

Name: _____

uniq name: _____

BIOSTAT 651
APPLIED STATISTICS II: EXTENSIONS OF LINEAR REGRESSION

Test #1
Wednesday, February 17, 2016
1:10-2:30 p.m.

Statistical table and blank paper are provided.

<u>Question</u>	<u>Points Possible</u>	<u>Points Received</u>
1	13	_____
2	12	_____
3	25	_____
Total	50	

1. (13 points, total) The Pareto distribution with a known scaling parameter $\alpha > 0$ is given

$$f(Z = z; \beta) = \frac{\beta \alpha^\beta}{z^{(\beta+1)}}, \quad z > \alpha, \beta > 0.$$

- (a) (4 points) Write out $P(Z = z; \beta)$ as an exponential family form (Note that $a(\phi) > 0$).

- (b) (4 points) Log transformation of Y, $Y = \log(Z/\alpha)$, is commonly used as an outcome in regression analysis. The density function of Y is given as

$$f(Y = y; \beta) = \beta e^{-\beta y}, \quad y > 0, \beta > 0.$$

Suppose that we have n independent observations (X_i, Y_i) . Write out the above density function as an exponential family form and determine the canonical link function $g(\mu_i) = X_i^T \beta$, where $\mu_i = E(Y_i)$. (Note that $a(\phi) > 0$)

- (c) (5 points) Determine the deviance. It should be expressed as a function of Y_i and $\hat{\mu}_i$ ($i = 1, \dots, n$), where $\hat{\mu}_i$ is the estimated μ_i under the model in (b).

2. (12 points, total) Suppose that the response Y_i follows a $\text{Poisson}(\lambda_i)$ distribution

$$f(Y_i|\lambda_i) = \frac{\lambda_i^{Y_i} e^{-\lambda_i}}{Y_i!}.$$

Assume that the only covariate, x_i , acts additively on λ_i , i.e., $\lambda_i = x_i\beta$.

- (a) (4 points) Derive score function, $U(\beta)$, and expected Fisher information, $I(\beta)$, as a function of Y_i , x_i and β .

(b) (4 points) Suppose now that the following data are observed:

x_i	1	2	3	4	5
Y_i	2	4	6	8	10

Compute the maximum likelihood estimate of β and the corresponding 95% confidence interval.

(c) (4 points) Carry out a score test for $H_0 : \beta = 2.5$ versus $H_1 : \beta \neq 2.5$.

3. (23 points, total) Researchers are interested in whether a new anti-viral drug can decrease the number of viral RNA copies. 30 patients were randomly assigned to drug ($X_1 = 1$) and placebo groups ($X_1 = 0$), and the viral RNA count per mL (Y) was measured using RNA-sequencing. In addition, age in years (X_2) was measured as a covariate. The following regression model has been proposed for this analysis

$$\log(\lambda_i) = \beta_0 + \beta_1 X_1 + \beta_2 X_2,$$

where Y_i follows the poisson distribution

$$f(Y_i|\lambda_i) = \frac{\lambda_i^{Y_i} e^{-\lambda_i}}{Y_i!}.$$

The estimate of the inverse of the Fisher information matrix is

$$I(\hat{\beta})^{-1} = \begin{pmatrix} 0.0008553 & -0.000054 & -0.000035 \\ -0.000054 & 0.0002458 & -3.745 \times 10^{-6} \\ -0.000035 & -3.745 \times 10^{-6} & 1.7028 \times 10^{-6} \end{pmatrix}$$

Subset of the SAS outputs are provided in a separate document.

- (a) (4 points) Interpret β_1 and β_2

- (b) (4 points) Estimate $\exp(\beta_1)$ and its 95% confidence interval.

(c) (4 points) Derive the likelihood ratio test for the drug effect.

- * Write out full and reduced models.

- * Calculate the likelihood ratio test (LRT) statistic using SAS outputs.

- * What is a conclusion based on the LRT?

(d) (5 points) Suppose that the first individual has $(Y_1, X_{i1}, X_{i2}) = (530, 0, 22)$ and the corresponding leverage $h = 0.3$. Obtain the standardized Pearson residual and the Cook's distance. Is this a high leverage observation (YES/NO)? or high influence observation (YES/NO)? Please justify your answer.

(e) (4 points) Carry out a Goodness of Fit test and state your conclusion.

(f) (4 points) Researchers want to investigate multicollinearity among covariates. Since proc genmod does not calculate the variance inflation factor (VIF), they plan to use proc reg with the following command. Is it a right way to calculate VIF in this model (YES/NO)? Please justify your answer.

[SAS Code]

```
proc reg data=ViralRNA;  
model Y = X1 X2/ vif;  
run;
```