

# Homework 8

Qianbo Wang  
uni: qw2180

## Problem 1

**Reading Assignment:** Ramsey and Schafer (2nd Ed), Chapter 22 (Log-Linear Regression for Poisson Counts)

## Problem 2

Consider the Valve characteristics data (Display 22.16, Ramsey and Schafer, 2nd Ed).

Valve Failure in Nuclear Reactors. Display 22.16 shows characteristics and numbers of failures observed in valve types from one pressurized water reactor. There are five explanatory factors:

- system (1 = containment, 2 = nuclear, 3 = power conversion, 4 = safety, 5 = process auxiliary);
- operator type (1 = air, 2 = solenoid, 3 = motor-driven, 4 = manual);
- valve type (1 = ball, 2 = butterfly, 3 = diaphragm, 4 = gate, 5 = globe, 6 = directional control);
- head size (1 = less than 2 inches, 2 = 210 inches, 3 = 1030 inches);
- operation mode (1 = normally closed, 2 = normally open).

The lengths of observation periods are quite different, as indicated in the last column, time. Using an offset for log of observation time, identify the factors associated with large numbers of valve failures.

**(a). Do Problem Number 24, Page 667, using the R function glm.**

First, use all of the variables to fit a log-linear poisson regression model with  $\log(\text{time})$  as offset. And the result is as follows:

### Log-linear Poisson Regression

```
Call: glm(formula = Failures ~ System + Operator + Valve + Size + Mode,
  family = "poisson", data = Nuclear, offset = log(Time))
```

Coefficients:

(Intercept)	System2	System3	System4	System5	Operator2	Operator3	Operator4	Valve2
-3.7687	0.9156	1.0188	1.2231	0.3329	0.7044	-1.1926	-2.4723	0.1853
Valve3	Valve4	Valve5	Valve6	Size2	Size3	Mode2		
0.6067	2.9589	1.7932	1.0089	-0.0122	1.6146	-0.2093		

Degrees of Freedom: 89 Total (i.e. Null); 74 Residual

Null Deviance: 386

Residual Deviance: 196 AIC: 332

Call:

```
glm(formula = Failures ~ System + Operator + Valve + Size + Mode,
  family = "poisson", data = Nuclear, offset = log(Time))
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-3.189	-1.007	-0.436	0.336	5.314

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	-3.7687	0.8194	-4.60	4.2e-06	***
System2	0.9156	0.5318	1.72	0.0852	.
System3	1.0188	0.5055	2.02	0.0439	*
System4	1.2231	0.5552	2.20	0.0276	*
System5	0.3329	0.5841	0.57	0.5687	
Operator2	0.7044	0.5667	1.24	0.2139	
Operator3	-1.1926	0.2485	-4.80	1.6e-06	***
Operator4	-2.4723	0.4766	-5.19	2.1e-07	***
Valve2	0.1853	0.7611	0.24	0.8076	
Valve3	0.6067	0.7811	0.78	0.4373	
Valve4	2.9589	0.6001	4.93	8.2e-07	***
Valve5	1.7932	0.6104	2.94	0.0033	**
Valve6	1.0089	0.9301	1.08	0.2780	
Size2	-0.0122	0.2834	-0.04	0.9657	
Size3	1.6146	0.3210	5.03	4.9e-07	***
Mode2	-0.2093	0.1903	-1.10	0.2714	

---

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

(Dispersion parameter for poisson family taken to be 1)

Null deviance: 385.53 on 89 degrees of freedom  
 Residual deviance: 195.68 on 74 degrees of freedom  
 AIC: 332

Number of Fisher Scoring iterations: 7

So the model is:

$$\begin{aligned}
 \log(\text{Failures}) = & -3.77 + 0.92\text{System2} + 1.02\text{System3} + 1.22\text{System4} + 0.33\text{System5} \\
 & + 0.70\text{Operator2} - 1.19\text{Operator3} - 2.47\text{Operator4} \\
 & + 0.19\text{Valve2} + 0.61\text{Valve3} + 2.96\text{Valve4} + 1.79\text{Valve5} + 1.01\text{Valve6} \\
 & - 0.01\text{Size2} + 1.1\text{Size3} \\
 & - 0.21\text{Mode2}
 \end{aligned}$$

And since there are several coefficients not significant, this means that these corresponding variables are not associated with Failures.

**(b). Interpret the estimated parameters.**

There are total 16 parameters here in the model.  $\beta_0$  as the intercept, i.e., when System, Operator, Valve, Size and Mode all equal 1, the  $= e^{-3.77}$ .  $\beta_1 - \beta_4$ , as the coefficients of 4 different System levels 2-5, respectively.  $\beta_5 - \beta_7$ , as the coefficients of 3 different Operator levels 2-4, respectively.  $\beta_8 - \beta_{12}$ , as the coefficients of 5 different Valve levels 2-6, respectively.  $\beta_{13} - \beta_{14}$ , as the coefficients of 2 different levels Size 2-3, respectively and  $\beta_{15}$ , as the coefficient of Mode level 2.

- System

When other explanatory variables are fixed, then System 2 will cause  $e^{0.92} = 2.509$  times as many Failures as System 1, and similarly System 3 will cause  $e^{1.02} = 2.773$  times as many Failures as System 1, System 4 will cause  $e^{1.22} = 3.387$  times as many Failures as System 1. System 5 will cause  $e^{0.33} = 1.391$  times as many Failures as System 1. And since under 0.05 significance level, the coefficients of System 2 and System 5 are not significant, then this means that System 2 and System 5 are not different from System 1 in the sense of Large number of Failures.

- Operator

When other explanatory variables are fixed, then Operator 2 will cause  $e^{0.70} = 2.014$  times as many Failures as Operator 1, and similarly Operator 3 will cause  $e^{-1.19} = 0.304$  times as many Failures as System 1, Operator

4 will cause  $e^{-2.47} = 0.085$  times as many Failures as Operator 1. And since under 0.05 significance level, the coefficient of Operator 2 is not significant, then this means that Operator 2 is not different from System 1 in the sense of Large number of Failures.

- Valve

When other explanatory variables are fixed, then Valve 2 will cause  $e^{0.19} = 1.20$  times as many Failures as Valve 1, and similarly Valve 3 will cause  $e^{0.61} = 1.84$  times as many Failures as Valve 1, Valve 4 will cause  $e^{2.96} = 19.298$  times as many Failures as Valve 1, Valve 5 will cause  $e^{1.79} = 5.990$  times as many Failures as Valve 1, Valve 6 will cause  $e^{1.01} = 2.746$  times as many Failures as Valve 1. And since under 0.05 significance level, the coefficients of Valve 2, Valve 3 and Valve 6 are not significant, then this means that Valve 2, Valve 3 and Valve 6 are not different from System 1 in the sense of Large number of Failures.

- Size

When other explanatory variables are fixed, then Size 2 will cause  $e^{-0.01} = 0.990$  times as many Failures as Size 1, and similarly Size 3 will cause  $e^{1.61} = 5.002$  times as many Failures as Size 1. And since under 0.05 significance level, the coefficient of Size 2 is not significant, then this means that Size 2 is not different from Size 1 in the sense of Large number of Failures.

- Mode

When other explanatory variables are fixed, then Mode 2 will cause  $e^{-0.20} = 0.819$  times as many Failures as Mode 1. And since under 0.05 significance level, the coefficient of Mode 2 is not significant, then this means that Mode 2 is not different from Mode 1 in the sense of Large number of Failures.

### (c). Assess the goodness of fit of the model

#### Goodness of fit test Result

##### Goodness of Fit Test

```
res.deviance  df  p.value
195.67809951  74  1.00000000
```

From result of Deviance Goodness of Fit Test, the p-value is approximate to 1 which means that either the model is adequate or that insufficient data are available to detect inadequacies.

And the Chi-square Anova Table of log-linear regression result is as follows:

#### Log-Linear ANOVA Table

##### Analysis of Deviance Table

Model: poisson, link: log

Response: Failures

Terms added sequentially (first to last)

	Df	Deviance	Resid.	Df	Resid.	Dev	Pr(>Chi)
NULL				89		386	
System	4	22.7		85		363	0.00015 ***
Operator	3	5.3		82		357	0.14882
Valve	5	109.9		77		248	< 2e-16 ***
Size	2	50.7		75		197	9.6e-12 ***
Mode	1	1.2		74		196	0.27084
---							
Signif. codes:	0	***	0.001	**	0.01	*	0.05 . 0.1 1

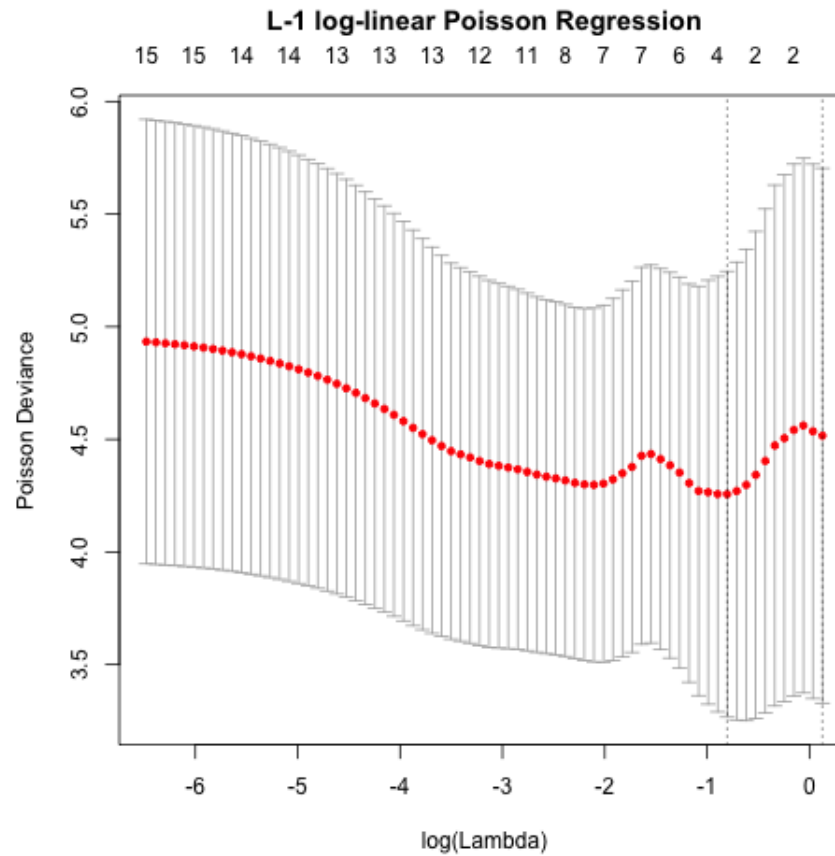
Since under 0.05 significance level, the Operator and Mode of explanatory variables are not significant. Then we can conclude that they are not strongly associated with Failures.

### Problem 3

Repeat 2(a) using the glmnet package and comment on the results.

Use cross validation and log linear poisson regression on the data.

The following plot shows the best lambda:



So the best lambda is 0.1644754 and the corresponding coefficients are:

Coefficients of log-linear poisson regression with lasso

16 x 1 sparse Matrix of class "dgCMatrix"

```

1
(Intercept) -1.42488
System2      .
System3      .
System4      .
System5      .
Operator2    .
Operator3    .
Operator4    .
Valve2       .
Valve3       .
Valve4       0.73104
Valve5       .
Valve6       .
Size2        .
Size3        0.55313
Mode2        .

```

And the model is:

$$\log(Failures) = -1.4249 + 0.7310Valve4 + 0.5531Size3$$

**Compare the results:**

Since from the two models we can find that, only Valve 4 and Size 3 are included in the model. Then this means only Valve 4 and Size 3 have significant effect on Failures. And compare the two models, we can find that these two models are quite different. And using Lasso with cross validation really depends on the seed because it depends on the training dataset. And lasso will reduce many variables, whereas, regular log-linear regression won't.

**R Code:**

```

rm(list=ls())
library(Sleuth2)
data(ex2224)
Nuclear<-ex2224
levels(Nuclear$System)<-c("containment","nuclear","power conversion","safety","process auxiliary")
levels(Nuclear$Operator)<-c("air","solenoid","motor-driven","manual")
levels(Nuclear$Valve)<-c("ball","butterfly","diaphragm","gate","globe","directional control")
levels(Nuclear$Size)<-c("<2","2-10","10-30")
levels(Nuclear$Mode)<-c("closed","open")
for (names in names(Nuclear)){
  if (class(Nuclear[[names]])!="numeric"){
    Nuclear[[names]]<-as.factor(as.numeric(Nuclear[[names]]))
  }
}

#loglinear poisson regression
glmPoisson<-glm(Failures~System+Operator+Valve+Size+Mode,offset=log(Time),data=Nuclear,family="poisson")
summary(glmPoisson)

sink('/Users/ramond/Drive/STAT W4201/HW8/glmPoisson.txt')
glmPoisson
summary(glmPoisson)
sink()

# goodness of fit test
Goodness_of_Fit<-cbind(res.deviance = sprintf("%.8f",glmPoisson$deviance),df = glmPoisson$df.residual, p.value = sprintf(
  "%.8f",(1-pchisq(glmPoisson$deviance,glmPoisson$df.residual,lower.tail = FALSE))))

sink('/Users/ramond/Drive/STAT W4201/HW8/glmGoodness.txt')
cat("Goodness of Fit Test \n \n")
write.table(Goodness_of_Fit,row.names = FALSE,quote = FALSE,sep=" ")
sink()

sink('/Users/ramond/Drive/STAT W4201/HW8/glmChisq.txt')
anova(glmPoisson,test="Chisq")
sink()

#lasso log linear poisson regression
#data matrix transformation of dummy variables
library(glmnet)

NuclearData<-subset(Nuclear,select=c("Failures","Time"))
for (names in names(Nuclear)[1:5]){
  for (factor in levels(Nuclear[[names]])[-1]){
    NuclearData[[paste(names,factor,sep="")]]<-as.numeric(Nuclear[[names]]==factor)
  }
}
NuclearMat<-as.matrix(NuclearData[,-c(1,2)])
set.seed(123)
glmPoisson.l1<-glmnet(NuclearMat,NuclearData$Failures,offset=log(NuclearData$Time),family="poisson")
glmPoisson.l1.cv<-cv.glmnet(NuclearMat,NuclearData$Failures,offset=log(NuclearData$Time),family="poisson")

#plot the best lambda
png(filename = "/Users/ramond/Drive/STAT W4201/HW8/glmPoisson.l1.cv.png")
plot(glmPoisson.l1.cv)
title(main = "L-1 log-linear Poisson Regression", line = 2.5)
dev.off()

#get model coefficients
lambda<-glmPoisson.l1.cv$lambda.min
model<-glmPoisson.l1.cv$glmnet.fit
coeff<-coef(model,lambda)
sink('/Users/ramond/Drive/STAT W4201/HW8/glmPoisson.l1.cv.txt')
coeff
sink()

```