

2132 Midterm

Orly Olbum

3/10/2021

Problem 1

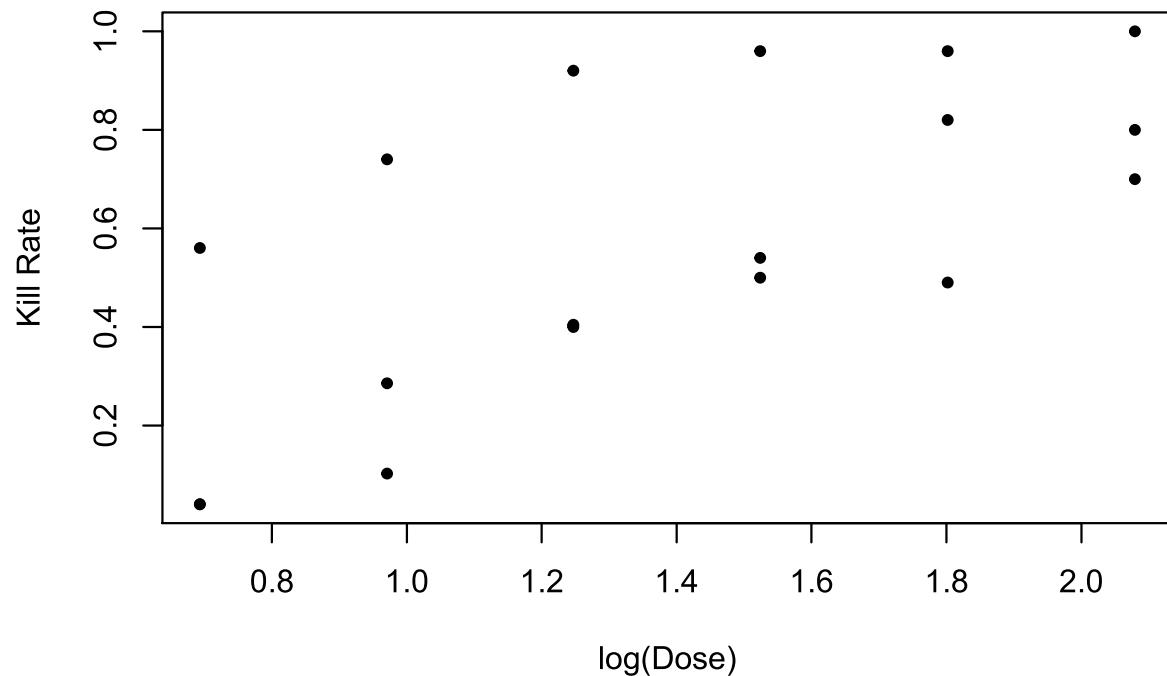
*Flour beetles *Tribolium castaneum* were sprayed with one of three insecticides in solution at different doses. The file “Insecticide.txt” contains data relating to the number of insects killed after a six-day period. The columns are:*

- Type: A factor variable with 3 levels giving the type of insecticide. 1:=DDT, 2:=gamma-BHC, 3:=DDT+gamma-BHC
- Dose: The dose of insecticide (in mg/10cm²).
- Trials: The total number of insects.
- Successes: The number of insects killed after a six-day period.

(a)

Investigate graphically the relationship between the dose, either in original units or in log units, and the kill rate. Here, kill rate is the probability an insect is killed.

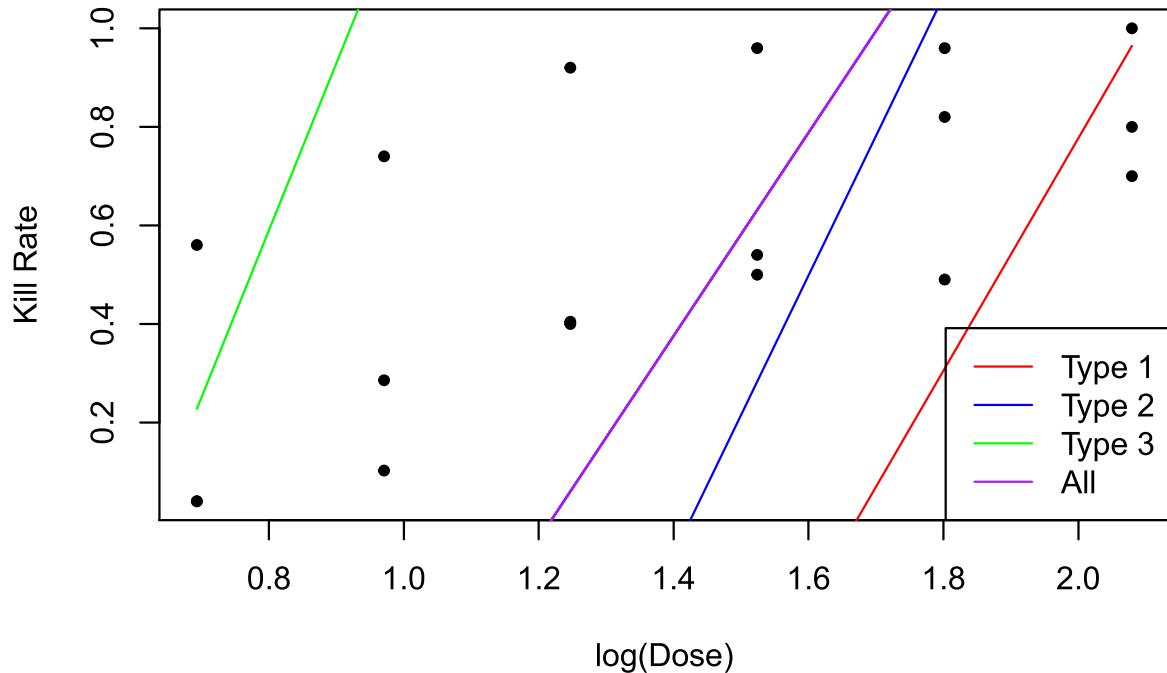
Log-Dose vs. Kill Rate



(b)

On the graph for part (a), plot the linear logistic fitted curve for each of the insecticides plus the combination.

Log-Dose vs. Kill Rate



(c)

Consider the two models, one in which the relationship is described by three parallel straight lines in the log dose and and one in which the three lines are straight but not parallel. Assess the evidence against the hypothesis of parallelism.

Parallel lines would indicate a good model fit when considering a proportional odds model, whereas if the lines are straight but not parallel, the slopes are different and this contradicts the proportional odds model fit. Fit1 does not have a significant coefficient for log-dose; Fit2 does not have a significant coefficient; Fit3 does not have a significant coefficient; the fit with all types considered has a significant coefficient at an alpha of .10. Although the lines are ~sort of~ parallel, since none of the individually fitted models have significant coefficients for log-dose, we may not be content with this model fit (proportional odds).

(d)

Let *Type* be a 3-level factor, and let *ldose* be the log dose. Explain the relationship between the regression coefficients in the model formulae '*Type* + *ldose*' and '*Type* + *ldose* - 1'.

```
##  
## Call:  
## glm(formula = killrate ~ Type + ldose, family = binomial(link = "logit"),  
##       data = insect)  
##
```

```

## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.28112 -0.11783 -0.05326  0.14433  0.33696
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) -4.4940    2.515   -1.787  0.0739 .
## Type2       -3.882    2.346   -1.654  0.0980 .
## Type3       -1.454    2.011   -0.723  0.4695
## ldose        2.725    1.522    1.790  0.0734 .
##
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 8.31082 on 17 degrees of freedom
## Residual deviance: 0.48343 on 14 degrees of freedom
## AIC: 18.27
##
## Number of Fisher Scoring iterations: 5

##
## Call:
## glm(formula = killrate ~ Type + ldose - 1, family = binomial(link = "logit"),
##      data = insect)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.28112 -0.11783 -0.05326  0.14433  0.33696
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## Type1     -4.494    2.515   -1.787  0.0739 .
## Type2     -3.882    2.346   -1.654  0.0980 .
## Type3     -1.454    2.011   -0.723  0.4695
## ldose      2.725    1.522    1.790  0.0734 .
##
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 8.66580 on 18 degrees of freedom
## Residual deviance: 0.48343 on 14 degrees of freedom
## AIC: 18.27
##
## Number of Fisher Scoring iterations: 5

```

In the second formula, we are omitting the intercept, which forces the model through the origin. We might do this if we are okay with (i.e., expecting to see) a zero, and in our dataset it is possible to have a 0 kill rate. When we leave in the intercept, we retain a degree of freedom because we code just two factors to indicate Type, whereas if we omit the intercept we lose a degree of freedom by adding another covariate as an indicator for the third Type (i.e., the intercept). However, we still have significant predictors in both models.

(e)

On the assumption that three parallel straight lines suffice, estimate the DDT+gamma-BHC dose required to give a 99% kill rate, and obtain a 90% confidence interval for this dose.

While the coefficients from above were not significant, the lines were ~sort of~ parallel, so we can use all of the data rather than just the Type 3 data to find a 90% CI for the designated dose.

The estimated log-dose required for 99% kill rate is 3.44. Since we are holding the assumption that the lines are parallel (sufficient for model assumptions), we use all of the data rather than just the third type of dose. The 90% CI is (1.44, 5.44).

Problem 2

MassSpec contains observed data, *Pep* contains expected data given *P* generated the spectrum in *MassSpec*. Let $n = \#\text{mass-to-charges in } \text{Pep}$; let $Y_r = 1$ if *P*'s r th mass-to-charge exists in observed spectrum (*MassSpec*) and $Y_r = 0$ if not for $r = 1, \dots, n$. Conditional on *P* having generated observed mass spectrum in *MassSpec*, assume Y_i 's are independent.

(a)

Suppose *P* generated the observed spectrum, and assume $\text{logit}\{E(Y_r)\} = B_0 + B_1 \log(x_r)$, where x_r is the expected relative intensity for mass-to-charge r (i.e. the r th entry of the second column in *Pep.txt*). Given what you know about the problem, what do you expect the sign of B_1 to be?

Based on the explanation of the problem and a rough understanding of the data, I would expect B_1 to be positive.

(b)

Compute the maximum likelihood estimate for B_1 .

The MLE estimate for B_1 is 198.2 with a standard error of 70482.6.

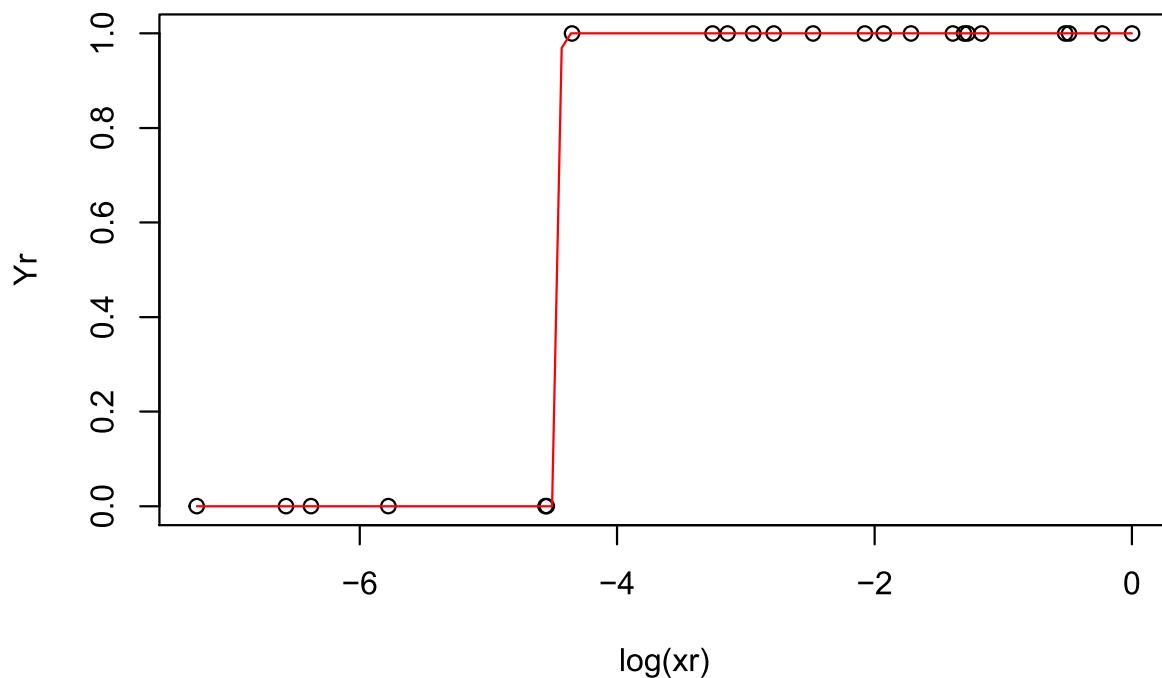
(i) Do you notice anything peculiar about the estimate?

The p-values associated with the coefficients are close to 1, and the coefficient values are extremely high with even more surprisingly high standard errors.

(ii) see attached

(iii) Plot Y_r as a function of $\log(x_r)$. In light of part (ii), are the observations you made in part (i) consistent with the plot? Explain.

Plot of Y_r as a function of $\log(xr)$



Considering Y_r can only take on 0 or 1, the plot does look as expected and seems like a good fit for the logit model and the observations from (i) are consistent.

Problem 2

- (b) (ii) Observe data $(z_1, w_1), \dots, (z_n, w_n)$
 where $z_i \in \{0, 1\}$, $w_i \in \mathbb{R}^p$
 $(p-1)$ -dim hyperplane $\gamma \in \mathbb{R}^p$ that perfectly
 partitions z_i

\rightarrow for $2z_i - 1 \in \{1, -1\}$,

$$(2z_i - 1)(w_i^\top \gamma) > 0, i=1, \dots, n$$

Show that MLE for η under
 $z_i \sim \text{Ber}(\pi_i)$, z_i ind.

$$\text{logit}(\pi_i) = w_i^\top \eta, \quad (=1, \dots, n).$$

has infinite form (that $\eta \rightarrow \infty$)

$$f(z_i) = \pi_i^{z_i} (1 - \pi_i)^{1-z_i}$$

$$\begin{aligned} \text{logit}(\pi_i) &= w_i^\top \eta \\ \rightarrow \pi_i &= \frac{\exp(w_i^\top \eta)}{1 + \exp(w_i^\top \eta)} \end{aligned}$$

$$(2z_i - 1)(w_i^\top \eta)$$

$$W = \left\{ 2(\pi_i^{z_i} (1 - \pi_i)^{1-z_i}) - 1 \right\} \cdot (w_i^\top \eta)$$

$$= 2(\pi_i^{z_i} (1 - \pi_i)^{1-z_i})(w_i^\top \eta) - (w_i^\top \eta)$$

For $(2z_i - 1)(w_i^\top \gamma) > 0$, we have

$f(z_i) \rightarrow 1$ which requires

$\pi_i = 1$ and therefore we need

$\eta \rightarrow \infty$ (infinite) because as $\eta \rightarrow \infty$,

$\pi_i \rightarrow 1$ and

$\text{logit}(\pi_i) \rightarrow 0$

$$\textcircled{C} \quad \beta = (\beta_0, \beta_1)^T$$

$$w_r = (1, \log(x_r))^T \in \mathbb{R}^2$$

$$Y = (Y_1, \dots, Y_n)^T \in \mathbb{R}^n$$

consider:

$$Y_{rl} | \Pi_r \sim \text{Ber}(\Pi_r), \quad Y_{il} | \Pi_i \text{ ind}$$

$$\text{logit}(\pi_r) | \beta = \omega_r^T \beta \quad , \quad r=1, \dots, n$$

$$\beta \sim N_2(\mu, \Sigma)$$

for some known constants $MER^2, \Sigma EIR^{2 \times 2}$
such that Σ is pos.

(i) Let $\text{pr}(\beta|Y_n)$ be pdf of posterior distribution

for $\beta \neq n$

Show that

$$\log \{ \Pr(\beta | Y_n) \} = l(\beta; Y_n) + g(\beta; \mu, \Sigma) + c(Y_n)$$

$\ell \rightarrow \log\text{-likelihood}$

$g \rightarrow$ function depends only on β, μ, Σ

\hookrightarrow function depends only on y_n

Derive expressions for λ_1 , λ_2

$$\text{Hint: Bayes' } \rightarrow \Pr(\beta|Y_n) = \frac{\Pr(Y_n|\beta) \Pr(\beta)}{\Pr(Y_n)}$$

$$\log \{ \Pr(\beta | Y_n) \} = \log \left(\frac{\Pr(Y_n | \beta) \Pr(\beta)}{\Pr(Y_n)} \right)$$

$$= \underbrace{\log(\Pr(Y_n | \beta))}_{\text{Term 1}} + \underbrace{\log(\Pr(\beta))}_{\text{Term 2}} - \underbrace{\log(\Pr(Y_n))}_{\text{Term 3}}$$

by definition, = l

only junction
of B

only
function
of y_n

(ii) Show that $\log \{ \Pr(\beta | Y_n) \}$ is concave

when treated as a function of p

say $L = \log \{\Pr(\beta | Y_n)\}$

$$\frac{\partial L}{\partial \beta} = 0 + g'(\beta; \mu, \sigma) - 0$$

$$\frac{\partial^2 L}{\partial \beta^2} = g''(\beta; \mu, \Sigma) > 0$$

So, L is concave

7

(iii) Any similarities between $\log \{\Pr(\beta | Y_n)\}$ and Ridge Regression from 2131?

It does seem similar \rightarrow in ridge regression, we considered a penalty when modeling and subsequently estimating β . $\rightarrow \hat{\beta}_{\text{ridge}} = (X^T X + \lambda I_p)^{-1} X^T Y$

which is similar to

$$\log \{\Pr(\beta | Y_n)\} = \ell + g + c$$

\downarrow penalty on γ
similar to $2\|\beta\|_2^2$ term

(iv) Instead of shrinking towards 0 like ridge regression, our "shrinkage estimator" for β , $g(\beta | \mu, \Sigma) = \log \{\Pr(\beta)\}$, will shrink towards μ .



(v)

$$\hat{\beta}^{(MAP)} = \underset{\beta \in \mathbb{R}^2}{\operatorname{argmax}} \log \{ \Pr(\beta | Y_n) \}$$

Show that $\hat{\beta}^{(MAP)}$ exists (ie, $\|\hat{\beta}^{(MAP)}\| < \infty$)
and use part (ii) to argue that it
is unique.

$$\begin{aligned} \hat{\beta}^{(MAP)} &= \underset{\beta \in \mathbb{R}^2}{\operatorname{argmax}} \frac{\Pr(\beta | Y_n)}{\Pr(Y_n | \beta) \Pr(\beta)} \\ &= \underset{\beta \in \mathbb{R}^2}{\operatorname{argmax}} \end{aligned}$$

and since $\Pr(Y_n)$ has nothing to do

with β_1 , we can drop it

$$\underset{\beta \in \mathbb{R}^2}{\operatorname{argmax}} \Pr(Y_n | \beta) \Pr(\beta)$$

$$\underset{\beta \in \mathbb{R}^2}{\operatorname{argmax}} \log \{ \Pr(Y_n | \beta) \} + \log \{ \Pr(\beta) \}$$

$$\underset{\beta \in \mathbb{R}^2}{\operatorname{argmax}} \sum_{i \in n} \log \{ \Pr(Y_i | \beta) \} + \log \{ \Pr(\beta) \}$$

$$Y_i | \beta \sim \text{Ber}(\beta) \rightarrow \Pr(Y_i | \beta) = \beta^{y_i} (1-\beta)^{1-y_i}$$

$$\begin{aligned} \sum \log \{ \Pr(Y_i | \beta) \} &= \sum \{ y_i \log \beta + (1-y_i) \log (1-\beta) \} \\ &= \sum y_i \log \beta + \sum (1-y_i) \log (1-\beta) \end{aligned}$$

$$\text{derivative wrt } \beta \rightarrow \frac{\sum y_i}{\beta} - \frac{2 - \sum y_i}{1-\beta} < \infty$$

$$\beta \sim N_2(\mu, \Sigma) \rightarrow \Pr(\beta) = \frac{1}{\sqrt{2\pi}^2} \exp \left\{ -\frac{1}{2\Sigma} (\beta - \mu)^2 \right\}$$

$$\log \{ \Pr(\beta) \} = -\log(\sqrt{2\pi}) - \log(\Sigma) - \frac{1}{2\Sigma} (\beta - \mu)^2$$

$$\text{derivative wrt } \beta \rightarrow -\frac{1}{2\Sigma} \cdot 2(\beta - \mu) = -\frac{(\beta - \mu)}{\Sigma}$$

so $\|\hat{\beta}^{(MAP)}\| < \infty$ and because the
second derivative of g (prior) is concave
from (ii), $\hat{\beta}^{(MAP)}$ is unique.

→

(vi) Find $\nabla_{\beta} \log \{ \text{pr}(\beta | Y_n) \}$ and $\nabla_{\beta}^2 \log \{ \text{pr}(\beta | Y_n) \}$, use to derive iterative algorithm to find $\hat{\beta}^{(\text{MAP})}$ that uses Newton-Raphson updates.

(ie, find vector function $U: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ and matrix function $H: \mathbb{R}^2 \rightarrow \mathbb{R}^{2 \times 2}$ such that

an iteration t can be written as

$$\beta^{(t)} = \beta^{(t-1)} + H \{ \beta^{(t-1)} \} U \{ \beta^{(t-1)} \}$$

$$\begin{aligned} \nabla_{\beta} \log \{ \text{pr}(\beta | Y_n) \} &= \nabla_{\beta} \log \{ \Pr(Y_n | \beta) \Pr(\beta) \} \\ &= \nabla_{\beta} \log \{ \Pr(Y_n | \beta) \} + \nabla_{\beta} \log \{ \Pr(\beta) \} \\ &= \frac{\sum y_i}{\beta} - \frac{2 - \sum y_i}{1 - \beta} - \frac{(\beta - \mu)}{\Sigma} \end{aligned}$$

$$\begin{aligned} \nabla_{\beta}^2 \log \{ \text{pr}(\beta | Y_n) \} &= \nabla_{\beta}^2 \log \{ \Pr(Y_n | \beta) \} + \nabla_{\beta}^2 \log \{ \Pr(\beta) \} \\ &= -\frac{\sum y_i}{\beta^2} - \frac{(1-\beta)(0) - (2-\sum y_i)(-1)}{(1-\beta)^2} + \frac{1}{\Sigma} \\ &= -\frac{\sum y_i}{\beta^2} + \frac{2 - \sum y_i}{(1-\beta)^2} - \frac{1}{\Sigma} \\ &= -\frac{1}{\beta} \cdot \frac{\sum y_i}{\beta} - \frac{1}{1-\beta} \cdot \frac{2 - \sum y_i}{1-\beta} - \frac{1}{\Sigma} \\ &= -\frac{1}{\Sigma} \cdot \frac{(\beta)(1-\beta)}{(\beta)(1-\beta)} - \frac{(1-\beta)}{\beta(1-\beta)} \cdot \frac{\sum y_i}{\beta} - \frac{\beta}{\beta(1-\beta)} \cdot \frac{2 - \sum y_i}{1-\beta} \\ &= \underbrace{\frac{1}{\beta(1-\beta)}}_{U \{ \beta^{(t-1)} \}} \left\{ -\frac{1}{\Sigma} (\beta)(1-\beta) - (1-\beta) \cdot \frac{\sum y_i}{\beta} - (\beta) \cdot \frac{2 - \sum y_i}{1-\beta} \right\} \underbrace{\beta^{(t-1)}}_{H \{ \beta^{(t-1)} \}} \end{aligned}$$

Problem 3

Consider the dataset “Fertility.txt”. Your goal is to understand which factors help determine the number of children a woman from Fiji will have in her lifetime. The data matrix consists of 5 columns: - YSFM: A factor variable giving the number of years since a woman’s first marriage. The levels are 1 (< 5 years), 2 (5-9 years), 3 (10-14 years), 4 (15-19 years), 5 (20-24 years), 6 (25+ years)

- Place: A factor variable with two levels (Urban and Rural), indicating if the mother lives in a urban or rural area.
- Education: A factor variable giving the mother’s education level. The levels are 1 (none), 2 (lower elementary), 3 (upper elementary), 4 (secondary or higher).
- Nwomen: The total number of women sampled.
- Average: The mean number of children born to a woman in that row.

(a)

Fit an appropriate model describing the how the number of children varies with marital age, mother’s abode and education. Explain the meaning of all parameters in your model, and comment on the major factors affecting fertility.

```
##  
## Call:  
## glm(formula = Average ~ YSFM + Place + Education, family = poisson(link = "log"),  
##       data = fert)  
##  
## Deviance Residuals:  
##      Min        1Q     Median        3Q       Max  
## -2.5176  -0.2105  -0.0239   0.1725   1.1070  
##  
## Coefficients:  
##             Estimate Std. Error z value Pr(>|z|)  
## (Intercept) 0.04765   0.39350   0.121  0.903625  
## YSFM2       0.99525   0.43022   2.313  0.020702 *  
## YSFM3       1.39972   0.41045   3.410  0.000649 ***  
## YSFM4       1.57115   0.40400   3.889  0.000101 ***  
## YSFM5       1.75061   0.39825   4.396  1.10e-05 ***  
## YSFM6       1.75178   0.39822   4.399  1.09e-05 ***  
## PlaceUrban  0.01840   0.14979   0.123  0.902251  
## Education2  0.04072   0.19602   0.208  0.835441  
## Education3 -0.06270   0.20119   -0.312  0.755317  
## Education4 -0.66436   0.24021   -2.766  0.005679 **  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ',' 1  
##  
## (Dispersion parameter for poisson family taken to be 1)  
##  
## Null deviance: 68.906  on 47  degrees of freedom  
## Residual deviance: 18.739  on 38  degrees of freedom  
## AIC: Inf  
##  
## Number of Fisher Scoring iterations: 5
```

Because we are looking at count data as a dependent variable, we’ll use a poisson glm model fit. YSFM is a factor variable with 6 variables. The model uses the first level as a baseline, and codes 5 others as

indicators for 1 if that level and all others are 0. If all 5 are 0, the intercept serves as the indicator for level 1. Place is a factor variable with Urban as 1 and Rural as 0. Education acts the same way as YSFM, with 4 different levels and if all are 0, the intercept acts as the indicator for level 1. The predictors with significant p-values (at 0.05) are YSFM3, YSFM4, YSFM5, YSFM6, and Education4. The intercept B0 is the log of the expected value of average children when all covariates are zero (or in this case, when a woman has level 1 YSFM, rural living, and level 1 Education). Exponential of the sum of all covariates is the expected value of average children for the covariates. Because we are using a log link, each of these coefficients represents the effect on the log of the average children a woman will have, not on the number of children itself.

(b)

Construct a 95% confidence interval for the mean number of children born to an urban woman with upper elementary education after ten years of marriage.

The 95% CI for mean number of children (4.07) with these levels is (2.58, 6.42).

(c)

Estimate the lifetime average number of children born to rural women with secondary education. Give 90% confidence limits.

Average children will be 1.13 children, with a confidence interval of (.73, 1.54). Since the highest level of YSFM has the most years since marriage, it is a good estimate for lifetime average across all women, based on the information provided.

```

## ----setup, include = FALSE, warning = FALSE, purl = TRUE-----
-----
knitr::opts_chunk$set(echo = FALSE)
setwd("C:/Users/orlyo/OneDrive/Desktop/Grad School/Spring 2021/2. STAT 2132 - Applied Stat.
Methods 2/Exams/Midterm")
# Problem 1
insect = read.csv("Insecticide.csv")
insect$killrate = insect$Successes/insect$Trials
insect$type = as.factor(insect$type)
insect$ldose = log(insect$dose)
# Problem 2
fert = read.csv("Fertility.csv")
fert$YSFM = as.factor(fert$YSFM)
fert$Place = as.factor(fert$Place)
fert$Education = as.factor(fert$Education)
# Problem 3
spec = read.csv("MassSpec.csv")
pep = read.csv("Pep.csv")
pep$Yr = as.numeric(pep$Exp_m.to.z %in% spec$Obs_m.to.z)
pep$lexp = log(pep$Exp_Int)

## ---- warning = FALSE-----
#
# plot(insect$dose, insect$killrate, main = "Dose vs. Kill Rate", xlab = "Dose", ylab = "Kill Rate")
plot(insect$ldose, insect$killrate, main = "Log-Dose vs. Kill Rate", xlab = "log(Dose)", ylab = "Kill Rate",
pch = 20)

## ---- warning = FALSE-----
#
fit1 = glm(killrate ~ ldose, family = binomial(link = "logit"), data = insect[insect$type == 1,])
fit2 = glm(killrate ~ ldose, family = binomial(link = "logit"), data = insect[insect$type == 2,])
fit3 = glm(killrate ~ ldose, family = binomial(link = "logit"), data = insect[insect$type == 3,])
fit_all = glm(killrate ~ ldose, family = binomial(link = "logit"), data = insect)
# summary(fit1)
# summary(fit2)
# summary(fit3)
# summary(fit_all)

# lines(insect$ldose[insect$type == 1], fit1$fitted.values, col = "red", pch = 20)
# lines(insect$ldose[insect$type == 2], fit2$fitted.values, col = "blue", pch = 20)
# lines(insect$ldose[insect$type == 3], fit3$fitted.values, col = "green", pch = 20)
# lines(insect$ldose, fit_all$fitted.values, col = "purple", pch = 20)
# legend("bottomright", legend = c("Type 1", "Type 2", "Type 3", "All"), col = c("red", "blue", "green",
"purple"), lty = 1)

```

```

plot(insect$Idose, insect$killrate, main = "Log-Dose vs. Kill Rate", xlab = "log(Dose)", ylab = "Kill Rate",
pch = 20)
lines(insect$Idose[insect$type == 1], log(fit1$fitted.values/(1 - fit1$fitted.values)), col = "red", pch = 20)
lines(insect$Idose[insect$type == 2], log(fit2$fitted.values/(1 - fit2$fitted.values)), col = "blue", pch = 20)
lines(insect$Idose[insect$type == 3], log(fit3$fitted.values/(1 - fit3$fitted.values)), col = "green", pch = 20)
lines(insect$Idose, log(fit_all$fitted.values/(1 - fit_all$fitted.values)), col = "purple", pch = 20)
legend("bottomright", legend = c("Type 1", "Type 2", "Type 3", "All"), col = c("red", "blue", "green",
"purple"), lty = 1)

## ---- warning = FALSE-----
-----
form1 = glm(killrate ~ Type + Idose, data = insect, family = binomial(link = "logit"))
summary(form1)
form2 = glm(killrate ~ Type + Idose - 1, data = insect, family = binomial(link = "logit"))
summary(form2)

## ---- warning = FALSE-----
-----
# fit3 = glm(killrate ~ Idose, family = binomial(link = "logit"), data = insect[insect$type == 3,])
fit_all = glm(killrate ~ Idose, family = binomial(link = "logit"), data = insect)
log.dose = (log(.99/(1-.99)) - summary(fit_all)$coefficients[1,1])/summary(fit_all)$coefficients[2,1]
# log.dose

lo = log.dose - 1.65*summary(fit_all)$coefficients[2, 2]
hi = log.dose + 1.65*summary(fit_all)$coefficients[2, 2]
ci = c(lo, hi)
# log.dose; ci

## ---- include = FALSE, warning = FALSE-----
-----
par(mfrow = c(1, 2))
plot(spec$Obs_m.to.z, spec$Obs_Int, type = "h", main = "observed", xlab = "mass/charge", ylab =
"Relative intensity")
plot(pep$Exp_m.to.z, pep$Exp_Int, type = "h", main = "expected", xlab = "mass/charge", ylab = "relative
intensity")

## ---- warning = FALSE-----
-----
n = sum(pep$Yr)
model = glm(Yr ~ log(Exp_Int), data = pep, family = binomial(link = "logit"))
# summary(model)

```

```

## ---- warning = FALSE-----
-----
par(mfrow = c(1, 1))
plot(pep$lexp, pep$Yr, main = "Plot of Yr as a funciton of log(xr)", xlab = "log(xr)", ylab = "Yr")
g = glm(Yr ~ lexp, family = binomial, pep)
curve(predict(g, data.frame(lexp = x), type = "resp"), add = TRUE, col = "red")

## ---- warning = FALSE-----
-----
model1 = glm(Average ~ YSFM + Place + Education, data = fert, family = poisson(link = "log"))
summary(model1)

## ---- warning = FALSE-----
-----
# exp(confint(model1))
log.mean = 0.04765 + 1.399525 + 0.01840 - 0.06270
new.mean = exp(log.mean)

new.data = data.frame(YSFM = "3", Place = "Urban", Education = "3")
pred = predict(model1, newdata = new.data, se.fit = TRUE, interval = "confidence")

upper = pred$fit + 1.96*pred$se.fit
lower = pred$fit - 1.96*pred$se.fit
interval = c(exp(lower), exp(upper))
# interval

## ---- warning = FALSE-----
-----
new = data.frame(YSFM = "6", Place = "Rural", Education = "4")
pred = predict(model1, newdata = new, se.fit = TRUE, interval = "confidence")
ci = c(pred$fit - 1.65*pred$se, pred$fit + 1.65*pred$se)

```