

Question	Out of	Question	Out of	Total	100
1	20	4	20	20	20
2	20	5	20	20	20
3	20				

- In the event of a fire alarm, do not check your cell phone when escorted outside.
- Please bring it to the attention of an Exam Facilitator so it can be recorded before leaving the exam hall.
- If you are feeling ill and unable to finish your exam, hand in.
- When you are done your exam, raise your hand and collect your bag and jacket before your exam is handed in.
- If you are feeling ill and unable to finish your exam, please bring it to the attention of an Exam Facilitator so it can be recorded before leaving the exam hall.
- Turn off and place all cell phones, smart watches, electronic devices, and unauthorized study materials in your bag under your desk. If it is left in your bag during an exam, you may be charged with an academic offence.
- As a student, you help create a fair and inclusive writing environment. If you possess an unauthorized aid during an exam, you may be charged with an academic offence.
- Do not begin writing the actual exam until the announcements have ended and the Exam Facilitator has started the exam.
- Fill out your name and student number on the top of this page.
- Fill out your name and student number on the top of this page.

U of T Email:

Student Number: 42

Last Name: MAUL

First Name: SOLUTIONS

Aids Allowed: Non-programmable calculator

Duration: 3 hours

Methods of Data Analysis II

STA303H1S

August 2019 Examinations

Faculty of Arts and Science

UNIVERSITY OF TORONTO

$$\begin{aligned}
 & \left(\frac{x}{\theta} - \frac{\theta}{2x} - \frac{1}{2} (\log(2\pi x^2) + \frac{1}{2}) \right) = \exp \left(\frac{x}{\theta} - \frac{\theta}{2x} - \frac{1}{2} (\log(2\pi x^2) + \frac{1}{2}) \right) \quad (1) \\
 & \left(\frac{x}{\theta} - \frac{\theta}{2x} - \frac{1}{2} (\log(2\pi x^2) + \frac{1}{2}) \right) = \exp \left(\frac{x/\theta^2 - \theta/x}{2} \right) = \\
 & \left(\frac{x}{\theta} - \frac{\theta}{2x} + \frac{\theta^2}{4x^2} - \frac{\theta^2}{4x^2} \right) = \exp \left(\frac{x}{\theta} - \frac{\theta}{2x} - \frac{\theta^2}{4x^2} \right) = \\
 & \exp \left(\frac{x}{\theta} - \frac{\theta}{2x} - \frac{\theta^2}{4x^2} \right) = f(x) = \exp \left(\frac{x}{\theta} - \frac{\theta}{2x} - \frac{\theta^2}{4x^2} \right) \quad (2)
 \end{aligned}$$

new + formula for $f(x)$ at (2)

for some $b(\theta)$ \leftarrow if (they don't fit \leftarrow $b(\theta) = 0$)

$\text{Exp form: } f(x|\theta, \phi) = \exp \left(\frac{\phi}{\theta} - b(\theta) + c(y|\phi) \right)$

- (a) (5 marks): Write down what it means for a distribution to be in the exponential family, and show that $\text{IG}(u, 1)$ is such a distribution by writing down the likelihood for this model in exponential family form.
- Explicitly specify the dispersion parameter and whichever functions $(b(\cdot), c(\cdot, \cdot), \text{etc.})$ you're using.

where θ_i is the canonical parameter and $\beta \in \mathbb{R}^p$ the parameter to be estimated.

$$\theta_i = \frac{\beta_i^T x}{L}$$

$$Y_i \sim \text{IG}(u_i, 1)$$

Let (Y_i, X_i) be independent response/covariate measurements, with $X_i \in \mathbb{R}^p$ and $Y_i \sim \text{IG}(u_i, 1)$. We propose a generalized linear model for Y_i ,

$$f_Y(x) = \sqrt{\frac{2\pi x^3}{(x - \mu)^2}} \exp \left(-\frac{2\mu^2 x}{(x - \mu)^2} \right), \quad x > 0.$$

- the density of Y is given by
1. **Exponential Families (20 marks)**: A random variable Y follows an $\text{IG}(u, 1)$ distribution with $u > 0$ if

$$?X(B_1)X + ?h) \stackrel{?}{=} \underline{\underline{Z}}$$

$$\textcircled{1} \quad \frac{\partial \Phi}{\partial x^j} = -\frac{1}{2} V_i - \frac{1}{2} \epsilon_{ijk} \partial_k \Phi$$

$$\textcircled{1} \quad x = d_1 \times \cancel{\left(\frac{d_2}{d_3} \right)} \cdot \frac{dp}{p} = \frac{dp}{p} ; \quad \frac{dp}{p} \cdot \frac{e^p}{e^p} = \frac{de}{e^p} \quad \textcircled{1}$$

$$\textcircled{1} \quad f(B) = \sum_{i=1}^n Y_i B - \alpha \sum_{i=1}^n \theta_i + C$$

$$= \exp\left(\frac{-2}{\sum_{i=1}^n y_i^2 E_i - \frac{1}{E_i}}\right) + \prod_{i=1}^n (y_i! E_i)$$

$$L(\beta, y|x) = \prod_{i=1}^n \exp\left(\frac{y_i}{\theta_i - \beta}\right) + C(y_i|\phi)$$

$$Q = \frac{\epsilon P}{(kT)e} \text{ where } Q \propto T^4$$

(c) (5 marks) Derive the score equations to which the maximum likelihood estimator $\hat{\theta}$ is the solution.

$$\textcircled{3} \quad \frac{1}{h_1^2} - \frac{M_3}{2} (h_1 - h_4) (\sin \gamma) =$$

$$\text{left} =$$

$$\text{Lin. Regr.: } z_i = g(\mu_1) + g(\mu_2) + \dots + g(\mu_n)$$

$$\textcircled{3} \quad \sin \theta = \frac{y}{r} = \frac{1}{\sqrt{2}} \Rightarrow \theta = 45^\circ$$

(b) (10 marks) Write down the link function, linearized response, and variance function for this model.

$$\text{Value function: } b_{ii} = \max_{\theta} \left(\theta^T \Theta_i - \frac{\theta^T \Xi_i}{2} \right)$$

A statistical model is a useful representation of the truth. How usually we ask the question "What is the most likely truth, given the data we observed?"

Bayesian statistics instead asks "What truths are consistent with the data we observed?" and then gives an answer by averaging over all such potential models, weighted according to how they give much more accurate and stable answers to the account of uncertainty in addition and fit the data. This approach is a special way in practice.

5) Anything reasonable gets full marks

(a) (5 marks): As a statistical analyst working at a bank, you've decided to build a Bayesian regression model to predict future losses that a portfolio will incur as a function of some covariates. However, your business executives are skeptical; they know little beyond the basics of frequentist statistics and linear/logistic regression, and they get anxious when they hear advanced words like "Bayesian" language; don't use any terminology or notation that they wouldn't understand.

where X is a known $n \times p$ matrix of full rank, $y \in \mathbb{R}^n$ is a vector of observations, $\beta \in \mathbb{R}^p$ is a vector of parameters to be estimated, and $\epsilon \sim N(0, \sigma^2 I)$.

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

model. That is, 2. Bayesian Theory (20 marks). In this question, we assume that our data follows a standard regression

(c) (2 marks): Circle true or false: a large choice of λ will usually lead to smaller coefficient estimates β .

(2)

$$\text{Q.E.D.} \quad h_T X_1 (X_1^T X_1 + \lambda I) = \beta \quad \text{or} \quad \Rightarrow h_T X = \beta (I + \lambda X_1^T X_1)$$

$$h_T X \beta = \beta I + \beta \lambda X_1^T X_1 \Leftarrow$$

$$S(\beta) = 0 \Leftarrow (X_1^T X_1 + \lambda I)^{-1} X_1^T Y = \beta$$

$$S(\beta) = \frac{\partial f}{\partial \beta} = (Y - X\beta)^T (Y - X\beta) + \lambda \beta^T \beta = 0$$

$$\text{Let } f(\beta) = (Y - X\beta)^T (Y - X\beta) + \lambda \beta^T \beta$$

Derive an explicit form for β . (Hint: this was done in lecture and on Test 2.)

$$\beta = \arg \min_{\beta} ((y - X\beta)^T (y - X\beta) + \lambda \|\beta\|_2^2)$$

(b) (5 marks): Recall that in the penalized regression version of ordinary least squares, called ridge regression, the estimated coefficients minimize the penalized residual sum of squares

on β

The two expressions are equal up to constants not depending

$$\text{Want: } \exp\left(\frac{\beta}{2}\right) (B - Za)^T (B - Za) = \exp\left(\frac{\beta}{2}\right) (B^T B - 2a^T B + a^T a)$$

$$\underbrace{(B + X^T X)}_{=Z} \underbrace{B^T B}_{=a^T a} - \underbrace{2a^T B}_{=2a^T \beta} = (B^T B) \exp(\beta)$$

$$\text{Likelihood: } \frac{1}{2} (y - XB)^T (y - XB) = L(\beta)$$

density is proportional to $\exp\left(\frac{\beta}{2}(y - Za)^T (y - Za)\right)$, where $Z = (X^T X + \lambda I)^{-1}$ and $a = X^T y$.
 the estimator you derived in part (b). Hint: by completing the square in B , show that the posterior theorem to prove that the mean of the posterior distribution of β given the data (X, y) coincides with

(d) (8 marks): Assume we put a $N(0, \sigma^2/\lambda)$ prior distribution independently on each β_j . Use Bayes'

3. **Linear Mixed Effects (20 marks)**: Recall from lecture the Rat Growth data: data on the weights of rats measured for 5 consecutive weeks. We are interested in modelling growth as a function of time.

$\epsilon_{ij} \sim N(0, \sigma^2_\epsilon)$: homoscedastic randomization/error term.

$\alpha_i + \beta_i w_{ij} \sim V_j$: random slope, $V_j \sim N(0, \sigma^2_{Vj})$. Penetration of each weight, growth of individual random effects means from the overall rate β_1 .

β_0, β_1 : population intercept and slope for rat weight/week.

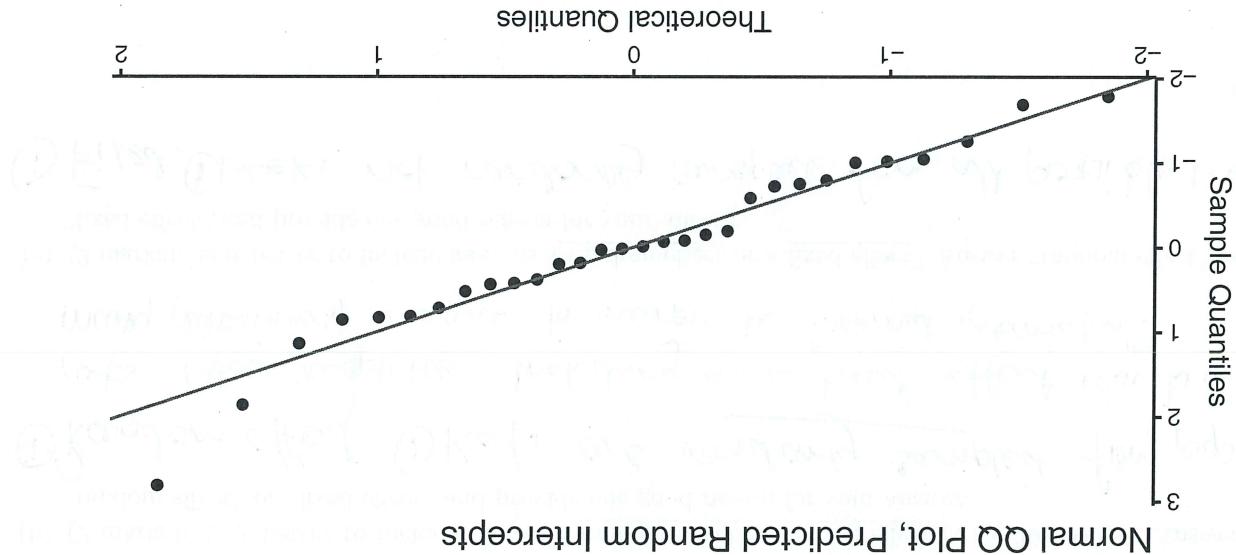
y_{ij} : weight measurement for rat in j^{th} week

Model 3: same as model 1; estimation procedure different.

$$\text{Model 2: } y_{ij} = \beta_0 + \beta_1 \text{week}_j + u_i + v_j \times \text{week}_j$$

$$\text{Model 1: } y_{ij} = \beta_0 + \beta_1 \text{week}_j + u_i + \epsilon_{ij}$$

(a) (4 marks): Write down the full statistical model for each of the three models, clearly defining all terms including all parameters and distributions (if you use the same term(s) in multiple models, you only need define them once).



① No ② newer fit with REML have in comparable likelihoods

(e) (3 marks): Consider the models "RAT MODEL 1" and "RAT MODEL 2". The latter contains one more term than the former; hence we wish to compare them using a procedure for nested model comparison (i.e. a likelihood ratio test). Can we do this? Answer "Yes" or "No" and briefly explain why or why not.

$$\textcircled{1} \quad \%St = \frac{\overbrace{Q_1 + Q_2}^{\text{Total}}}{\overbrace{Q_1}^{\text{Initial}}} \quad \Leftrightarrow \quad Q_2 = 64.2a$$

(d) (1 mark): For the model “RAT MODEL 1”, what is the proportion of total variance explained by rat?

① Fixed ① weeks not randomly sampled from all possible weeks

(c) (2 marks): Is it better to include week as a random effect or a fixed effect? Answer “random effect” or “fixed effect” and provide one good reason for your answer.

① Random effect: ① Reefs are randomly sampled from lots of reefs. Also acceptable: including as a fixed effect reefs in lots many parameters, so want to incorporate "external information".

(b) (2 marks): Is it better to include rat as a random effect or a fixed effect (with 30 levels)? Answer „random effects“ or „fixed effect“ and provide one good reason for your answer.

- (e) (2 marks): Consider the plot titled "Normal QQ Plot, Predicted Random Intercept", which is a normal Q-Q-plot of the predicted random intercepts from "RAT MODEL 1". What model assumption is being tested here?
- $Q_i \sim N(\bar{Q}_i)$, normality of intercepts
- (f) (2 marks): Consider the plot titled "Normal QQ Plot, Predicted Random Intercept", which is a normal Q-Q-plot of the predicted random intercepts from "RAT MODEL 1". Why are the estimated variances different?
- Model 1 has RML \Rightarrow OLS unadjusted
Model 2 has RML \Rightarrow OLS biased.
- (g) (2 marks): Consider the models "RAT MODEL 1" and "RAT MODEL 3". Why are the estimated variances different?
- (h) (2 marks): If we wanted to predict the weight in a given week after birth for a given rat from the 30 included in the dataset using "RAT MODEL 1", what would be the prediction equation?
- $y_{ij} = \hat{\alpha}_{1j} + b_{1j} z_{ij}$
- (i) (2 marks): If we wanted to predict the weight in a given week after birth for a new rat that was not one of the 30 included in the dataset using "RAT MODEL 1", what would be the prediction equation?
- $y_{ij} = \hat{\alpha}_{1j} + b_{1j} z_{ij}$

```

## Number of equivalent replicates : 3.972
## Expected number of effective parameters(std dev) : 37.77(3.744)
## Precision for id 10.39 4.181   4.798   9.569   20.81 8.271
## mean    sd 0.025quant 0.5quant 0.975quant mode
## Model hyperparameters:
##      #
##      # id IID model
##      # Name Model
## Random effects:
##      #
##      # conc:brood3 -1.8479 0.2461 -2.3357 -1.8433 0
##      # conc:brood2 -1.6885 0.2498 -2.1832 -1.6841 0
##      # brood3  1.3576 0.1355 1.0961 1.3561 1.6278 1.3531 0
##      # brood2  1.1747 0.1382 0.9076 1.1733 1.4496 1.1706 0
##      # conc   -0.0509 0.2268 -0.5001 -0.0499 0.3919 -0.0479 0
##      # (Intercept) 1.6376 0.1414 1.3554 1.6391 1.9117 1.6420 0
##      # mean    sd 0.025quant 0.5quant 0.975quant mode kld
##      # Fixed effects:
##      #
##      # Time used:
##      # Pre-processing Running inLA Post-processing Total
##      # 1.7597 0.2016 0.0937 2.0550
##      # Call:
##      # C("inLA(formula = live ~ conc * brood + f(id, model = "iid", prior = "pc.proc", param
##      # Model formula: live ~ conc * brood +
##      # prior = "pc.proc", param = c(3, .75))
##      # f(id, model = "iid",
##      # Model formula by 300 (the maximum value is 310), you fit a model using INLA:
## The counts are naturally grouped, with counts from the same zooplankton expected to be more similar
## than counts from different zooplankton. We wish to model the linked mean with one random-effect for each
## population, and linear terms corresponding to the available covariates and their interaction. After scaling
## the concentration by dividing by 300 (the maximum value is 310), you fit a model using INLA:
## $ live <dbl> 3, 14, 10, 5, 12, 15, 6, 11, 17, 6, 12, 15, 6, 15, 5...
## $ brood <chr> "1", "2", "3", "4", "5", "6", "7", ...
## $ id    <int> 1, 1, 2, 2, 3, 3, 4, 4, 5, 5, 6, 6, 7, ...
## $ conc  <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ...
## Variables: 4
## Observations: 150
## Offspring in each brood was recorded. The data is as follows:
## different concentrations of nitrofen. Each animal then gave birth to three broods, and the number of live
## offspring in each brood was recorded. The data is as follows:
## Nitrofen, a herbicide. Specifically, 50 zooplankton were split into 10 groups of 5 each, and exposed to
## different concentrations of nitrofen, which consists of living offspring of zooplankton exposed to various concentrations
## of nitrofen data, you have been provided the

```

4. GLMMs and INLA (20 marks). As part of your new job as an agriculturist, you have been provided the

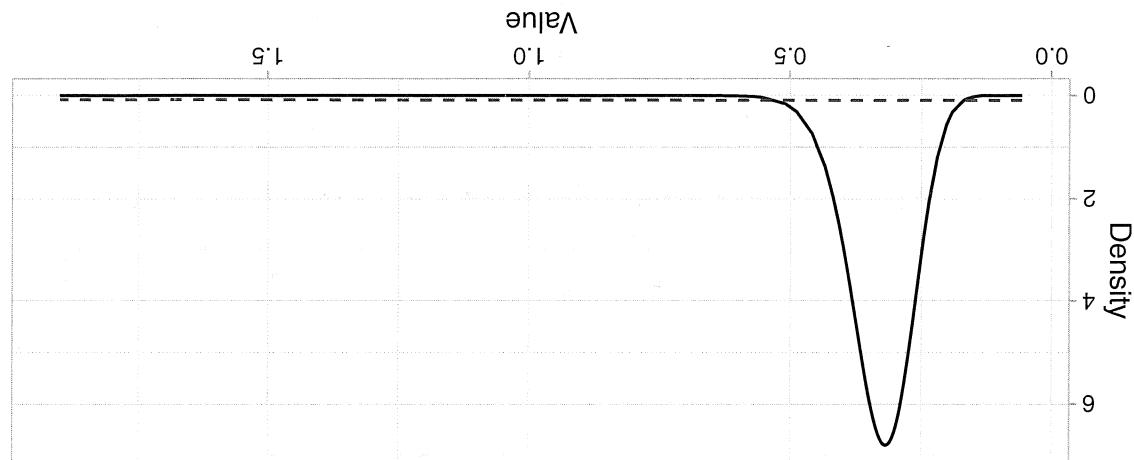
2. $j = 1, \dots, 3$ represents the j_{th} brood for each zooplankton,

1. $i = 1, \dots, 50$ represents the i_{th} zooplankton,

where

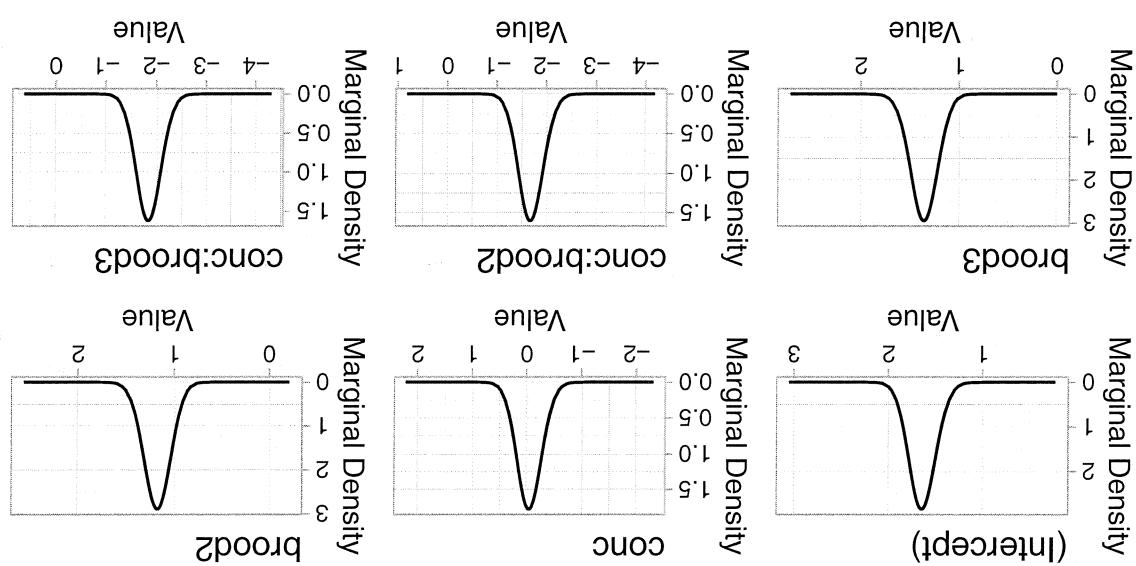
$$(0.1) \quad \begin{aligned} u_i &\sim \text{Normal}(0, \sigma_u^2) \\ \log A_{ij} &= \eta_{ij} = \beta_0 + \beta_1 x_i + \beta_2 b_j + \beta_3 x_i \times b_j + u_i \\ Y_{ij} &\sim \text{Poisson}(\lambda_{ij}), i = 1, \dots, 50; j = 1, \dots, 3 \end{aligned}$$

The full hierarchical model takes the form



Red: prior. Black: posterior

Posterior standard deviation of random effect



```
## Marginal Log-Likelihood: -431.88
## SD for id 0.3269879 0.06006813 0.2194122 0.3231583 0.4557495 0.3163052
##   mean    sd   q0.025   q0.5   q0.975   mode
## 0.3269879 0.06006813 0.2194122 0.3231583 0.4557495 0.3163052
```

and interaction. including λ_{ij} , ② for y_{ij} , and ② for properly scaling b_{ij} , because not enough info. Give ② for known formula for x_{ij}^B ; ② for

$$\lambda_{ij} = 1.6376 - 0.0509 + 1.1417 - 1.6885 \cdot t_{ij}$$

$$\text{Corr} = 1 - \lambda_{ij} = 1.6376 - 0.0509 + t_{ij}$$

$$\lambda_{ij} = 1.6376 + 1.1417 + t_{ij}, \text{ etc}$$

95% CI = $\lambda_{ij} \pm \dots$ not enough information.

$$\text{Corr} = 0 : \lambda_{ij} = 1.6376 + t_{ij}, \quad y_{ij} = \exp(\lambda_{ij})$$

estimates with 95% confidence intervals on the natural scale.
zero concentration of the herbicide? How about when concentration increases to 1? Provide point

(c) (8 marks): Based on the model, what happens to the average numbers of offspring for each brood at

but we they can result the numbers off the
offspring

there's about a 75% chance that $a_u > 3$. Write down the corresponding parameters for PC-prec.
standard deviation of flea-means (a very rough proxy for the scale of parameter a_u), you decided that

(b) (3 marks): You used a PC prior for the random effect standard deviation. Based on the observed

in a brood for each zoonplankton varies due to the zoonplankton.

③ Q. describes how much the average number of offspring

the sole hyperparameter in the model to them in one sentence.
banking to agriculture, and they're clues about mixed effects models. Explain the interpretation of

(a) (3 marks): Unfortunately, the businesses executives from Problem 2a have also switched industries from

zooplankton.

6. u_i is a random effect designed to capture correlation in offspring counts across broods from the same

5. b_j is an indicator variable representing a fixed intercept for the j th brood,

4. x_i is the concentration of nitrogen to which the i th zooplankton was exposed,

3. X_u is the count of live offspring in the j th brood of the i th zooplankton,

④ Just recall the quantities directly: (Q_{25}, Q_{75})

Give 2 marks if they do $Q_5 + 2SE(Q_5)$

(d) (4 marks): INLA reports summary statistics of the posterior of the log precision, which is $-2\log(\sigma_u)$, rather than that of σ_u (which is what we prefer). Provide a 95% confidence interval for σ_u .

(e) (2 marks): Does the individual variation in α_d outweigh the broad and concentration effects? Answer "Yes" or "No" and justify based on the INLA output.

This question is intentionally vague and full marks should be given for any well thought out answer that reflects the posterior of Q_5 , and that the inferences

of the β_s should be given for any well thought out answer that reflects the posterior of Q_5 , and that the inferences

① ② ③

Residuals = 0. Perfect fit.

“Yes” or “No” and justify based on the residual deviance.

(a) (5 marks): Coefficients and standard errors notwithstanding, does the model fit the data well? Answer

Look at those coefficient estimates! Look at those standard errors! Clearly something strange has happened.

```

## Number of Fisher Scoring iterations: 23
##
## AIC: 4
##
## Residual deviance: 3.3510e-10 on 4 degrees of freedom
## Null deviance: 8.3178e+00 on 5 degrees of freedom
## Dispersion parameter for binomial family taken to be 1
##
## x
## (Intercept) -1.828 46931.107 0 1
## Estimate Std. Error z value Pr(>|z|)
## Coefficients:
## Deviance Residuals:
## glm(formula = y ~ x, family = "binomial", data = mystery)
## Call:
## 
```

Surprise, R gives you a warning (which has been suppressed) and outputs this model summary:

Since y is binary and x is continuous, you try to fit a logistic regression model to predict y using R. To your

```

## $ y <dbl> 1, 1, 0, 0, 0
## $ x <dbl> -4.757, -3.941, -4.413, 4.453, 3.349, 4.899
## Variables: 2
## Observations: 6
## 
```

wisely). The mystery data contains six observations as follows:

5. Logistic Regression Mystery (20 marks). (Note: this question is challenging. Manage your time

- ② $\hat{\beta} = \arg \min_{\beta} \sum_i \log \left(\frac{1}{1 + e^{-x_i^T \beta}} \right)$
- From the graph, we can see that the function $\log \left(\frac{1}{1 + e^{-x_i^T \beta}} \right)$ is strictly increasing and convex. As $x_i^T \beta \rightarrow -\infty$, the value of the function approaches $-\infty$. Hence $\hat{\beta}$ does not exist.
- $\Rightarrow \log \left(\frac{1}{1 + e^{-x_i^T \beta}} \right) \rightarrow -\infty$ as $|x_i^T \beta| \rightarrow \infty$ and when $y_i=0$, $x_i^T \beta > 0$.
- ③ $\hat{\beta}$ pointing in a direction such that when $y_i=1$, $x_i^T \beta < 0$
- Because the objective is linearly separable, you can pick

$$(2) \quad \ell(\beta) = \sum_{i=1}^{n=0} \log \left(\frac{1}{1 + e^{-x_i^T \beta}} \right) - \sum_{i=1}^{n=1} \log \left(1 - \frac{1}{1 + e^{-x_i^T \beta}} \right)$$

one for the $y_i = 1$ data. What happens when you try to maximize $\ell(\beta)$?
 Use this to explain mathematically why the coefficient estimates have such a large magnitude and why the estimation procedure failed. (Hint: start by splitting the above into two sums, one for the $y_i = 0$ data and one for the $y_i = 1$ data. What happens when you try to maximize $\ell(\beta)$?)

$$\ell(\beta) = \sum_u \left[y_i \cdot \log \left(\frac{1}{1 + e^{-x_i^T \beta}} \right) - (1 - y_i) \cdot \log \left(1 - \frac{1}{1 + e^{-x_i^T \beta}} \right) \right]$$

- (b) (6 marks): Recall that the log-likelihood for binary logistic regression takes the general form

(note: formally, calling $-X\|B\|^2$ makes this function convex, but students don't know how to say this to get marks).

Larger values of B are now penalized. ②

$$\hat{B} = \text{argmax}_B (L(B) - \lambda \|B\|_1) \text{ for some } \lambda \text{ would solve this.}$$

identifed in (b)).

for an estimator of B , but you do need to convince us that your approach mitigates the problem(s) an approach has already come up elsewhere on this exam. No need to derive an explicit expression and show that your idea keeps the magnitude of the coefficient estimates from blowing up. (Hint: such an alternative regression approach that does not involve removing any variables or modifying the data, alternaative regression approach that does not involve removing any variables or modifying the data,

(d) (5 marks): Use your knowledge of Bayesian statistics and/or penalized likelihood to propose an

The fitted probabilities \hat{P}_i are all nearly 0 or 1, so $\hat{P}_i(1-\hat{P}_i) \approx 0$ and $(X^T W X)^{-1}$

why the standard errors are so large.

where the weight matrix W is diagonal with entries $W_{ii} = \hat{P}_i \cdot (1-\hat{P}_i)$. Use this to explain mathematically

$$\text{Var}(g) = (X^T W X)^{-1}$$

(c) (5 marks): Recall that R computes approximate standard errors for the coefficient estimates using

THIS PAGE IS FOR ROUGH WORK. NOTHING ON THIS PAGE WILL BE MARKED.

THIS PAGE IS FOR ROUGH WORK. NOTHING ON THIS PAGE WILL BE MARKED.

THIS PAGE IS FOR ROUGH WORK. NOTHING ON THIS PAGE WILL BE MARKED.

THIS PAGE IS FOR ROUGH WORK. NOTHING ON THIS PAGE WILL BE MARKED.
