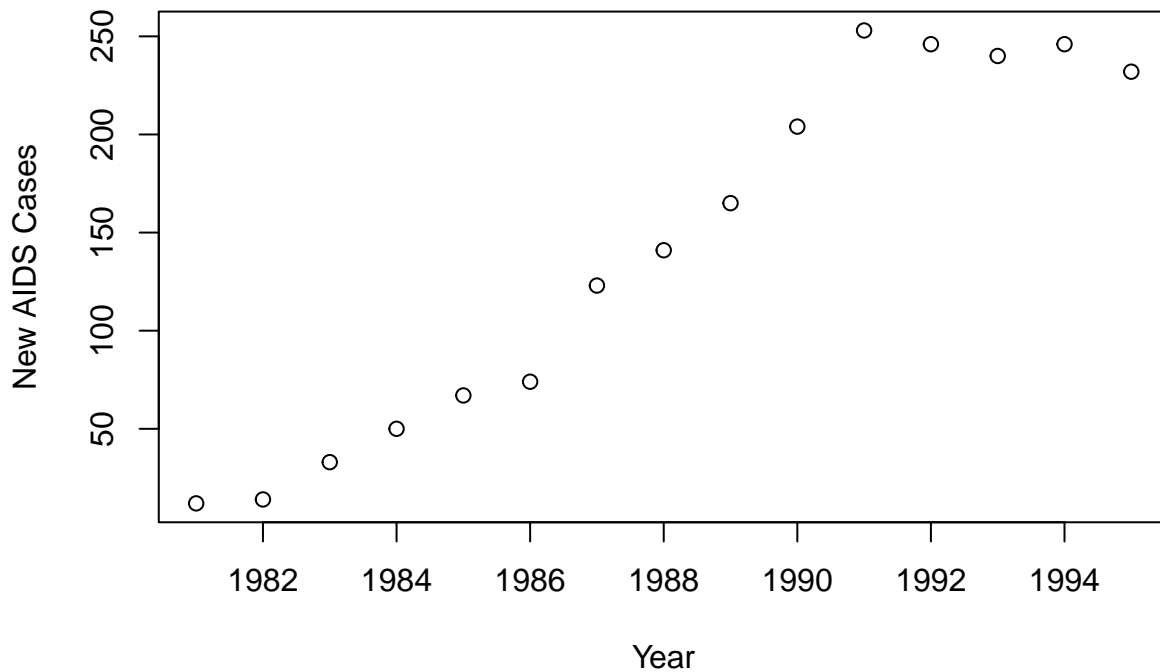# STAC51A3

Yulun Wu, Raymond Chan

2023-03-05

## Q2

```
AIDs = data.frame(Year=c(1981,1982,1983,1984,1985,1986,1987,1988,1989,1990,
1991,1992,1993,1994,1995),
AIDS=c(12,14, 33, 50, 67, 74,123, 141, 165, 204, 253, 246, 240, 246, 232))
```

### (a)

```
plot(AIDs$Year,AIDs$AIDS,main="New AIDS Cases v.s. Year",ylab="New AIDS Cases",xlab="Year")
```



### (b)

```
AIDs$t=AIDs$Year-1980
model = glm(AIDS~t,family=poisson,data=AIDs)
summary(model)
```
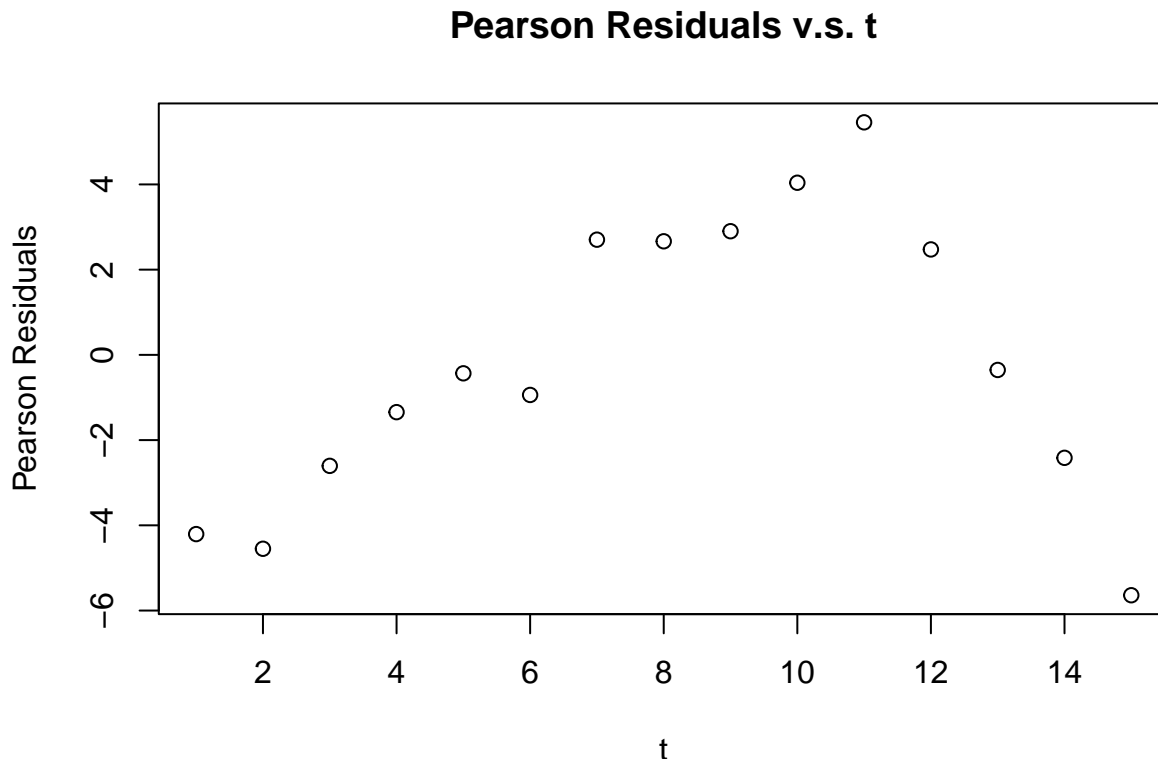
```
##
```

```
## Call:
## glm(formula = AIDS ~ t, family = poisson, data = AIDs)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -5.9751  -2.6345  -0.4367   2.5776   5.1378
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) 3.478884   0.064975   53.54   <2e-16 ***
## t           0.155739   0.005735   27.16   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##     Null deviance: 1010.27  on 14  degrees of freedom
## Residual deviance:  173.33  on 13  degrees of freedom
## AIC: 273.65
##
## Number of Fisher Scoring iterations: 4
```

The estimated model is $log(\mu(t)) = 3.478884 + 0.155739t$, standard error for $\alpha$ is 0.064975, p-value for $\alpha$ is <2e-16; standard error for $\beta$ is 0.005735, p-value for $\beta$ is <2e-16, thus both coefficients are statistically significant.

**(c)**

```
Pres = residuals(model,type = "pearson")
plot(AIDs$t,Pres,xlab="t",ylab="Pearson Residuals",main="Pearson Residuals v.s. t")
```

## Pearson Residuals v.s. t

Since when t = 1, 2, 10, 11, 15, |Pearson residual| > 3, we consider these data points as potential outliers.

**(d)**

```
t_bar = mean(AIDs$t)
t_st = AIDs$t-t_bar
model2 = glm(AIDS~t_st+I(t_st^2),family=poisson,data=AIDs)
summary(model2)
```

```
##
## Call:
## glm(formula = AIDS ~ t_st + I(t_st^2), family = poisson, data = AIDs)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.45122  -0.54143  0.03733  0.56349  1.54168
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)  4.982168   0.032545  153.08   <2e-16 ***
## t_st         0.214565   0.008816   24.34   <2e-16 ***
## I(t_st^2)   -0.021221   0.001775  -11.96   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##     Null deviance: 1010.2722  on 14  degrees of freedom
## Residual deviance:    9.2446  on 12  degrees of freedom
## AIC: 111.56
##
## Number of Fisher Scoring iterations: 4
```

The estimated model is $log(\mu(t)) = 4.982168 + 0.214565(t - \bar{t}) - 0.021221(t - \bar{t})^2$, standard error for $\alpha$ is 0.032545, p-value for $\alpha$ is <2e-16; standard error for $\beta_1$ is 0.008816, p-value for $\beta_1$ is <2e-16; standard error for $\beta_2$ is 0.001775, p-value for $\beta_2$ is <2e-16, thus all 3 coefficients are statistically significant.

**(e)**

```
anova(glm(AIDS~t_st,family=poisson,data=AIDs),model2,test="Chisq")
```

```
## Analysis of Deviance Table
##
## Model 1: AIDS ~ t_st
## Model 2: AIDS ~ t_st + I(t_st^2)
##   Resid. Df Resid. Dev Df Deviance  Pr(>Chi)
## 1        13    173.335
## 2        12      9.245  1   164.09 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
LRT_ts = 173.335-9.245
pchisq(LRT_ts,df=1,lower.tail=F)
```

```
## [1] 1.445775e-37
```

P-value = 1.445775e-37 < 0.05, so we reject the null hypothesis that $\beta_2 = 0$, this means that the quadratic Poisson model fits the data better. This is agreed with by comparing the residual standard deviance of both.

## Q3

### (a)

```
crab = read.table("horseshoecrabs.dat", header=TRUE)
model = glm(Satellites~Weight,family=poisson,data=crab)
summary(model)
```

```
##
## Call:
## glm(formula = Satellites ~ Weight, family = poisson, data = crab)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -2.9307  -1.9981  -0.5627   0.9298   4.9992
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) -0.42841    0.17893  -2.394   0.0167 *
## Weight       0.58930    0.06502   9.064   <2e-16 ***
## ---
## 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: 560.87  on 171  degrees of freedom
## AIC: 920.16
##
## Number of Fisher Scoring iterations: 5
```

$log(\mu(x)) = -0.42841 + 0.58930x$.

### (b)

```
CI_l = 0.58930-qnorm(1-0.05/2,0,1)*0.06502
CI_u = 0.58930+qnorm(1-0.05/2,0,1)*0.06502
cbind(CI_l,CI_u)
```

```
##           CI_l      CI_u
## [1,] 0.4618631 0.7167369
```

```
exp(cbind(CI_l,CI_u))
```

```
##          CI_l     CI_u
## [1,] 1.587028 2.04774
```

A 95% Wald confidence interval for $\beta$ is (0.4618631,0.7167369). Take exponent of this CI, we get (1.587028,2.04774). The expected log count for each 1 satellites increase in number of satellites given weight is from 0.4618631 and to 0.7167369. There is a 95% confidence that the multiplicative effect on odds of a 1kg increase on weight will cause from 1.587028 and to 2.04774 increase in in number of satellites.

**(c)**

```
anova(glm(Satellites~1,family=poisson,data=crab),model,test="Chisq")
```

```
## Analysis of Deviance Table
##
## Model 1: Satellites ~ 1
## Model 2: Satellites ~ Weight
##   Resid. Df Resid. Dev Df Deviance  Pr(>Chi)
## 1       172     632.79
## 2       171     560.87  1   71.925 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```
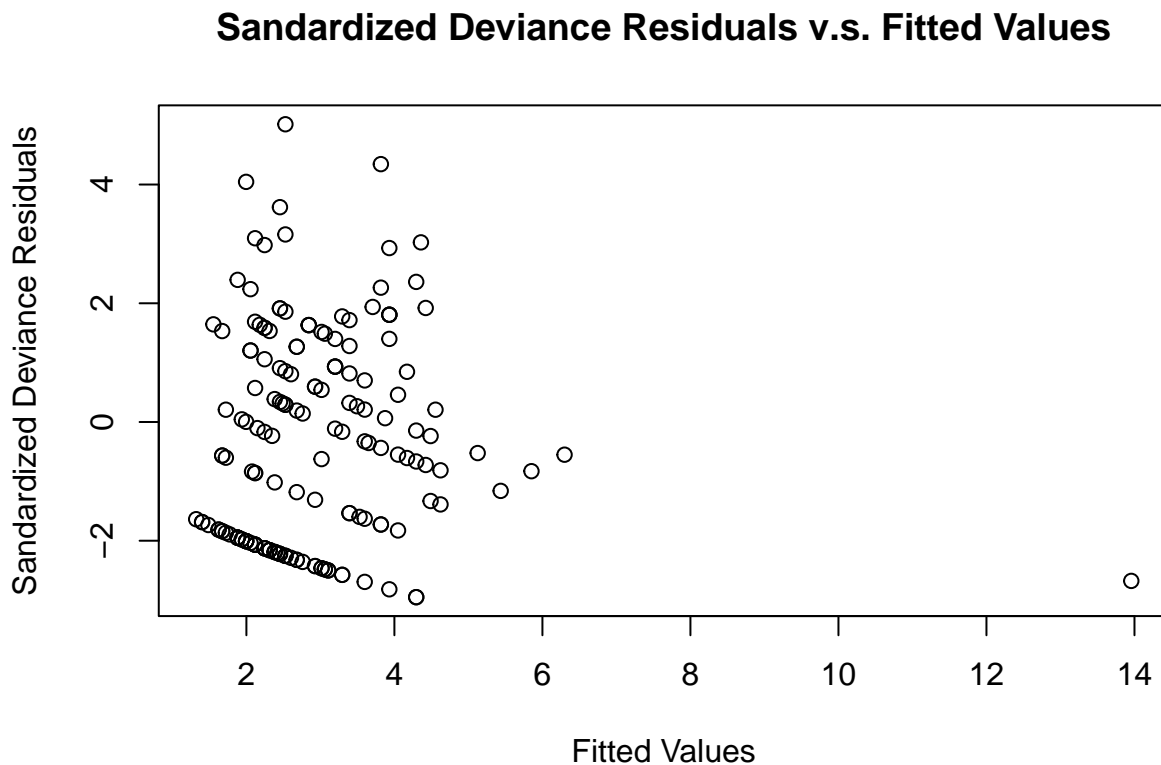
```
LRT_ts = 632.79-560.87
pchisq(LRT_ts,df=1,lower.tail=F)
```

```
## [1] 2.24101e-17
```

The residual deviance is given by 560.87 with 171 degrees of freedom, the null deviance with 632.79 and 172 degrees of freedom, p-value = 2.24101e-17 < 0.05, so we reject the null hypothesis that weight effect doesn't exist. We can conclude based on this test that our model fits the data well.

**(d)**

```
sdr = rstandard(model)
crab_fitted = model$fitted.values
plot(crab_fitted,sdr,ylab="Sandardized Deviance Residuals",xlab="Fitted Values",main="Sandardized Deviar
```

## Sandardized Deviance Residuals v.s. Fitted Values



We have a couple outliers which has |sandardized deviance residual| > 3., we also have a fitted value of 14 on one data point which may be considered an different type of outlier since its standardized deviance residual

looks close to the other data.

## (e)

```
model$deviance
```

```
## [1] 560.8664
```

```
model$df.residual
```

```
## [1] 171
```

Deviance of the model is 560.8664 and degree of freedom is 171. We cannot use Goodness-of-Fit test because the data is not equally dispersed (mean,variance are equal). A deviance of the model being lower indicates that there is a better model fit, we may compare to other models in Q4.

# Q4

## (a)

```
library(MASS)
# Give interval a label from 1 to 11
crab$level = cut(crab$Weight, breaks = c(0,seq(1.5,3.3,0.2),Inf),labels = 1:11)
# Original interval
crab$interval = cut(crab$Weight, breaks = c(0,seq(1.5,3.3,0.2),Inf))
# Function to calculate information needed in each row of table
cal_info = function(level){
  ncase = round(sum(as.numeric(crab$level==level)),0)
  selected_row_y = crab[crab$level==level,"Satellites"]
  mu = round(mean(selected_row_y),2)
  sigma = round(var(selected_row_y),2)
  return(cbind(ncase,mu,sigma))
}
info_matrix = mapply(cal_info,1:11) # result is a 3 by 11 matrix
crab_info_df = as.data.frame(cbind(unique(crab$interval)[order(unique(crab$interval))],t(info_matrix)))
colnames(crab_info_df) = c("Weight (kg)","No. Cases","Sample Mean","Sample Variance")
crab_info_df$"Weight (kg)" = unique(crab$interval)[order(unique(crab$interval))] # Do this beacause cbi
crab_info_df
```

```
##     Weight (kg) No. Cases Sample Mean Sample Variance
## 1       (0,1.5]         5        0.80            3.20
## 2     (1.5,1.7]        11        0.82            1.56
## 3     (1.7,1.9]        16        1.25            7.93
## 4     (1.9,2.1]        25        2.56            7.01
## 5     (2.1,2.3]        29        2.83           12.86
## 6     (2.3,2.5]        13        3.00            6.33
## 7     (2.5,2.7]        20        2.55            6.05
## 8     (2.7,2.9]        17        3.06            5.93
## 9     (2.9,3.1]        17        5.59           17.01
## 10    (3.1,3.3]        13        4.69           13.56
## 11    (3.3,Inf]         7        4.00            2.67
```

We can see from that table that when the interval is (1.7,1.9], (1.9,2.1], (2.1,2.3], (2.9,3.1], (3.1,3.3], sample means are quit different from sample variances, so there is evidence of overdispersion.

## (b)

```
model.nb = glm.nb(Satellites~Weight,data=crab)
summary(model.nb)
```

```
##
## Call:
## glm.nb(formula = Satellites ~ Weight, data = crab, init.theta = 0.9310592338,
##     link = log)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -1.8394  -1.4122  -0.3247   0.4744   2.1279
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -0.8647     0.4048  -2.136   0.0327 *
## Weight        0.7603     0.1578   4.817 1.45e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(0.9311) family taken to be 1)
##
##     Null deviance: 216.43  on 172  degrees of freedom
## Residual deviance: 196.16  on 171  degrees of freedom
## AIC: 754.64
##
## Number of Fisher Scoring iterations: 1
##
##
##             Theta:  0.931
##          Std. Err.:  0.168
##
##  2 x log-likelihood:  -748.644
```

```
model.nb$theta # theta
```

```
## [1] 0.9310592
```

```
model.nb$SE.theta # SE of theta
```

```
## [1] 0.1675963
```

We have negative binomial log linear defined with $log(\mu(x)) = -0.8647 + 0.7603x$, $\theta = 0.9310592$, SE of $\theta$ is 0.1675963. The dispersion parameter for theta is .931 with a SE of 0.168.

## (c)

```
CI_l = 0.7603-qnorm(1-0.05/2,0,1)*0.1578
CI_u = 0.7603+qnorm(1-0.05/2,0,1)*0.1578
cbind(CI_l,CI_u) # beta CI
```
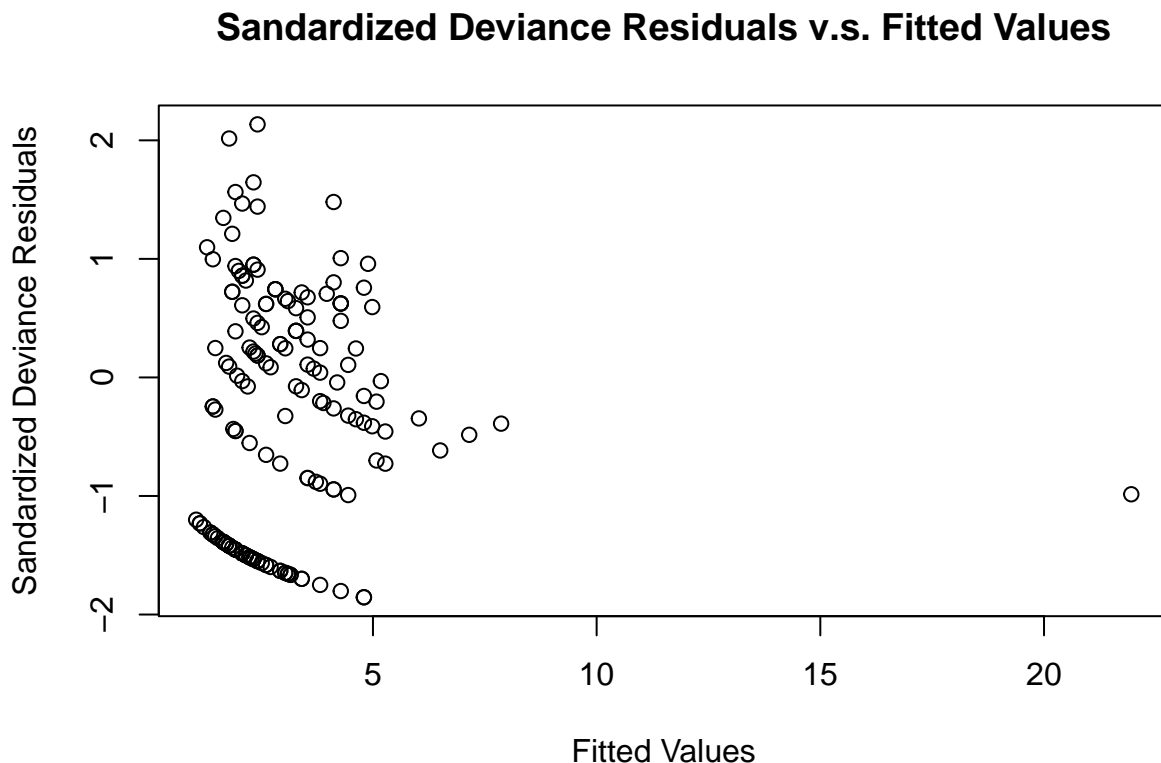
```
##           CI_l     CI_u
## [1,] 0.4510177 1.069582
```

```
exp(cbind(CI_l,CI_u)) # exp(beta) CI
```

```
##             CI_l     CI_u
## [1,] 1.569909 2.914162
```

The 95% Wald confidence interval for $exp(\beta)$ in Q3 (b) is (1.587028,2.04774). The 95% Wald confidence interval for $exp(\beta)$ with negative binomial model is (1.569909,2.914162). The interval is wider with the negative binomial model because SE is under estimated in poisson model. I.e. the negative binomial accounts for the overdispersion of the data, while the poisson does not.

**(d)**

```
sdr = rstandard(model.nb) # Standardized deviance residuals
crab_fitted = model.nb$fit
plot(crab_fitted,sdr,ylab="Sandardized Deviance Residuals",xlab="Fitted Values",main="Sandardized Devian
```



**Sandardized Deviance Residuals v.s. Fitted Values**

There is no outlier if we use |Sandardized Deviance Residual| > 3 as standard for detecting outlier. If we use |Sandardized Deviance Residual| > 2 as standard for detecting outlier, there is 1 outlier. The large fitted value is now deviates less from the model and looks better

# Q5

**(a)**

```
PayYes = c(24,10,5,16,7,47,45, 57,54,59)
PayNo = c(9,3,4,7,4,12,8,9,10,12)
District = rep(c("NC", "NE", "NW", "SE", "SW"),2)
Race = c(rep("Blacks",5), rep("Whites",5))
q5_df = data.frame(Race, District, PayYes, PayNo)
model.logit = glm(cbind(PayYes,PayNo)~Race+District,family=binomial,data=q5_df)
summary(model.logit)
```

```
##
## Call:
## glm(formula = cbind(PayYes, PayNo) ~ Race + District, family = binomial,
##     data = q5_df)
##
## Deviance Residuals:
##        1          2          3          4          5          6          7          8
##  0.60191    0.30311   -0.97042   -0.09608   -0.30707   -0.53319   -0.18583    0.47422
##        9         10
##  0.07216    0.15054
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.74947    0.29581   2.534  0.01129 *
## RaceWhites   0.79129    0.28532   2.773  0.00555 **
## DistrictNE   0.25837    0.42067   0.614  0.53909
## DistrictNW   0.13836    0.40517   0.341  0.73273
## DistrictSE   0.12087    0.37287   0.324  0.74581
## DistrictSW   0.00445    0.38486   0.012  0.99077
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 10.665  on 9  degrees of freedom
## Residual deviance:  2.071  on 4  degrees of freedom
## AIC: 49.437
##
## Number of Fisher Scoring iterations: 4
```

```r
# Wald test
p_value_Wald = 0.00555 # copy from summary

# LRT
drop1(model.logit,test="Chisq")
```

```
## Single term deletions
##
## Model:
## cbind(PayYes, PayNo) ~ Race + District
##         Df Deviance    AIC    LRT Pr(>Chi)
## <none>       2.0710 49.437
## Race      1  9.4624 54.828 7.3915 0.006553 **
## District  4  2.5876 41.953 0.5167 0.971859
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```r
p_value_LRT = 0.006553 # copy from drop1
```

P-value for Wald test is 0.00555, P-value for Likelihood Ratio test is 0.006553, so we reject the null hypothesis that the merit pay increase is independent of race, when holding the district unchanged in two test with $\alpha = 0.05$.

**(b)**

```
common_or = exp(0.79129) # Estimated common odd ratio, which is exp(beta_1)
common_or
```

```
## [1] 2.206241
```

```
# 95% Wald CI for exp(beta_1)
CI_l_w = 0.79129-qnorm(1-0.05/2,0,1)*0.28532
CI_u_w = 0.79129+qnorm(1-0.05/2,0,1)*0.28532
exp(cbind(CI_l_w,CI_u_w)) # exp(beta_1) CI
```

```
##        CI_l_w   CI_u_w
## [1,] 1.261212 3.859381
```

```
# 95% LR CI for exp(beta_1)
exp(confint(model.logit))
```

```
## Waiting for profiling to be done...
```

```
##                  2.5 %   97.5 %
## (Intercept) 1.2015339 3.853896
## RaceWhites  1.2519737 3.844984
## DistrictNE  0.5755038 3.033447
## DistrictNW  0.5224551 2.585144
## DistrictSE  0.5437179 2.363742
## DistrictSW  0.4725999 2.154622
```

The 95% Wald confidence interval for the common odds ratio between Merit Pay and Race is (1.261212,3.859381). Since 1 is not in Ci, Merit Pay and Race have significant association. We are 95% confident that the odds of White people got a merit pay increase is 1.261212 to 3.859381 factor odds of Black people got a merit pay increase.

The 95% likelihood ratio confidence interval for the common odds ratio between Merit Pay and Race is (1.2519737,3.844984). Since 1 is not in Ci, Merit Pay and Race have significant association. We are 95% confident that the odds of White people got a merit pay increase is 1.2519737 to 3.844984 factor odds of Black people got a merit pay increase.

**(c)**

```
# Fit model with interaction
model.int = glm(cbind(PayYes,PayNo)~Race*District,family=binomial,data=q5_df)
anova(model.logit,model.int,test="Chisq")
```

```
## Analysis of Deviance Table
##
## Model 1: cbind(PayYes, PayNo) ~ Race + District
## Model 2: cbind(PayYes, PayNo) ~ Race * District
##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1         4      2.071
## 2         0      0.000  4    2.071   0.7227
```

```
test_stat = 2.071-0
test_stat
```

```
## [1] 2.071
```

```
df=4-0
df
```

```
## [1] 4
```

```
p_value=0.7227
p_value
```

```
## [1] 0.7227
```

Since p-value = 0.7227>0.05, we fail to reject the null hypothesis that the simpler model (model without interaction) fits the data well, which means there is homogeneous association between merit pay decision and race across the five districts. Test statistic is 2.071, degrees of freedom is 4.

# Q6

## (a)

```
length = rep(c("short", "short","long","long"),3)
height = rep(c("low","high"),6)
time = c(rep(c("early"),4),rep(c("mid"),4),rep(c("late"),4))
y = c(54,44,25,18,77,63,64,21,22,25,13,5)
n = c(67,49,43,19,98,67,97,26,36,38,24,10)
lizards = data.frame(length, height, time, y,n)

model.full = glm(cbind(y,n-y)~length*height*time,family = binomial) # 3-way interaction model
anova(glm(cbind(y,n-y)~length*height+height*time+length*time,family = binomial),model.full,test="Chisq")
```

```
## Analysis of Deviance Table
##
## Model 1: cbind(y, n - y) ~ length * height + height * time + length *
##     time
## Model 2: cbind(y, n - y) ~ length * height * time
##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1         2     3.6927
## 2         0     0.0000  2   3.6927   0.1578
```

I will remove the 3-way interaction term. because it has highest order in the full model and 2-way interaction model fits data well according to p-value = 0.1578 > 0.05 from LRT of 2-way interaction model v.s. 3-way interaction model, so we fail to reject the null hypothesis that 2-way interaction model fits data well and hence can drop the 3-way interaction term.

## (b)

```
step(model.full,test="Chisq")
```

```
## Start:  AIC=68.21
## cbind(y, n - y) ~ length * height * time
##
##                       Df Deviance    AIC    LRT Pr(>Chi)
## - length:height:time   2   3.6927 67.904 3.6927   0.1578
## <none>                     0.0000 68.212
##
## Step:  AIC=67.9
## cbind(y, n - y) ~ length + height + time + length:height + length:time +
```

```
##      height:time
##
##                Df Deviance    AIC    LRT Pr(>Chi)
## - length:time   2   4.1858 64.397 0.4931  0.78149
## - length:height 1   3.6940 65.906 0.0013  0.97118
## <none>              3.6927 67.904
## - height:time   2   8.5909 68.802 4.8982  0.08637 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=64.4
## cbind(y, n - y) ~ length + height + time + length:height + height:time
##
##                Df Deviance    AIC    LRT Pr(>Chi)
## - length:height 1   4.2121 62.424 0.0263  0.87110
## <none>              4.1858 64.397
## - height:time   2   9.6543 65.866 5.4685  0.06494 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Step:  AIC=62.42
## cbind(y, n - y) ~ length + height + time + height:time
##
##              Df Deviance    AIC     LRT  Pr(>Chi)
## <none>            4.2121 62.424
## - height:time 2   9.8815 64.093  5.6694 0.0587351 .
## - length      1  15.2549 71.467 11.0428 0.0008903 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

##
## Call:  glm(formula = cbind(y, n - y) ~ length + height + time + height:time,
##     family = binomial)
##
## Coefficients:
##       (Intercept)        lengthshort          heightlow          timelate
##            1.8823             0.6859            -1.3407           -1.9070
##           timemid  heightlow:timelate   heightlow:timemid
##           -0.1027             1.2986             0.2025
##
## Degrees of Freedom: 11 Total (i.e. Null);  5 Residual
## Null Deviance:      54.04
## Residual Deviance: 4.212     AIC: 62.42
```

```r
model.final = glm(cbind(y,n-y)~length + height + time + height:time,family = binomial) # Final model
summary(model.final)
```

```
##
## Call:
## glm(formula = cbind(y, n - y) ~ length + height + time + height:time,
##     family = binomial)
##
## Deviance Residuals:
##        1         2         3         4         5         6         7         8
##  0.64895  -0.79146  -0.68477   1.14640  -0.11340   0.58998   0.09797  -0.66784
```

```
##         9        10        11        12
## -0.48553  -0.02114   0.57197   0.03907
##
## Coefficients:
##                    Estimate Std. Error z value Pr(>|z|)
## (Intercept)          1.8823     0.4470   4.211 2.54e-05 ***
## lengthshort          0.6859     0.2069   3.315 0.000918 ***
## heightlow           -1.3407     0.4803  -2.792 0.005245 **
## timelate            -1.9070     0.5258  -3.627 0.000287 ***
## timemid             -0.1027     0.5557  -0.185 0.853433
## heightlow:timelate   1.2986     0.6265   2.073 0.038179 *
## heightlow:timemid    0.2025     0.6177   0.328 0.743091
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 54.0430  on 11  degrees of freedom
## Residual deviance:  4.2121  on  5  degrees of freedom
## AIC: 62.424
##
## Number of Fisher Scoring iterations: 4
```

```
exp(model.final$coefficients) # exp(coefficients of final model)
```

```
##      (Intercept)        lengthshort          heightlow           timelate
##        6.5684972          1.9854995          0.2616652          0.1485257
##          timemid heightlow:timelate  heightlow:timemid
##        0.9024382          3.6642726          1.2244145
```

Due to the result of backward elimination, the best model is $logit(\pi) = \alpha + \beta_{short}c_{short} + \beta_{low}c_{low} + \beta_{late}c_{late} + \beta_{mid}c_{mid} + \gamma_{low\&late}c_{low}c_{late} + \gamma_{low\&mid}c_{low}c_{mid} = 1.8823 + 0.6859c_{short} - 1.3407c_{low} - 1.9070c_{late} - 0.1027c_{mid} + 1.2986c_{low}c_{late} + 0.2025c_{low}c_{mid}$.

Final stepped model has a lower AIC, 62.42 compared to 68.212 which indicates better fit. Lower is better for AIC, when comparing models. For the stepped model, more degrees of freedom of 5 indicates less overfitting, with a residual deviance of 4.212 compared to 3.6927 for the full model.

Interpretation:

$e^{\beta_{short}}$: Estimated odds of a lizard caught at a site will be grahami for observing short length lizard are 1.9854995 times the estimated odds for observing long length lizard of the same height and time of day.

$e^{\beta_{low}}$: Estimated odds of a lizard caught at a site will be grahami for observing low height lizard at early time of day are 0.2616652 times the estimated odds for observing high height lizard at early time of day of the same length.

$e^{\beta_{late}}$: Estimated odds of a lizard caught at a site will be grahami for observing high height lizard at late time of day are 0.1485257 times the estimated odds for observing high height lizard at early time of day of the same length.

$e^{\beta_{mid}}$: Estimated odds of a lizard caught at a site will be grahami for observing high height lizard at mid time of day are 0.9024382 times the estimated odds for observing high height lizard at early time of day of the same length.

$e^{\beta_{low}+\gamma_{low\&late}}$: Estimated odds of a lizard caught at a site will be grahami for observing low height lizard at late time of day are 0.2616652*3.6642726=0.9588126 times the estimated odds for observing high height lizard at late time of day of the same length.

13

$e^{\beta_{low}+\gamma_{low\&mid}}$: Estimated odds of a lizard caught at a site will be grahami for observing low height lizard at mid time of day are 0.2616652*1.2244145=0.3203867 times the estimated odds for observing high height lizard at mid time of day of the same length.

$e^{\beta_{late}+\gamma_{low\&late}}$: Estimated odds of a lizard caught at a site will be grahami for observing low height lizard at late time of day are 0.1485257*3.6642726=0.5442387 times the estimated odds for observing low height lizard at early time of day of the same length.

$e^{\beta_{mid}+\gamma_{low\&mid}}$: Estimated odds of a lizard caught at a site will be grahami for observing low height lizard at mid time of day are 0.9024382*1.2244145=1.104958 times the estimated odds for observing low height lizard at early time of day of the same length.

## (c)

```
model.reduced = glm(formula = cbind(y, n - y) ~ length + height + time+ height:time, family = binomial)
anova(model.reduced,model.full,test="Chisq")
```

```
## Analysis of Deviance Table
##
## Model 1: cbind(y, n - y) ~ length + height + time + height:time
## Model 2: cbind(y, n - y) ~ length * height * time
##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1         5     4.2121
## 2         0     0.0000  5   4.2121   0.5193
```

Based on the the goodness-of-fit test the p-value is $0.5193 > 0.05$, so we fail to reject the null hypothesis that my reduced model fits data as well as saturated model, hence my reduced model is better than the saturated model since it has less parameters.