

### Question 1

A1 Consider the analysis of covariance model without interaction, denoted by  $1 + X + A$  (that is,  $X$  is a continuous covariate while  $A$  is a factor with  $r$  levels).

- (a) Write the formula for the model in such a way that the parameters are *NOT* identifiable. Show the corresponding model (i.e. design) matrix.

let  $\alpha_i \in \{ \text{group } \# 2, \text{ group } \# 3, \dots, \text{ group } \# r \}$

Let  $Y_i = X_i\beta + \varepsilon_i = \beta_0 + \beta_1 z_i + \beta_2 \cdot 1_{(\alpha_i = \text{group } \# 2)} + \dots + \beta_r \cdot 1_{(\alpha_i = \text{group } \# r)}$

and  $\beta = (\beta_0, \beta_1, \dots, \beta_r)$

then

$$X = \begin{pmatrix} 1 & z_1 & 1 & 0 & 0 & 0 \\ 1 & z_2 & 0 & 1 & 0 & 0 \\ 1 & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & z_r & 0 & 0 & 0 & 1 \end{pmatrix}_{p \times r}$$

- (b) For the parameters in (a), give an example of a characteristic that is (i) estimable,  
(ii) not estimable.

(i) estimable :  $\beta_1$  is estimable

Now consider

$$X = \begin{pmatrix} 1 & z_1 & 1 & 0 & \dots & 0 \\ 1 & z_1+1 & 1 & 0 & \dots & 0 \\ 1 & z_2 & 0 & 1 & 0 & 0 \\ 1 & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & z_r & 0 & 0 & 0 & 1 \end{pmatrix}_{p \times (r+1)}$$

$$l^T \beta = \beta_1 = \alpha^T E[y] = \begin{bmatrix} -1 \\ 1 \\ \vdots \\ 0 \end{bmatrix}^T \begin{bmatrix} \beta_0 + \beta_1 z_1 + \beta_2 \\ \beta_0 + \beta_1 (z_1+1) + \beta_2 \\ \beta_0 + \beta_1 z_2 + \beta_3 \\ \beta_0 + \beta_1 z_3 + \beta_4 \\ \vdots \end{bmatrix}$$

not estimable : for example  $\beta_2$

- (c) Now express the model so that the parameters are identifiable. Explain how to interpret them. Show the model matrix when  $A$  divides the sample into three groups (i.e.  $r = 3$ ), each containing two observations.

Now by deleting the parameter  $\beta_p$ , the model becomes:

$$\begin{pmatrix} 1 & z_1 & 1 & 0 & \dots & 0 \\ 1 & z_2 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \vdots & & 1 \\ 1 & z_r & 0 & 0 & & 0 \end{pmatrix}_{(p-1) \times r}$$

The model matrix can be interpreted as  $A$  representing a factor with  $p-3$  levels. If  $z_i = z_j$  for some  $i \neq j \neq r$  and they are not from the same group, then the difference between  $E[z_i]$  and  $E[z_j]$  is the difference between  $\beta_{i+2}$  and  $\beta_{j+2}$ .

- Now when  $A$  divides the sample into 3 groups, each containing 2 observation

$$X = \begin{pmatrix} 1 & z_1 & 1 & 0 \\ 1 & z_2 & 1 & 0 \\ 1 & z_3 & 0 & 1 \\ 1 & z_4 & 0 & 1 \\ 1 & z_5 & 0 & 0 \\ 1 & z_6 & 0 & 0 \end{pmatrix}$$

## Question 2

A2 The inverse Gaussian distribution has density of the form

$$f(y; \nu, \lambda) = \left( \frac{\lambda}{2\pi y^3} \right)^{1/2} \exp \left( -\frac{\lambda(y-\mu)^2}{2\mu^2 y} \right)$$

for  $y > 0$ , with parameters  $\mu > 0$  and  $\lambda > 0$ .

- (a) Show that the family of inverse Gaussian distributions is an exponential dispersion family. Identify the functions  $b(\cdot)$ ,  $c(\cdot)$  as well as the canonical and the dispersion parameters.

Expand  $f(y; \nu, \lambda)$  to the form of  $\exp \left\{ \frac{y\theta - b(\theta)}{a(\theta)} + c(y, \phi) \right\}$

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

$$\text{then } b(\theta) = -\frac{2}{\mu} = -2\sqrt{\theta} \quad c(y, \phi) = \ln \left( \frac{\lambda}{2\pi y^3} \right)^{1/2} - \frac{\lambda}{2y} = \ln \left( \frac{z/\phi}{2\pi y^3} \right)^{1/2} - \frac{z/\phi}{2y}$$

canonical parameter  $\theta = -\frac{1}{\mu^2}$ , dispersion parameter  $\phi = z/\lambda$

$\Rightarrow$  The family of inverse Gaussian distribution is an exponential dispersion family.

- (b) Compute the mean and the variance of an inverse Gaussian random variable  $Y$  and identify the mean-variance relationship.

$$\text{mean : } E[Y] = b'(\theta) = (-\theta)^{-\frac{1}{2}} = (-\theta)^{-\frac{1}{2}} = \left(\frac{\mu}{\lambda}\right)^{-\frac{1}{2}}$$

$$\begin{aligned} \text{variance : } \text{Var}[Y] &= b''(\theta) \cdot a(\phi) = \frac{1}{2} (-\theta)^{-\frac{3}{2}} \cdot \frac{2}{\lambda} \\ &= (-\theta)^{-\frac{3}{2}} \cdot \frac{1}{\lambda} \\ &= \left(\frac{1}{\mu^2}\right)^{-\frac{3}{2}} \cdot \frac{1}{\lambda} \\ &= \frac{\mu^3}{\lambda} \end{aligned}$$

$$\text{mean variance relation : } V(\mu) = b''(\theta) = \frac{1}{2} (-\theta)^{-\frac{3}{2}}$$

- (c) Identify the canonical link for a GLM with inverse Gaussian responses. Do you think this link is sensible? What other link functions might be appropriate?

$$\text{Since } b'(\theta) = (-\theta)^{-\frac{1}{2}} \Rightarrow b'(\theta)^2 = -\frac{1}{\theta} \Rightarrow \theta = -\frac{1}{b'(\theta)^2}$$

$$\Rightarrow g(x) = (b')^{-1}(x) = -\frac{1}{x^2}$$

$$\Rightarrow g(E[Y|X]) = -\frac{1}{E[Y|X]^2} = x_B \Rightarrow \frac{1}{E[Y|X]^2} = -x_B$$

which can constraint the output of  $x_B$  to be negative, but its probably that in some cases  $x_B$  are positive and  $E[Y|X]$  may be an estimation of probability, Then this link is not sensible. The identity link might be more appropriate here.

- (d) Consider a GLM with inverse Gaussian responses and the canonical link. Write down the score equations. How do they change when the log-link is used?  $\Delta$

The score equation is  $\sum_{i=1}^n \frac{y_i - b'(\theta_i)}{a(\phi)b''(\theta_i)} \cdot \frac{1}{g'(\mu_i)} \cdot x_{ij} = 0$

$$b'(\theta_i) = \left(\frac{2}{\mu_i}\right)^{-\frac{1}{2}} \quad b''(\theta_i) = \mu_i^3/2 \quad a(\phi) = \frac{2}{\lambda}$$

- Consider the canonical link :  $g(\mu_i) = -\frac{1}{\mu_i^2} \Rightarrow g'(\mu_i) = 2\frac{1}{\mu_i^3}$

$$\Rightarrow \sum_{i=1}^n \frac{y_i - \left(\frac{2}{\mu_i}\right)^{-\frac{1}{2}}}{\frac{2}{\lambda} \cdot \frac{\mu_i^3}{2}} * \frac{\mu_i^3}{2} \cdot x_{ij} = 0$$

$$\sum_{i=1}^n \left(y_i - \left(\frac{2}{\mu_i}\right)^{-\frac{1}{2}}\right) * \frac{\lambda}{\mu_i^3} \cdot \frac{\mu_i^3}{2} \cdot x_{ij} = 0$$

$$\sum_{i=1}^n \frac{\lambda}{2} \left(y_i - \left(\frac{2}{\mu_i}\right)^{-\frac{1}{2}}\right) x_{ij} = 0$$

The optimal solution is  $y_i = \left(\frac{2}{\mu_i}\right)^{-\frac{1}{2}}$

- When using the log link, we will have a different function  $g(\mu_i) = \log \mu_i$ ,  $g'(\mu_i) = \frac{1}{\mu_i}$

$$\Rightarrow \sum_{i=1}^n \frac{y_i - \left(\frac{2}{\mu_i}\right)^{-\frac{1}{2}}}{\frac{2}{\lambda} \cdot \frac{\mu_i^3}{2}} * \mu_i * x_{ij} = 0$$

$$\sum_{i=1}^n \left(y_i - \left(\frac{2}{\mu_i}\right)^{-\frac{1}{2}}\right) * \frac{\lambda}{\mu_i^3} \cdot \mu_i \cdot x_{ij} = 0$$

$$\sum_{i=1}^n \left(y_i - \left(\frac{2}{\mu_i}\right)^{-\frac{1}{2}}\right) \frac{\lambda}{\mu_i^2} x_{ij} = 0$$

Then the optimal solution appears when  $y_i = \left(\frac{2}{\mu_i}\right)^{-\frac{1}{2}}$

# Question 3

## Question 3. R excercise.

Load the data set crabs available on MyCourses under Content -> Datasets, viz. This data set contains the number of satelites (males attached to the female's nest) for 173 female horseshoe crabs. The weight (in g), carapace width (in cm), spine condition and color of the crab are also recorded; weight and width are continuous whereas color (1=medium light, 2=medium, 3=medium dark, 4=dark) and spine condition (1= both good, 2 = one worn or broken, 3 = both worn or broken) are categorical. It is of interest to understand how the width, weight, color and spine condition affect the number of satelites.

```
crabs <- read.table("crabs.txt", header=TRUE)
attach(crabs)
head(crabs)

##   color spine width satell weight
## 1     3      3  28.3      8    3050
## 2     4      3  22.5      0    1550
## 3     2      1  26.0      9    2300
## 4     4      3  24.8      0    2100
## 5     4      3  26.0      4    2600
## 6     3      3  23.8      0    2100

library(faraway)
```

(a). Fit a linear regression model to the data with width and spine condition as main effects, using lm. Plot the data along with the fitted regression line(s) and interpret the model.

Fit the linear regression model using lm

```
width<-crabs$width
spine<-crabs$spine
# creating the factor variable
satell<-crabs$satell
fit.crabs<-lm(satell~factor(spine)+width)
summary(fit.crabs)

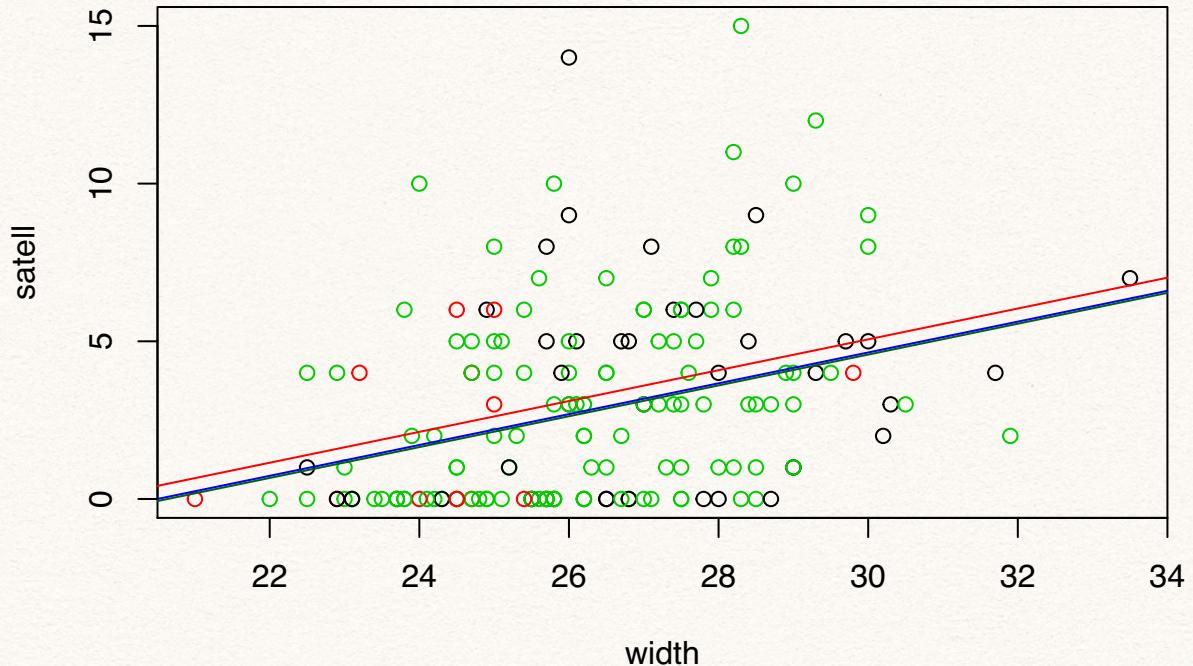
##
## Call:
## lm(formula = satell ~ factor(spine) + width)
##
## Residuals:
##    Min     1Q Median     3Q    Max 
## -4.426 -2.250 -0.734  1.846 11.185 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) -9.6081    3.0878  -3.112  0.00218 **  
## factor(spine)2 -0.4828    0.9514  -0.508  0.61245    
## factor(spine)3 -0.4156    0.5686  -0.731  0.46586    
## width         0.4890    0.1125   4.348 2.36e-05 ***
```

```

## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.982 on 169 degrees of freedom
## Multiple R-squared: 0.1185, Adjusted R-squared: 0.1028
## F-statistic: 7.57 on 3 and 169 DF, p-value: 8.816e-05

#plot the data
plot(width,satell,col=factor(spine))
# Line for width =1
# factor value of spine(1,2,3) just specifying the category so it will be described by indicator
# X2=0,X3=0 represents spine=1
# a represents the intercept and b represent the slope
abline(a=coef(fit.crabs)[1],b=coef(fit.crabs)[4],col="red")
#line for width=2
# X2=1,X3=0 represents spine=2
abline(a=coef(fit.crabs)[1]+coef(fit.crabs)[2],b=coef(fit.crabs)[4],col="darkgreen")
#line for width=3
# X2=0,X3=1 represent spine=3
abline(a=coef(fit.crabs)[1]+coef(fit.crabs)[3],b=coef(fit.crabs)[4],col="blue")

```



The above figures show that a better spins condition can cause a larger number of satellites. According to the line trend, we can conclude that higher width can cause larger number of satellites.

(b). Fit the same model as in part (b), but using `glm` with `family = gaussian`. If you leave the link unspecified, what link is being used? Using the function `summary`, compare the estimated coefficients and standard errors to those obtained in part (a). What do you think the quantity given in the output line (Dispersion parameter for gaussian family taken to be ...) estimates? How could you alternatively estimate the same quantity from the output in part (a)?

Fit the model in part(a) with unspecified link, then the canonical link is being used.

```

gaussianmod<-glm(satell~factor(spine)+width,family=gaussian)
summary(gaussianmod)

## 
## Call:
## glm(formula = satell ~ factor(spine) + width, family = gaussian)
## 
## Deviance Residuals:
##      Min     1Q Median     3Q    Max 
## -4.426 -2.250 -0.734  1.846 11.185 
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) -9.6081    3.0878  -3.112 0.00218 **  
## factor(spine)2 -0.4828    0.9514  -0.508 0.61245  
## factor(spine)3 -0.4156    0.5686  -0.731 0.46586  
## width         0.4890    0.1125   4.348 2.36e-05 *** 
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
## 
## (Dispersion parameter for gaussian family taken to be 8.892932)
## 
## Null deviance: 1704.9 on 172 degrees of freedom 
## Residual deviance: 1502.9 on 169 degrees of freedom 
## AIC: 874.96 
## 
## Number of Fisher Scoring iterations: 2

```

The coefficient and standard error obtained in part(a) are the same as those in part(b). As the dispersion parameter  $\phi$  is taken to be 8.892932 in part(b) and  $\phi = \sigma^2$  in Gaussian family. The estimated value of  $\sigma$  is equal to 2.982 in part(a), so both of these 2 parts gain the same value for the standard errors and dispersion parameter.

(c). Fit a Poisson glm to the data, again with width and spine condition as main effects, using the canonical link. Plot the data along with the fitted regression curve(s) and compare the plot to the one in part (a).

```

# Fit the model with poisson family
poissonmod<-glm(satell~factor(spine)+width,family=poisson)
summary(poissonmod)

## 
## Call:
## glm(formula = satell ~ factor(spine) + width, family = poisson)
## 
## Deviance Residuals:
##      Min     1Q Median     3Q    Max 
## -2.9509 -1.9740 -0.4963  1.0832  4.7173 
## 
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)    
## (Intercept) -3.02570   0.59421  -5.092 3.54e-07 ***  
## factor(spine)2 -0.19932   0.20983  -0.950   0.342  

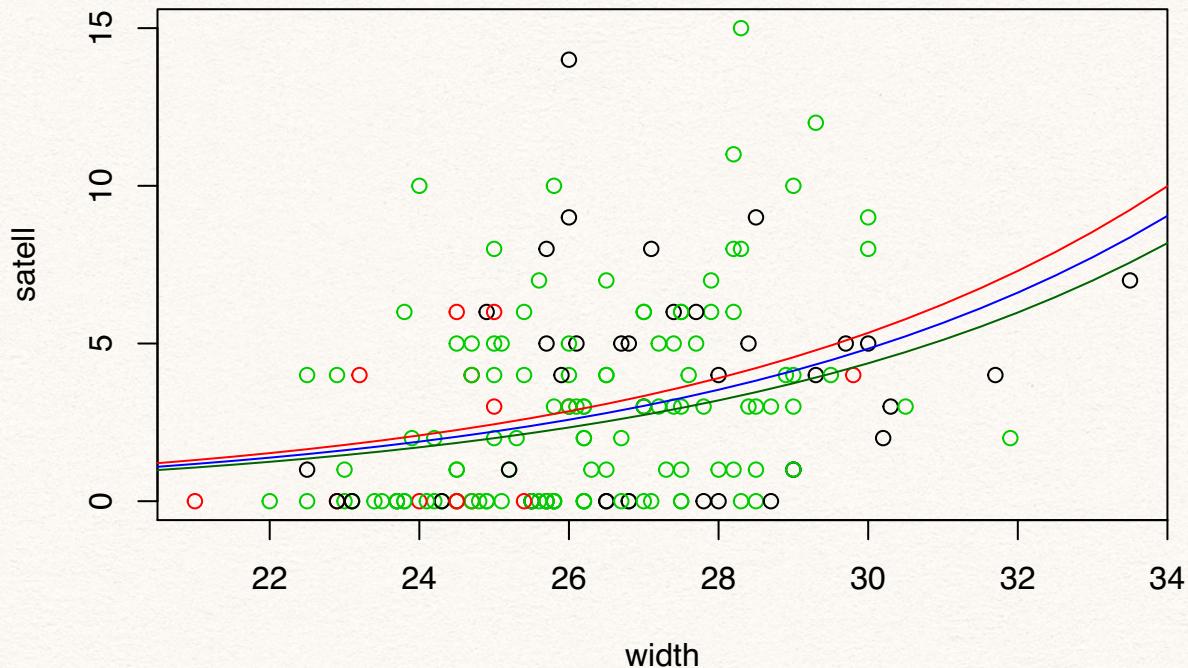
```

```

## factor(spine)3 -0.09899    0.10490   -0.944     0.345
## width           0.15668    0.02096    7.476 7.69e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
## Null deviance: 632.79  on 172  degrees of freedom
## Residual deviance: 566.60  on 169  degrees of freedom
## AIC: 929.9
##
## Number of Fisher Scoring iterations: 6

# plot both graph
plot(width,satell,col=factor(spine))
#  $g(x)=e^x$ 
x<- seq(from=20,to=34,by=0.5)
lines(x,exp(coef(poissonmod)[1]+coef(poissonmod)[4]*x),col="red")
lines(x,exp(coef(poissonmod)[1]+coef(poissonmod)[2]+coef(poissonmod)[4]*x),col="darkgreen")
lines(x,exp(coef(poissonmod)[1]+coef(poissonmod)[3]+coef(poissonmod)[4]*x),col="blue")

```



Comparing to the graph in part(a), they both have the increasing trend as width raising where part(a) are growing linearly while lines in this graph are growing exponentially. Better spines condition can still have a larger number of satell.

(d). Calculate the estimated number of satellites for a female crab with carapace width 28cm, and one worn or broken spine using (i) the linear model in part (a) and (ii) the Poisson glm in part (c).

Since the female crab has one worn or broken spin, then the spin condition is 2.

```

# Calculate the linear model
coef(fit.crabs)[1]+coef(fit.crabs)[2]+coef(fit.crabs)[4]*28

```

```

## (Intercept)
##      3.600611
# Calculate the poisson GLM model
exp(coef(poissonmod)[1]+coef(poissonmod)[2]+coef(poissonmod)[4]*28)

## (Intercept)
##      3.19653

```

The above result show that the estimated value for Linear model in part(a) will be  $3.601 \approx 4$  satellites and the estimated result for poisson glm model will be  $3.197 \approx 3$  satellites.

(e). Fit a glm to the data with only spine condition as a predictor, and (i) family=gaussian and (ii) family =poisson(link=identity). Compare the estimated coefficients, and explain why they are the same (you need to provide a proof by carrying out appropriate calculation(s)). Compare the estimated standard errors. Why do you think they differ (an explanation in words suffices)? Which standard errors do you think are more trustworthy?

```

canomod <- glm(satell ~ factor(spine), family=gaussian)
summary(canonmod)

##
## Call:
## glm(formula = satell ~ factor(spine), family = gaussian)
##
## Deviance Residuals:
##      Min        1Q     Median        3Q       Max
## -3.6486   -2.8099   -0.8099    2.0000   12.1901
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 3.6486    0.5154   7.079 3.66e-11 ***
## factor(spine)2 -1.6486    0.9597  -1.718  0.0876 .
## factor(spine)3 -0.8387    0.5890  -1.424  0.1563
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 9.829768)
##
## Null deviance: 1704.9 on 172 degrees of freedom
## Residual deviance: 1671.1 on 170 degrees of freedom
## AIC: 891.3
##
## Number of Fisher Scoring iterations: 2
identmod<- glm(satell ~ factor(spine), family=poisson(link=identity))
summary(identmod)

##
## Call:
## glm(formula = satell ~ factor(spine), family = poisson(link = identity))
##
## Deviance Residuals:
##      Min        1Q     Median        3Q       Max
## -2.7014   -2.3706   -0.5097    1.1252   5.0859

```

```

## 
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) 3.6486    0.3140 11.619 < 2e-16 ***
## factor(spine)2 -1.6486    0.4816 -3.423 0.000619 ***
## factor(spine)3 -0.8387    0.3490 -2.403 0.016265 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## (Dispersion parameter for poisson family taken to be 1)
## 
##     Null deviance: 632.79 on 172 degrees of freedom
## Residual deviance: 621.16 on 170 degrees of freedom
## AIC: 982.46
## 
## Number of Fisher Scoring iterations: 3

```

Proofing the same estimated coefficients:

For(i) family = Gaussian that is  $g^{-1}(x) = b'(x) = x$ :

$$\begin{aligned} g(E[Y|X]) &= \beta X \\ b'^{-1}(E[Y|X]) &= \beta X \\ E[Y|X] &= \beta X \end{aligned}$$

For(ii) family =poisson(link=identity) that is  $g(x) = x$ :

$$\begin{aligned} g(E[Y|X]) &= \beta X \\ E[Y|X] &= \beta X \end{aligned}$$

Also  $\beta$  can be calculated through the score function. First consider the Gaussian with identity link:

$$\begin{aligned} \frac{\delta l}{\delta \beta_j} &= \sum_{i=1}^n \frac{y_i - b'(\theta)}{a(\phi)b''(\theta)} * \frac{1}{b'(\theta_i)} * x_{ij} \\ \frac{\delta l}{\delta \beta_j} &= \sum_{i=1}^n \frac{y_i - \mu_i}{\sigma^2} * \frac{1}{\mu_i} * x_{ij} \\ \sum_{i=1}^n y_i x_{ij} &= \sum_{i=1}^n \mu_i x_{ij} \\ \sum_{i=1}^n y_i x_{ij} &= \sum_{i=1}^n X_i \beta_1 x_{ij} \\ X^T Y &= X^T X \beta_1 \\ \beta_1 &= (X^T X)^{-1} X^T Y \end{aligned}$$

Consider the poisson with identity link, then  $g'(x) = 1$ :

$$\begin{aligned} \frac{\delta l}{\delta \beta_j} &= \sum_{i=1}^n \frac{y_i - b'(\theta)}{a(\phi)b''(\theta)} * x_{ij} \\ \sum_{i=1}^n \frac{y_i - \lambda_i}{\lambda_i} * x_{ij} &= 0 \\ \sum_{i=1}^n y_i x_{ij} &= \sum_{i=1}^n \lambda_i x_{ij} \\ \beta_2 &= (X^T X)^{-1} X^T Y \end{aligned}$$

Then both of these 2 model will give the same estimation on  $\beta$ . Since both models are fitting the data by  $E[Y|X] = X\beta$ , so they obatain the same estimated coefficients.

For (i), we have that the dispersion parameter  $\phi = \sigma^2 = 9.830$  then the estimated standard error is 3.135 while (ii)'s dispersion parameter is 1 and  $var[Y|X] = b''(\theta)a(\phi) = e^\theta = E[Y|X]$ . Since they have different distriburion , so they have different estimated standard error.I think the standard error of (ii) is more trustworthy.