What is $b(\theta)$?

Answer:

Maximum marks: 2

3 Canonical Link Function

What is the canonical link function?

Select one alternative:

- identity
- log
- logit
- inverse
- exponential
- something else
- probit
- cloglog

Maximum marks: 1

A few years ago a group of researchers investigated the numbers of deaths caused by hurricanes, and concluded that hurricanes with female sounding names were deadlier.

Of course, one has to control for the strength of the hurricane, so a variable called NDAM (=“normalised damage”) was used.

We can assume that the number of deaths follows a Poisson distribution, and fit a GLM.

When the model is fitted, we get the following output.

Call:

```
glm(formula = alldeaths ~ NDAM + Minpressure + Gender, family = "poisson",
     data = Himmicanes)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	2.88299	0.16187	17.810	< 2e-16	***
NDAM	0.32597	0.07708	4.229	2.35e-05	***
Minpressure	-0.43588	0.14002	-3.113	0.00185	**
GenderMale	-0.58112	0.29224	-1.989	0.04675	*

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for poisson family taken to be 28.03668)

Null deviance: 4031.9 on 91 degrees of freedom
 Residual deviance: 2467.2 on 88 degrees of freedom
 AIC: 2800.1

Number of Fisher Scoring iterations: 6

NDAM is the amount of damage done by the hurricane (adjusted for inflation): a stronger hurricane should do more damage. It is standardised to mean 0, variance 1.

Minpressure is the minimum pressure of the hurricane: a stronger hurricane has a lower pressure. It is standardised to mean 0, variance 1.

Gender is the gender of the name assigned to the hurricane: it is a factor with two levels (*Male* and *Female*)

4 Gender Effect

If two hurricanes had the same NDAM and minpressure, by how many times would the predicted deaths be higher if the gender was female?

Answer to 2 decimal places: .

Maximum marks: 1

5 Gender Effect Calculations

Show your working for this problem

Fill in your answer here

Maximum marks: 3

A few years ago a group of researchers investigated the numbers of deaths caused by hurricanes, and concluded that hurricanes with female sounding names were deadlier.

Of course, one has to control for the strength of the hurricane, so a variable called NDAM (=“normalised damage”) was used.

We can assume that the number of deaths follows a Poisson distribution, and fit a GLM.

When the model is fitted, we get the following output.

Call:

```
glm(formula = alldeaths ~ NDAM + Minpressure + Gender, family = "poisson",
     data = Himmicanes)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	2.88299	0.16187	17.810	< 2e-16	***
NDAM	0.32597	0.07708	4.229	2.35e-05	***
Minpressure	-0.43588	0.14002	-3.113	0.00185	**
GenderMale	-0.58112	0.29224	-1.989	0.04675	*

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for poisson family taken to be 28.03668)

Null deviance: 4031.9 on 91 degrees of freedom
 Residual deviance: 2467.2 on 88 degrees of freedom
 AIC: 2800.1

Number of Fisher Scoring iterations: 6

NDAM is the amount of damage done by the hurricane (adjusted for inflation): a stronger hurricane should do more damage. It is standardised to mean 0, variance 1.

Minpressure is the minimum pressure of the hurricane: a stronger hurricane has a lower pressure. It is standardised to mean 0, variance 1.

Gender is the gender of the name assigned to the hurricane: it is a factor with two levels (*Male* and *Female*)

6 Confidence Interval

Calculate the lower bound of a 95% confidence interval for this estimate.

Answer to 2 decimal places: .

Maximum marks: 1

7 Confidence Interval Working

Show your working for this problem

Fill in your answer here

Maximum marks: 3

A few years ago a group of researchers investigated the numbers of deaths caused by hurricanes, and concluded that hurricanes with female sounding names were deadlier.

Of course, one has to control for the strength of the hurricane, so a variable called NDAM (=“normalised damage”) was used.

We can assume that the number of deaths follows a Poisson distribution, and fit a GLM.

When the model is fitted, we get the following output.

Call:

```
glm(formula = alldeaths ~ NDAM + Minpressure + Gender, family = "poisson",
     data = Himmicanes)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	2.88299	0.16187	17.810	< 2e-16	***
NDAM	0.32597	0.07708	4.229	2.35e-05	***
Minpressure	-0.43588	0.14002	-3.113	0.00185	**
GenderMale	-0.58112	0.29224	-1.989	0.04675	*

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for poisson family taken to be 28.03668)

Null deviance: 4031.9 on 91 degrees of freedom
 Residual deviance: 2467.2 on 88 degrees of freedom
 AIC: 2800.1

Number of Fisher Scoring iterations: 6

NDAM is the amount of damage done by the hurricane (adjusted for inflation): a stronger hurricane should do more damage. It is standardised to mean 0, variance 1.

Minpressure is the minimum pressure of the hurricane: a stronger hurricane has a lower pressure. It is standardised to mean 0, variance 1.

Gender is the gender of the name assigned to the hurricane: it is a factor with two levels (*Male* and *Female*)

8 Overdispersion

The output states ‘(Dispersion parameter for poisson family taken to be 19.28544)’. The dispersion parameter would normally be around 1.

What does a value of 19.3 suggest about the data and the assumption that it follows a Poisson distribution? What might cause this?

Fill in your answer here

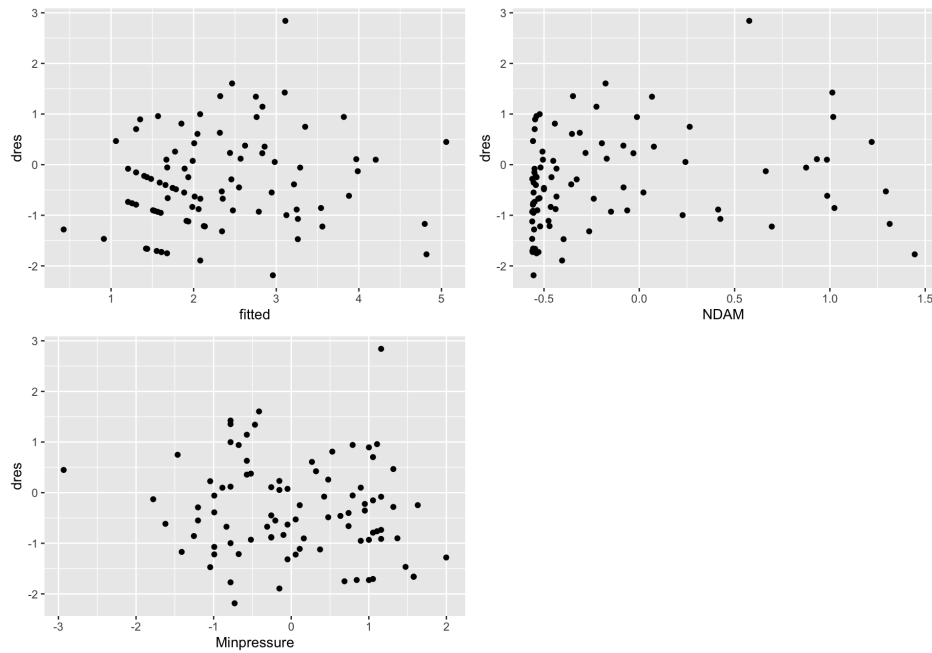
Maximum marks: 4

9 Overdispersion Calculation

How could this value be calculated?

Fill in your answer here

Maximum marks: 2



The data analysis was re-examined and residual plots suggested there was a problem with the model. Here are the plots, of the deviance residuals against the fitted values and continuous covariates.

10 How good is the model?

Describe (briefly!) what these plots tells you about the model fit.

Fill in your answer here

Maximum marks: 4

11 Model Comparison

After looking at this, a better model was suggested, by adding terms into the model, so that the old model was nested within the new one. The two models were compared, giving the following output

```
Analysis of Deviance Table
Model 1: alldeaths ~ ...
Model 2: alldeaths ~ ....
  Resid. Df Resid. Dev Df Deviance
1         88      2467.2
2         86      2034.9  2    432.28
```

Is there evidence the new model is better? Explain which statistics you used to reach this conclusion.

Fill in your answer here

Maximum marks: 4

12 LMMS

Note: Answers for this question must be written on handwritten sheets.

In a linear mixed model, the joint distribution of \mathbf{Y} and $\boldsymbol{\gamma}$ is:

$$\begin{pmatrix} \mathbf{Y} \\ \boldsymbol{\gamma} \end{pmatrix} \sim N\left(\begin{pmatrix} \mathbf{X}\boldsymbol{\beta} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{UGU}^T + \sigma^2\mathbf{I} & ? \\ ?^T & \mathbf{G} \end{pmatrix}\right)$$

We want to look at the off-diagonal part, i.e. $\text{Cov}(\mathbf{Y}, \boldsymbol{\gamma})$, for a model with a random intercept, i.e.

$$Y_{ij} = \beta_0 + \beta_1 x_{ij} + \gamma_{0i} + \varepsilon_{ij}$$

where $i = 1, \dots, N$, $j = 1, \dots, n_i$, $\gamma_{0i} \sim N(0, \tau_0^2)$ and $\varepsilon_{ij} \sim N(0, \sigma^2)$

- Show that when $i = k$, $\text{Cov}(y_{ij}, \gamma_{0k}) = E(\gamma_{0i}^2) = \text{Var}(\gamma_{0i})$.
- Derive the expression for when $i \neq k$.
- From this, explain what the matrix $?$ looks like (or draw it if that is easier!)

We can look at this model:

$$y_{ij} = x_{ij}\boldsymbol{\beta} + \gamma_{0i} + \gamma_{1i}x_{ij} + \varepsilon_{ij} \text{ with } \gamma_{1i} \sim N(0, \tau_1^2) \text{ and } \text{Corr}(\gamma_{0i}, \gamma_{1i}) = \rho,$$

having a single covariate and a random slope for the effect of that covariate that varies between groups.

- Derive $\text{Cov}(y_{ij}, y_{ik})$.

Maximum marks: 14

Biologists have collected data on the sizes of birds and the sizes of their eggs: an obvious expectation is that larger birds have larger eggs.

Species are grouped into genus (i.e. each genus contains more than one species), so we can look at whether the change between genera is greater than between species.

The following model was fitted

```
Egg1 <- lmer(logEggMass ~ logMassSpecies + logMassGenus + (1|Genus),
data=Eggs)
```

where

- **logEggMass** is the log of the average egg mass (base e) for each species
- **logMassSpecies** is the log of the average body mass (base e) for each species, after correcting for genus mass (this correction is only important for the final question)
- **logMassGenus** is the log of the average body mass (base e) for each genus
- **Genus** is the genus that a species belongs to

13 Equation

Write down the equation for this model (either in matrix or scalar form).

(if you struggle to use the maths tool, write it on paper)

Maximum marks: 3

14 Explain the equation

Explain how each part of the equation relates to the model in R code.

Fill in your answer here

Maximum marks: 3

When the model was fitted, the following summary was obtained

```
Linear mixed model fit by REML ['lmerMod']
Formula: logEggMass ~ MassSpecies + MassGenus + (1 | Genus)
Data: Eggs
```

REML criterion at convergence: 98.4

Scaled residuals:

Min	1Q	Median	3Q	Max
-2.8334	-0.3899	-0.0264	0.3860	4.7210

Random effects:

Groups	Name	Variance	Std.Dev.
Genus	(Intercept)	0.19338	0.4397
	Residual	0.02745	0.1657

Number of obs: 237, groups: Genus, 101

Fixed effects:

	Estimate	Std. Error	t value
(Intercept)	-1.17468	0.15623	-7.519
MassSpecies	0.54835	0.03261	16.817
MassGenus	0.75012	0.02484	30.193

Correlation of Fixed Effects:

	(Intr)	MssSpc
MassSpecies	0.000	
MassGenus	-0.956	0.000

15 Egg Size Effects

There is a theory that the effect of mass on egg size should be lower within a species than at the genus level, i.e. the coefficient for *MassSpecies* < *MassGenus*.

Test whether the data support this theory using the information provided in the model output. If you cannot do the test, explain how you could do it if you had time/the correct information.

Fill in your answer here (or write your answer on a piece of paper)

Maximum marks: 5

16 Likelihood Ratio Tests

It is possible to do a likelihood ratio test, with some work:

- (a) how would you arrange the data to do the test, and what models would you fit?
- (b) if you managed to obtain the correct models, how would you perform the test?

(you should be able to answer part b even if you cannot do part a)

Fill in your answer here

Maximum marks: 8

