

Homework 3: Binary and Count Regressions

Solution Set

This homework set will cover problems concerning binary and count regression models. Point totals for specific problems are given, 10 points will be reserved for correct submission of the homework assignment.

Problem 1 [10 points]: Manually write your own Fisher scoring algorithm which maximizes the logistic regression log likelihood for the CCSO example in the notes. Report $\hat{\beta}$ and reproduce the summary table (up to convergence tolerance differences) without using the `glm` or `summary` commands. You can ignore deviance residuals.

Solution 1:

The log-likelihood is

$$\begin{aligned}l(\beta) &= \sum_{i=1}^n y_i x_i^T \beta - \log(1 + \exp(x_i^T \beta)) \\l'(\beta) &= \sum_{i=1}^n \left(y_i x_i - \frac{x_i}{1 + \exp(x_i^T \beta)} \exp(x_i^T \beta) \right) \\&\implies l'(\beta) = X^T (Y - \pi)\end{aligned}$$

where $\pi_i = \frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)}$

$$\begin{aligned}l''(\beta) &= - \sum_{i=1}^n \left(\frac{x_i^2}{(1 + \exp(x_i^T \beta))^2} \exp(x_i^T \beta) \right) \\&\implies l''(\beta) = -X^T W X\end{aligned}$$

where $W = \text{diag}(\pi_i(1 - \pi_i))$

Thus the Fisher scoring algorithm:

$$\beta^{(t+1)} = \beta^{(t)} + (X^T W^{(t)} X)^{-1} X^T (Y - \pi)$$

```
# Reading in the data
library(data.table)
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.2 --
## v ggplot2 3.4.0      v purrr 0.3.5
## v tibble 3.1.8       v dplyr 1.0.10
## v tidyr 1.2.1        v stringr 1.4.1
## v readr 2.1.3        v forcats 0.5.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::between()   masks data.table::between()
## x dplyr::filter()    masks stats::filter()
## x dplyr::first()     masks data.table::first()
## x dplyr::lag()       masks stats::lag()
## x dplyr::last()      masks data.table::last()
## x purrr::transpose() masks data.table::transpose()

library(MASS)

##
## Attaching package: 'MASS'
##
## The following object is masked from 'package:dplyr':
##
##      select

CCSO = fread("https://uofi.box.com/shared/static/9elozjsg99bgcb7gb546wlfr3r2gc9b7.csv")

# Pre-Processing the data
CCSO_small <- CCSO %>% rename(Days = "Days in Jail", Age = "Age at Arrest",
Date = "BOOKING DATE", Sex = "SEX", Race = "RACE",
Crime = "CRIME CODE") %>%
mutate(atleastone = ifelse(Days > 0,1,0)) %>%
filter(Crime == "OTHER TRAFFIC OFFENSES") %>%
filter(Race %in% c("Asian/Pacific Islander","Black","White","Hispanic")) %>%
filter(Sex %in% c("Female","Male")) %>%
dplyr::select(atleastone, Age, Sex, Date, Race) %>%
mutate(Race = fct_drop(Race), Sex = fct_drop(Sex))
CCSO_small <- CCSO_small[complete.cases(CCSO_small), ]
head(CCSO_small)

##      atleastone Age      Sex      Date      Race
## 1:             0  22    Male 1/1/2011    White
## 2:             0  26    Male 1/1/2011    White
## 3:             0  32 Female 1/1/2011    White
## 4:             0  22    Male 1/2/2011    White
## 5:             0  35    Male 1/2/2011 Hispanic
## 6:             0  35    Male 1/2/2011 Hispanic

#Creating the model matrix
X = model.matrix(atleastone ~ -1 + Race + Sex + Age,data = CCSO_small)
n = nrow(X)
p = ncol(X)
Y = CCSO_small$atleastone

# Initializing the beta
beta = matrix(rep(0,6))

#Running the Fisher scoring iterations
for(t in 1:10)
```

```
{
  pi = exp(X%%beta)/(1+exp(X%%beta))
  W = diag(c(pi*(1-pi)))
  beta = beta + solve(t(X) %%% W %%% X)%*%t(X)%*(Y - pi)
}

# The final values
pi_CCS0 = exp(X%%beta)/(1+exp(X%%beta))
W_CCS0 = diag(c(pi_CCS0*(1-pi_CCS0)))
var_matrix_CCS0 = solve(t(X) %%% W_CCS0 %%% X)
sd_beta = sqrt(diag(var_matrix_CCS0))
z_Val = beta/sd_beta
pvalue = 2*(1 - pnorm(abs(z_Val)))

#Deviance results
deviance_res = -2*(t(Y)%*X%%beta - sum(log(1+exp(X%%beta))))
beta_0 = 0
deviance_null = -2*(beta_0*sum(Y) - n*log(1+exp(beta_0)))
AIC = deviance_res + 2*p
```

The summary table:

```
tab = data.frame("Estimate" = beta,"Std.Error" = sd_beta,"z value" = z_Val, "Pvalue" = pvalue)
list("Coefficients" = tab, "Null deviation" = deviance_null, "Residual Deviance" = deviance_res, "Null deviance" = deviance_null)

## $Coefficients
##              Estimate Std. Error z.value Pvalue
## RaceAsian/Pacific Islander -4.38986549 0.523613419 -8.383791 0.000000e+00
## RaceBlack -1.87654964 0.144600944 -12.977437 0.000000e+00
## RaceHispanic -2.80454861 0.173349092 -16.178617 0.000000e+00
## RaceWhite -3.04322627 0.147160127 -20.679693 0.000000e+00
## SexMale 0.73983403 0.105379772 7.020646 2.208456e-12
## Age 0.00770504 0.003186262 2.418207 1.559721e-02
##
## $`Null deviation`
## [1] 8201.317
##
## $`Residual Deviance`
## [1] 4668.728
##
## $`Null df`
## [1] 5916
##
## $`Residual df`
## [1] 5910
##
## $AIC
## [1] 4680.728
```

Problem 2 [10 points]: Manually write your own Fisher scoring algorithm which maximizes the Poisson regression log likelihood for the Galapagos example in the notes. Report $\hat{\beta}$ and reproduce the summary

table (up to convergence tolerance differences) without using the `glm` or `summary` commands. You can ignore deviance residuals.

Solution 2:

The log-likelihood is

$$l(\beta) = \sum_{i=1}^n (y_i x_i^T \beta - \exp(x_i^T \beta))$$

$$l'(\beta) = \sum_{i=1}^n (y_i x_i - x_i \exp(x_i^T \beta))$$

$$\Rightarrow l'(\beta) = X^T(Y - \pi)$$

where $\pi_i = \exp(x_i^T \beta)$

$$l''(\beta) = - \sum_{i=1}^n x_i^2 \exp(x_i^T \beta)$$

$$\Rightarrow l''(\beta) = -X^T W X$$

where $W = \text{diag}(\pi_i)$

Thus the Fisher scoring algorithm:

$$\beta^{(t+1)} = \beta^{(t)} + (X^T W^{(t)} X)^{-1} X^T (Y - \pi)$$

```
library(faraway)
```

```
#Pre-processing the data
```

```
data(gala)
```

```
gala <- gala %>%
```

```
mutate(Size = as.factor(1 + ifelse(Area > 1,1,0) + ifelse(Area > 25,1,0)))
```

```
head(gala)
```

```
##           Species Endemics  Area Elevation Nearest Scrutz Adjacent Size
## Baltra          58        23 25.09        346      0.6   0.6      1.84    3
## Bartolome       31        21  1.24         109      0.6  26.3    572.33    2
## Caldwell        3         3  0.21          114      2.8  58.7      0.78    1
## Champion       25         9  0.10           46      1.9  47.4      0.18    1
## Coamano         2         1  0.05           77      1.9   1.9    903.82    1
## Daphne.Major    18        11  0.34          119      8.0   8.0      1.84    1
```

```
# Creating the model matrix
```

```
X = model.matrix(Species ~ Elevation + Nearest + Scrutz + Adjacent + Size, data = gala)
```

```
n = nrow(X)
```

```
p = ncol(X)
```

```
Y = gala$Species
```

```
# Initialising the beta
```

```
beta = c(0, apply(X, 2, mean)[-1]/apply(X, 2, var)[-1])
```

```
#Running the fisher scoring iterations
```

```
for(t in 1:10)
```

```
{
```

```
  pi = exp(X%*%beta)
```

```
  W = diag(c(pi))
```

```
  beta = beta + solve(t(X) %*% W %*% X)%*%t(X)%*(Y - pi)
```

```
}
```

```
# The final values
```

```
pi = exp(X%*%beta)
```

```
W = diag(c(pi))
```

```
var_matrix = solve(t(X) %*% W %*% X)
```

```
sd_beta = sqrt(diag(var_matrix))
```

```
z_Val = beta/sd_beta
```

```
pvalue = 2*(1 - pnorm(abs(z_Val)))
```

```
#Deviance results
```

```
deviance_res = -2*(t(Y)%*%X%*%beta - sum(pi)+ sum(Y - Y*log(Y)))
```

```
beta_0 = log(mean(Y))
```

```
deviance_null = -2*(beta_0*sum(Y) - n*exp(beta_0) + sum(Y - Y*log(Y)))
```

```
AIC = -2*sum(dpois(Y,pi,log = TRUE)) + 2*(p)
```

The summary table:

```
tab = data.frame("Estimate" = beta,"Std.Error" = sd_beta,"z value" = z_Val, "Pvalue" = pvalue)
```

```
list("Coefficients" = tab, "Null deviation" = deviance_null, "Residual Deviance" = deviance_res, "Null d
```

```
## $Coefficients
```

```
##           Estimate      Std.Error    z.value      Pvalue
```

```
## (Intercept)  2.7897964692 8.107802e-02 34.408787 0.0000000000
```

```
## Elevation    0.0009360990 5.402069e-05 17.328527 0.0000000000
```

```
## Nearest      0.0064693041 1.747557e-03  3.701912 0.0002139805
```

```
## Scrutz      -0.0062664946 6.268336e-04 -9.997063 0.0000000000
```

```
## Adjacent    -0.0002857805 2.960795e-05 -9.652152 0.0000000000
```

```
## Size2       1.1276155415 9.535272e-02 11.825730 0.0000000000
```

```
## Size3       2.0586771315 9.419392e-02 21.855732 0.0000000000
```

```
##
```

```
## $`Null deviation`
```

```
## [1] 3510.729
```

```
##
```

```
## $`Residual Deviance`
```

```
##           [,1]
```

```
## [1,] 594.1753
```

```
##
```

```
## $`Null df`
```

```
## [1] 29
```

```
##
```

```
## $`Residual df`
```

```
## [1] 23
```

```
##
```

```
## $AIC
```

```
## [1] 769.0063
```

Problem 3 [5 points]: Derive the log-odds ratio of $x + 1$ to x when $Y = 1$, and observe that the log-odds ratio does not depend on x . Comment on this finding.

Solution 3:

We know that in the logistic regression model we have

$$\text{logit}(\pi(x)) = x\beta \quad \text{where } \pi(x) = E(Y|X = x) = P(Y = 1|X = x)$$

Thus the log of their ratio

$$\begin{aligned} \log \left(\frac{\text{odds}(Y = 1|X = x + 1)}{\text{odds}(Y = 1|X = x)} \right) &= \log \left(\frac{P(Y = 1|X = x + 1)/P(Y = 0|X = x + 1)}{P(Y = 1|X = x)/P(Y = 0|X = x)} \right) \\ &= \log \left(\frac{\pi(x + 1)/(1 - \pi(x + 1))}{\pi(x)/(1 - \pi(x))} \right) \\ &= \log(\pi(x + 1)/(1 - \pi(x + 1))) - \log(\pi(x)/(1 - \pi(x))) \\ &= \text{logit}(\pi(x + 1)) - \text{logit}(\pi(x)) \\ &= (x + 1)\beta - x\beta \\ &= \beta \end{aligned}$$

Thus it does not depend on X . This shows that a unit increase in X leads to an increase of beta in the log-odds ratio i.e a unit increase in x leads to the ratio of their odds increasing by e^β

Problem 4 [10 points]: Complete the following parts:

- Explain important findings and model information from the summary table produced by a call to `summary(m1)` in the CCSO example in the logistic regression notes. Keep in mind that we restricted attention to "other traffic offenses" in the CCSO example, and that this data is observational.
- Explain important findings and model information from the summary table produced by a call to `summary(m1)` in the Galapagos islands example in the count regression notes.

Solution 4:

```
# Creating the logistic regression model for the CCSO model
m1 <- glm(atleastone ~ -1 + Race + Sex + Age, data = CCSO_small,
family = "binomial", x = "TRUE")
summary(m1)

##
## Call:
## glm(formula = atleastone ~ -1 + Race + Sex + Age, family = "binomial",
##      data = CCSO_small, x = "TRUE")
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -0.9393 -0.5485 -0.4817 -0.3391 2.6527
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## RaceAsian/Pacific Islander -4.389865 0.523612 -8.384 < 2e-16 ***
## RaceBlack -1.876550 0.144601 -12.977 < 2e-16 ***
## RaceHispanic -2.804549 0.173349 -16.179 < 2e-16 ***
## RaceWhite -3.043226 0.147160 -20.680 < 2e-16 ***
## SexMale 0.739834 0.105380 7.021 2.21e-12 ***
## Age 0.007705 0.003186 2.418 0.0156 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 8201.3 on 5916 degrees of freedom
## Residual deviance: 4668.7 on 5910 degrees of freedom
## AIC: 4680.7
##
## Number of Fisher Scoring iterations: 6
```

The estimate column gives the estimate for β for the logistic model

$$\text{logit}(E(Y|X)) = X\beta$$

A unit increase in the predictor variable X_j corresponds to an increase of β_j (estimated by $\hat{\beta}_j$) in the log-odds ratio with everything else being held fixed. A simpler interpretation is that $\hat{\beta}_j > 0$ can be interpreted as: An increase in X_j implies that $P(Y = 1|X = x)$ increases.

Thus

- Race is significant when testing at reasonable significance levels. We observe that Black individuals are estimated to have comparatively larger propensity of incarcerations lasting one day or longer for “other traffic offenses”. We would need to look into other factors such as socio-economic status, repeat offenders, and multiple offenses before we could conclude that race is the driver of longer incarcerations.
- Sex being Male is estimated to increase the propensity of incarcerations lasting longer than one day for “other traffic offenses”.
- Age increasing also is estimated to increase the propensity of incarcerations lasting longer than one day for “other traffic offenses”.

The standard error column gives the standard error of the estimate of the β coefficients. The Z-value and P-value help in detecting the significance of the covariates. At a level of $\alpha = 0.05$ we can see that all the covariates are significant.

The null deviance and residual deviance give information about the goodness of fit of the null model (with no covariates) and the submodel we consider respectively. To check if the sub-model is better than the saturated model we can do a χ^2 test because under H_0 (The submodel is a better fit)

$$D(y; \hat{\mu}) \sim \chi^2_{n-p},$$

where D is the deviance.

```
pchisq(m1$deviance, df = m1$df.residual, lower = FALSE)
```

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

Since the p-value is 1 this shows that the submodel is indeed a good fit to the data. We can also check if the null model (M_0) is better than the submodel (M_1) we choose.

$$H_0 : M_0 \text{ true} \quad H_a : M_1 \text{ true, but not } M_0$$

Then

$$D(y; \hat{\mu}_0) - D(y; \hat{\mu}_1) \sim \chi^2_{p_0 - p_1}$$

where $D(y; \hat{\mu}_0)$ is the null deviance and $D(y; \hat{\mu}_1)$ is the residual deviance.

```
pchisq(m1$null.deviance - m1$deviance, df = m1$df.null - m1$df.residual,
lower = FALSE)
```

```
## [1] 0
```

Since the pvalue is 0 this means that the submodel we choose is a better fit than the null model.

- (b) Summarize the summary tables produced by a call to `summary(m1)` in the Galapagos islands example in the count regression notes.

```
m2 <- glm(Species ~ Elevation + Nearest + Scruz + Adjacent + Size,
family = "poisson", data = gala, x = TRUE)
summary(m2)
```

```
##
## Call:
## glm(formula = Species ~ Elevation + Nearest + Scruz + Adjacent +
##       Size, family = "poisson", data = gala, x = TRUE)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -10.3723  -3.5214  -0.9947   1.7193  10.6627
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  2.790e+00  8.108e-02  34.410 < 2e-16 ***
## Elevation    9.361e-04  5.402e-05  17.329 < 2e-16 ***
## Nearest      6.469e-03  1.748e-03   3.702 0.000214 ***
## Scruz       -6.266e-03  6.268e-04  -9.997 < 2e-16 ***
## Adjacent    -2.858e-04  2.961e-05  -9.652 < 2e-16 ***
## Size2        1.128e+00  9.535e-02  11.826 < 2e-16 ***
## Size3        2.059e+00  9.419e-02  21.856 < 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: 3510.73  on 29  degrees of freedom
## Residual deviance:  594.18  on 23  degrees of freedom
## AIC: 769.01
##
## Number of Fisher Scoring iterations: 5
```

The estimate column gives the estimate for β for the logistic model

$$\log(E(Y|X)) = X\beta$$

. A unit increase in the predictor variable X_j corresponds to an increase of β_j (estimated by $\hat{\beta}_j$) in the log of the mean response with everything else being held fixed. A simpler interpretation is that $\hat{\beta}_j > 0$ can be interpreted as: An increase in X_j implies that the mean response increases.

Thus

- Elevation is estimated to increase the expected number of plant species found on each island
- Distance to the nearest island is estimated to increase the expected number of plant species found on each island
- Distance to Scrub is estimated to decrease the number of expected plant species found on each island
- Larger adjacent islands are also estimated to decrease the number of expected plant species found on each island
- Medium and large islands are estimated to have more expected plant species found on each island

The standard error column gives the standard error of the estimate of the β coefficients. The Z-value and P-value help in detecting the significance of the covariates. At a level of $\alpha = 0.05$ we can see that all the covariates are significant.

The null deviance and residual deviance give information about the goodness of fit of the null model (with no covariates) and the submodel we consider respectively. To check if the sub-model is better than the saturated model we can do a χ^2 test because under H_0 (The submodel is a better fit)

$$D(y; \hat{\mu}) \sim \chi^2_{n-p},$$

where D is the deviance.

```
pchisq(m2$deviance, df = m2$df.residual, lower = FALSE)
```

```
## [1] 7.617409e-111
```

Since the p-value is very small this shows that the submodel not really a good fit to the data. We prefer the saturated model over it. We can also check if the null model (M_0) is better than the submodel (M_1) we choose.

$$H_0 : M_0 \text{ true} \quad H_a : M_1 \text{ true, but not } M_0$$

Then

$$D(y; \hat{\mu}_0) - D(y; \hat{\mu}_1) \sim \chi^2_{p_0 - p_1}$$

where $D(y; \hat{\mu}_0)$ is the null deviance and $D(y; \hat{\mu}_1)$ is the residual deviance.

```
pchisq(m2$null.deviance - m2$deviance, df = m2$df.null - m2$df.residual,
lower = FALSE)
```

```
## [1] 0
```

Since the pvalue is 0 this means that the submodel we choose is a better fit than the null model.

Problem 5 [10 points]: Derive expressions and compute standard errors $se(\hat{\mu})$ in the logistic regression CCSO example without using `predict.glm`. Then construct Wald based confidence intervals for the estimated mean

value parameters. Also construct confidence intervals $(g(\hat{\beta} - z_{\alpha/2}se(\hat{\beta})), g(\hat{\beta} + z_{\alpha/2}se(\hat{\beta})))$. Comment on the differences between these two confidence intervals for $\hat{\mu}$.

###solution 5

We know that

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, \Sigma^{-1})$$

where $\Sigma = X^T W X$, $W = \text{diag}(\mu_i(1 - \mu_i))$ and $\mu = e^{X\beta} / (1 + e^{X\beta})$

Now $\mu_i = g(\beta) = e^{M_i^T \beta} / (1 + e^{M_i^T \beta}) \implies g'(\beta) = M_i * \mu_i(1 - \mu_i)$

Thus by the delta method

$$\sqrt{n}(\hat{\mu}_i - \mu_i) \xrightarrow{d} N(0, \hat{\mu}_i^2(1 - \hat{\mu}_i)^2 M_i^T \Sigma^{-1} M_i)$$

Thus our estimation of the SE is:

```
#Creating the model matrix
X = model.matrix(atleastone ~ -1 + Race + Sex + Age, data = CCSO_small)
#Calculating the se
se_pihat = sqrt(apply(X, 1, function(j) t(j) %*% var_matrix_CCSO %*% j)) * pi_CCSO * (1 - pi_CCSO)
se_pihat[1:20]

## [1] 0.006773575 0.006392931 0.005695285 0.006773575 0.013586275 0.013586275
## [7] 0.011956810 0.013052798 0.005695285 0.006269441 0.010557231 0.010557231
## [13] 0.009911238 0.009911238 0.012797440 0.014138891 0.006552717 0.013037138
## [19] 0.006472744 0.006370288
```

The Wald based confidence intervals for the estimated mean value parameters are

$$\hat{\mu}_i(x) \pm z_{1-\alpha/2} \sigma_i \quad \text{where } \sigma_i = \hat{\mu}_i(1 - \hat{\mu}_i) \sqrt{M_i^T \Sigma^{-1} M_i}$$

```
#Creating the conf intervals
conf_lower = pi_CCSO - qnorm(0.975) * se_pihat
conf_upper = pi_CCSO + qnorm(0.975) * se_pihat
waldci = cbind(conf_upper, conf_lower)
waldci = waldci %>% as.data.frame() %>%
  mutate(length = conf_upper - conf_lower)
head(waldci)
```

```
##           V1           V2      length
## 1 0.11912305 0.09257112 0.02655193
## 2 0.12132956 0.09626973 0.02505983
## 3 0.06866721 0.04634211 0.02232511
## 4 0.11912305 0.09257112 0.02655193
## 5 0.16908475 0.11582753 0.05325722
## 6 0.16908475 0.11582753 0.05325722
```

And, for the other confidence interval type $g(\hat{\beta} + z_{\alpha/2}se(\hat{\beta}))$:

```
m1 <- glm(atleastone ~ -1 + Race + Sex + Age, data = CCSO_small,
          family = "binomial", x = "TRUE")

betahat = m1$coefficients
M = m1$x
alpha = 0.025
```

```

z = qnorm(p = 1-alpha)
n = nrow(M)
pici_upper = 1/(1 + exp( - M %*% (betahat + z*diag(vcov(m1))) ))
pici_lower = 1/(1 + exp( - M %*% (betahat - z*diag(vcov(m1))) ))
pici = cbind(pici_upper, pici_lower)
pici = pici %>% as.data.frame() %>%
  mutate(length = pici_upper - pici_lower)
head(pici)

```

```

##           V1           V2      length
## 1 0.11212335 0.09988262 0.012240731
## 2 0.11523654 0.10268058 0.012555961
## 3 0.05988461 0.05521374 0.004670872
## 4 0.11212335 0.09988262 0.012240731
## 5 0.15268704 0.13280332 0.019883714
## 6 0.15268704 0.13280332 0.019883714

```

Then, the average length of the Wald and plug-in approaches are given as:

```

avg_length_wald = round(mean(waldci$length), digits=4)
avg_length_wald

```

```
## [1] 0.0392
```

```

avg_length_pi = round(mean(pici$length), digits=4)
avg_length_pi

```

```
## [1] 0.0157
```

Then, the average Wald CI is somewhat larger than the average plug-in CI.

Problem 6 [10 points]: Construct a nonparametric bootstrap procedure that estimates the uncertainty associated with both estimates of the average treatment effect (ATE) of online learning in the logistic regression notes. Do the conclusions change when we factor in the uncertainty obtained from the nonparametric bootstrap procedure? Explain.

Solution 6:

Here we are estimating the ATE for the scores earned by students in online learning as opposed to in-person learning.

```

#Reading in the data

dat = read.csv("../notes/3_binary_response/online.csv")
dat_small <- dat %>% dplyr::select(Online, ACTMath, ACTMajor, ACT, Gender,
International, F17, S18, S19, Fa19, FR, SO, JR)

ATE_alt = NULL

#Taking 1000 bootstrap samples
for(i in 1:1000)
{
  rand = sample(nrow(dat_small),replace = T)
  m <- glm(Online ~., data = dat_small[rand,], family = "binomial")
  trt <- dat_small[rand,]$Online

```

```

preds <- predict(m, type = "response")
weights_alt_trt <- 1 / sum(trt / preds) * trt / preds
weights_alt_notrt <- 1 / sum((1 - trt) / (1 - preds)) * (1 - trt) / (1 - preds)
dat_new <- data.frame(dat[rand,], weights = weights_alt_trt - weights_alt_notrt)
ATE_alt <- c(ATE_alt, sum(weights_alt_trt * dat_new$ObjExam) -
sum(weights_alt_notrt * dat_new$ObjExam))
}

```

```
mean(ATE_alt)
```

```
## [1] 0.5886974
```

```
var(ATE_alt)
```

```
## [1] 0.45145
```

```
quantile(ATE_alt, prob = c(0.025, 0.975))
```

```
##      2.5%      97.5%
```

```
## -0.6685877  1.8941839
```

Since the mean and variance are small and the confidence interval contains 0 we can still conclude that there is no difference between the two types of learning.

```

ATE_DR = NULL
for(i in 1:1000)
{
  rand = sample(nrow(dat), replace = T)
  trt = dat_small[rand,]$Online
  m <- glm(Online ~ ., data = dat_small[rand,], family = "binomial")
  preds <- predict(m, type = "response")
  dat_boot = dat[rand,]
  m_trt <- lm(ObjExam ~ ACTMath + ACTMajor + ACT + International + Gender +
FR + SO + JR + F17 + S18 + S19,
data = dat_boot[trt == 1, ])
  Y_trt <- predict(m_trt, newdata = dat_boot)
  m_notrt <- lm(ObjExam ~ ACTMath + ACTMajor + ACT + International + Gender +
FR + SO + JR + F17 + S18 + S19,
data = dat_boot[trt == 0, ])
  Y_notrt <- predict(m_notrt, newdata = dat_boot)
  ATE_DR <- c(ATE_DR, mean( (dat_boot$ObjExam * trt - (trt - preds) * Y_trt) / preds -
(dat_boot$ObjExam * (1 - trt) + (trt - preds)*Y_notrt) / (1 - preds)))
}

```

```
mean(ATE_DR)
```

```
## [1] 0.4155045
```

```
var(ATE_DR)
```

```
## [1] 0.3864671
```

```
quantile(ATE_DR, prob = c(0.025, 0.975))
```

```
##      2.5%      97.5%
```

```
## -0.8000527  1.6897111
```

Since the mean and variance are small even for the robust estimate and the confidence interval contains 0 we

can still conclude that there is no difference between the two types of learning.

Solution 7

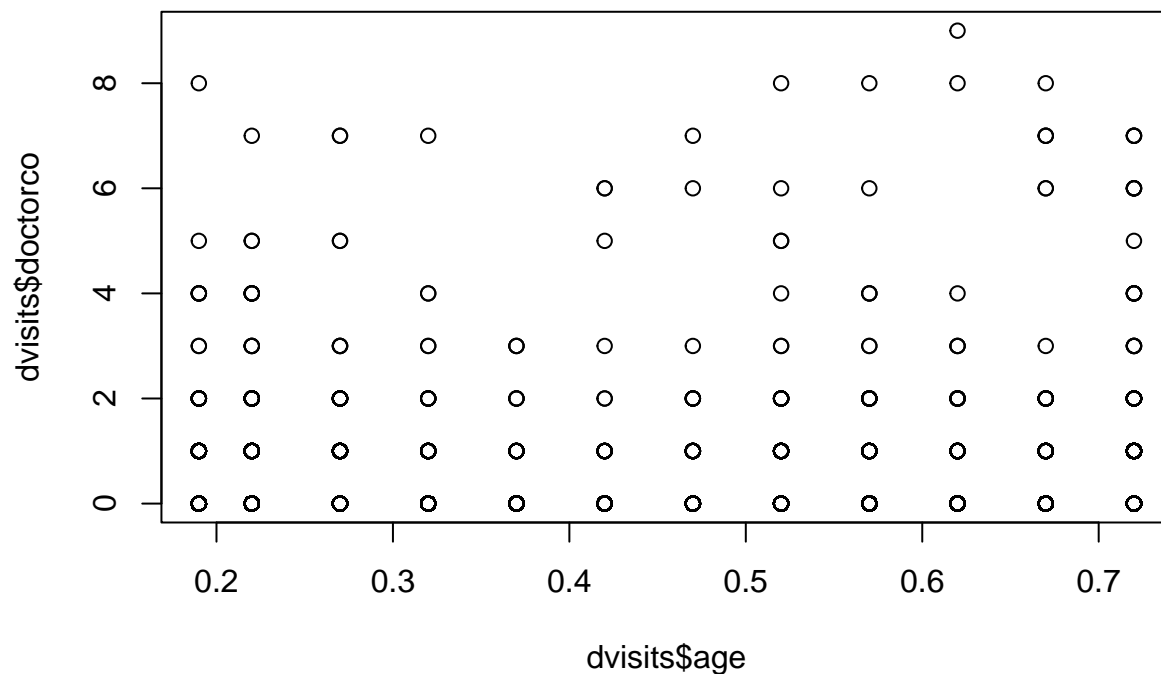
Problem 7 [15 points]: Use the `dvisits` data in the `faraway` package to answer the follow parts:

- (a) Make plots which show the relationship between the response variable, `doctorco`, and the potential predictors, `age` and `illness`.
- (b) Combine the predictors `chcond1` and `chcond2` into a single three-level factor. Make an appropriate plot showing the relationship between this factor and the response. Comment.
- (c) Build a Poisson regression model with `doctorco` as the response and `sex`, `age`, `agesq`, `income`, `levyplus`, `freepoor`, `freerepa`, `illness`, `actdays`, `hscore` and the three-level condition factor as possible predictor variables. Considering the deviance of this model, does this model fit the data?
- (d) Plot the residuals and the fitted values — why are there lines of observations on the plot? Make a QQ plot of the residuals and comment.
- (e) Use a stepwise AIC-based model selection method. What sort of person would be predicted to visit the doctor the most under your selected model?
- (f) For the last person in the dataset, compute the predicted probability distribution for their visits to the doctor, i.e., give the probability they visit 0, 1, 2, etc. times.
- (g) Tabulate the frequencies of the number of doctor visits. Compute the expected frequencies of doctor visits under your most recent model. Compare the observed with the expected frequencies and comment on whether it is worth fitting a zero-inflated count model.
- (h) Fit a comparable (Gaussian) linear model and graphically compare the fits. Describe how they differ.

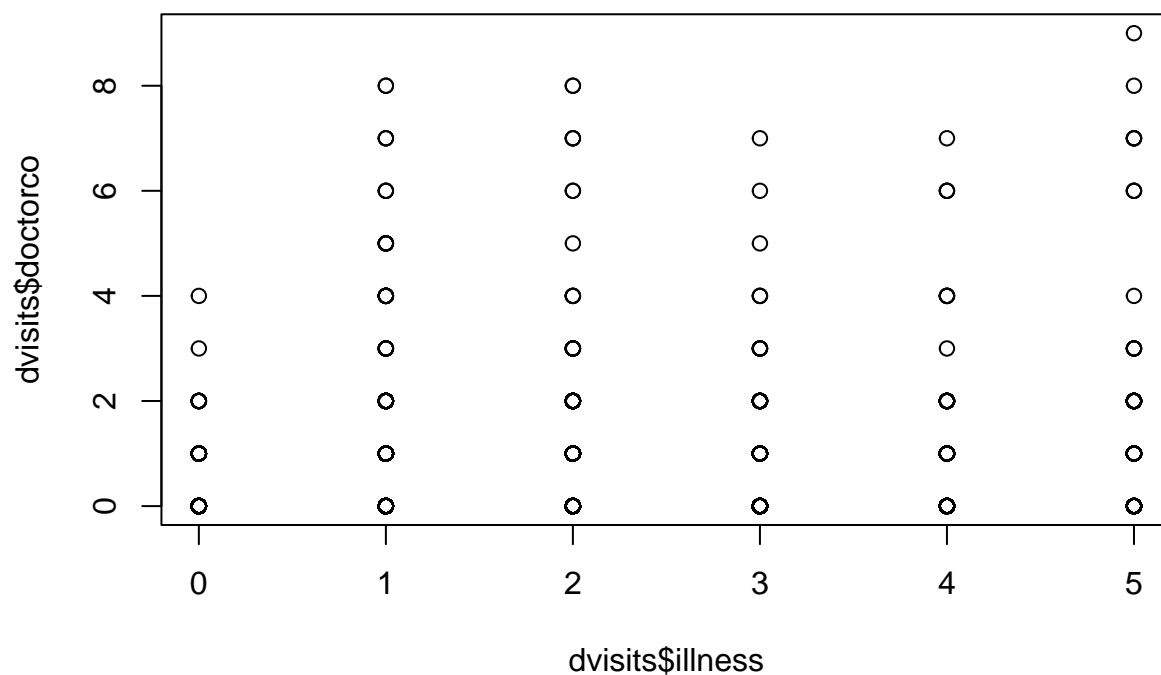
Solution 7:

- (a) Make plots which show the relationship between the response variable, `doctorco`, and the potential predictors, `age` and `illness`.

```
data(dvisits)
plot(dvisits$age, dvisits$doctorco)
```



```
plot(dvisits$illness,dvisits$doctorco)
```

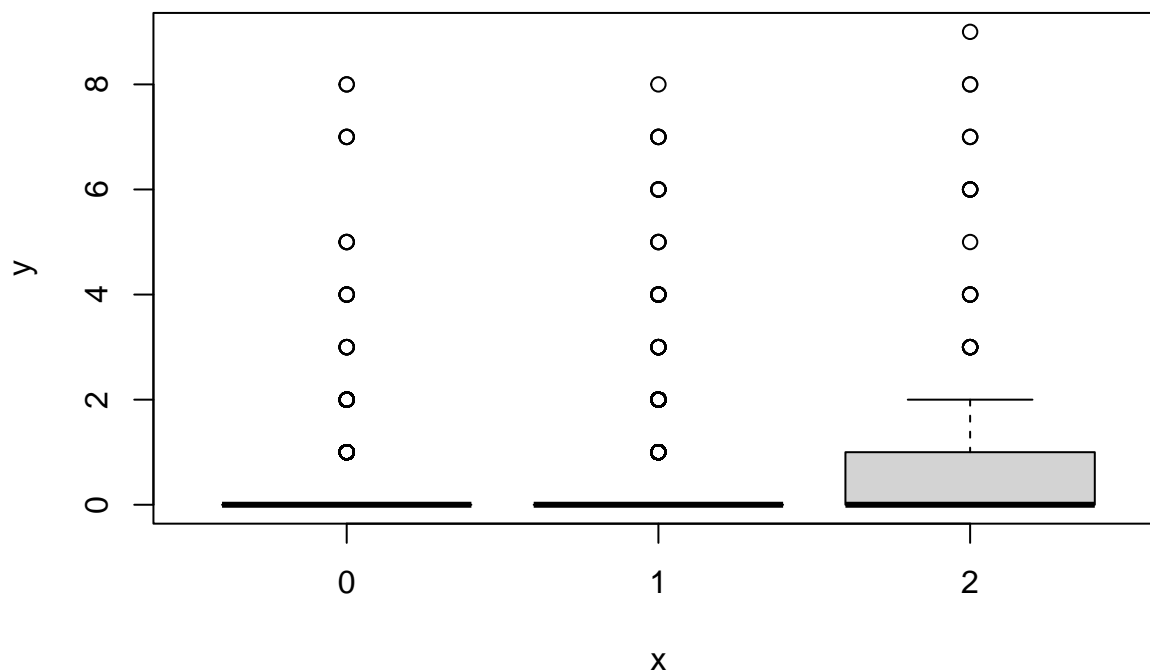


- (b) Combine the predictors chcond1 and chcond2 into a single three-level factor. Make an appropriate plot showing the relationship between this factor and the response. Comment.

We create a new variable `chcond` which takes 3 factor values

- 1 if the patient has a chronic condition(s) but is not limited in activity
- 2 if the patient has a chronic condition(s) but is limited in activity
- 0 Otherwise

```
chcond = as.factor(dvisits$chcond1+2*dvisits$chcond2)
plot(chcond,dvisits$doctorco)
```



From the plot we can see that patients with chronic conditions which limit activity have more visits to the doctor. This is possibly a factor which influences the response variable.

- (c) Build a Poisson regression model with doctorco as the response and sex, age, agesq, income, levyplus, freepoor, freerepa, illness, actdays, hscore and the three-level condition factor as possible predictor variables. Considering the deviance of this model, does this model fit the data?

```
dat = dvisits %>% dplyr::select(doctorco, sex, age, agesq, income,
levyplus, freepoor, freerepa, illness, actdays, hscore)
dat = cbind(dat, "chcond" = chcond)
head(dat)
```

```
## doctorco sex age agesq income levyplus freepoor freerepa illness actdays
## 1 1 1 0.19 0.0361 0.55 1 0 0 1 4
## 2 1 1 0.19 0.0361 0.45 1 0 0 1 2
## 3 1 0 0.19 0.0361 0.90 0 0 0 3 0
## 4 1 0 0.19 0.0361 0.15 0 0 0 1 0
## 5 1 0 0.19 0.0361 0.45 0 0 0 2 5
## 6 1 1 0.19 0.0361 0.35 0 0 0 5 1
## hscore chcond
## 1 1 0
## 2 1 0
## 3 0 0
## 4 0 0
## 5 1 1
## 6 9 1
```

```
#Creating the model
mod = glm(doctorco~., data = dat, family = "poisson", x = TRUE)
summary(mod)
```

```
##
## Call:
```

```
## glm(formula = doctorco ~ ., family = "poisson", data = dat, x = TRUE)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.9170  -0.6862  -0.5743  -0.4839   5.7005
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -2.223848   0.189816 -11.716  <2e-16 ***
## sex          0.156882   0.056137   2.795   0.0052 **
## age          1.056299   1.000780   1.055   0.2912
## agesq       -0.848704   1.077784  -0.787   0.4310
## income      -0.205321   0.088379  -2.323   0.0202 *
## levyplus     0.123185   0.071640   1.720   0.0855 .
## freepoor    -0.440061   0.179811  -2.447   0.0144 *
## freerepa     0.079798   0.092060   0.867   0.3860
## illness      0.186948   0.018281  10.227  <2e-16 ***
## actdays     0.126846   0.005034  25.198  <2e-16 ***
## hscore       0.030081   0.010099   2.979   0.0029 **
## chcond1      0.114085   0.066640   1.712   0.0869 .
## chcond2      0.141158   0.083145   1.698   0.0896 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 5634.8  on 5189  degrees of freedom
## Residual deviance: 4379.5  on 5177  degrees of freedom
## AIC: 6737.1
##
## Number of Fisher Scoring iterations: 6
```

```
#Testing for the goodness of fit of the model against the saturated model
pchisq(mod$deviance, df = mod$df.residual, lower = FALSE)
```

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

Since the pvalue is 1 we can conclude that the model is indeed a better fit than the saturated model.

```
pchisq(mod$null.deviance - mod$deviance, df = mod$df.null - mod$df.residual,
lower = FALSE) %>% round(4)
```

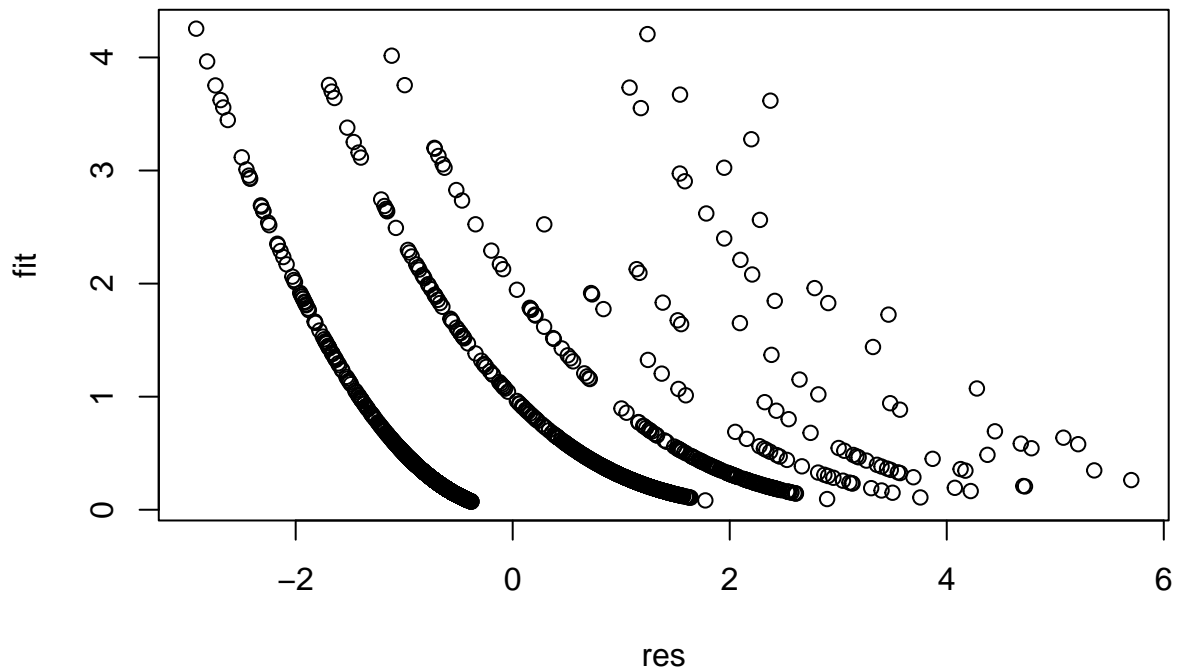
```
## [1] 0
```

Since the pvalue is nearly 0 we can conclude that the model is also better than the null model. Thus it is a appropriate fit to the data.

- (d) Plot the residuals and the fitted values — why are there lines of observations on the plot? Make a QQ plot of the residuals and comment.

```
res = residuals(mod)
fit = fitted(mod)

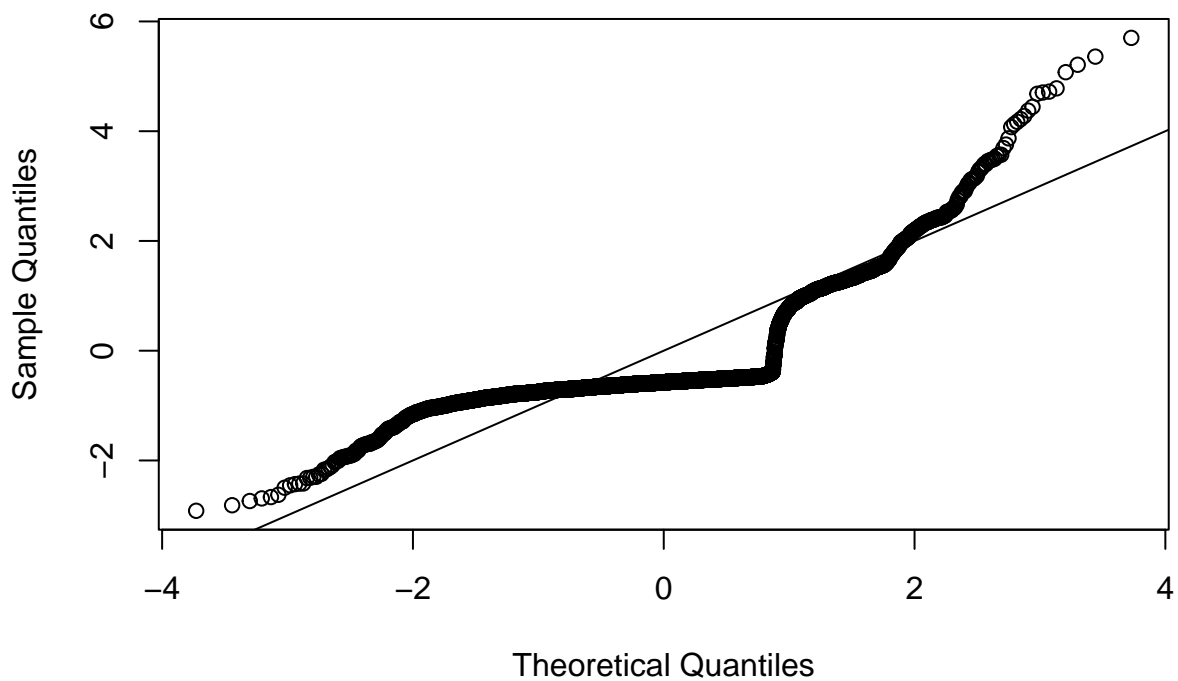
plot(res,fit)
```

We observe lines of observations because most of the variables are factor variables with a small number of levels.

```
qqnorm(res)
abline(0, 1)
```

Normal Q-Q Plot



The QQ-plot shows that the residuals do not follow a normal distribution very well indicating the normality of residuals assumption is not reasonable.

- (e) Use a stepwise AIC-based model selection method. What sort of person would be predicted to visit the

doctor the most under your selected model?

#Selecting the model

```
library(MASS)
```

```
mod_Select = stepAIC(mod)
```

```
## Start: AIC=6737.08
```

```
## doctorco ~ sex + age + agesq + income + levyplus + freepoor +  
##     freerepa + illness + actdays + hscore + chcond
```

```
##
```

	Df	Deviance	AIC
## - agesq	1	4380.1	6735.7
## - freerepa	1	4380.3	6735.8
## - age	1	4380.6	6736.2
## - chcond	2	4383.2	6736.7
## <none>		4379.5	6737.1
## - levyplus	1	4382.5	6738.1
## - income	1	4385.0	6740.5
## - freepoor	1	4386.2	6741.8
## - sex	1	4387.4	6743.0
## - hscore	1	4388.1	6743.7
## - illness	1	4481.8	6837.4
## - actdays	1	4917.1	7272.7

```
##
```

```
## Step: AIC=6735.7
```

```
## doctorco ~ sex + age + income + levyplus + freepoor + freerepa +  
##     illness + actdays + hscore + chcond
```

```
##
```

	Df	Deviance	AIC
## - freerepa	1	4381.0	6734.5
## <none>		4380.1	6735.7
## - chcond	2	4384.2	6735.8
## - age	1	4383.0	6736.5
## - levyplus	1	4383.3	6736.9
## - income	1	4385.0	6738.6
## - freepoor	1	4386.8	6740.4
## - sex	1	4388.0	6741.5
## - hscore	1	4389.1	6742.7
## - illness	1	4481.9	6835.4
## - actdays	1	4917.1	7270.7

```
##
```

```
## Step: AIC=6734.53
```

```
## doctorco ~ sex + age + income + levyplus + freepoor + illness +  
##     actdays + hscore + chcond
```

```
##
```

	Df	Deviance	AIC
## <none>		4381.0	6734.5
## - levyplus	1	4383.4	6735.0
## - chcond	2	4385.5	6735.0
## - income	1	4386.7	6738.2
## - age	1	4387.1	6738.7
## - freepoor	1	4389.1	6740.6
## - sex	1	4389.5	6741.0
## - hscore	1	4390.2	6741.8
## - illness	1	4482.7	6834.2

```
## - actdays 1 4917.6 7269.2
```

```
#Outputting the best models
```

```
mod_Select$anova
```

```
## Stepwise Model Path
```

```
## Analysis of Deviance Table
```

```
##
```

```
## Initial Model:
```

```
## doctorco ~ sex + age + agesq + income + levyplus + freepoor +
```

```
## freerepa + illness + actdays + hscore + chcond
```

```
##
```

```
## Final Model:
```

```
## doctorco ~ sex + age + income + levyplus + freepoor + illness +
```

```
## actdays + hscore + chcond
```

```
##
```

```
##
```

```
## Step Df Deviance Resid. Df Resid. Dev AIC
```

```
## 1 5177 4379.515 6737.083
```

```
## 2 - agesq 1 0.6180113 5178 4380.133 6735.701
```

```
## 3 - freerepa 1 0.8279216 5179 4380.961 6734.529
```

```
# The beta coefficients of our best model
```

```
beta = mod_Select$coefficients
```

```
beta
```

```
## (Intercept) sex age income levyplus freepoor
```

```
## -2.08906349 0.16199995 0.35513074 -0.19980641 0.08368852 -0.46959634
```

```
## illness actdays hscore chcond1 chcond2
```

```
## 0.18610078 0.12661065 0.03111559 0.12110045 0.15889355
```

We can see that the number of doctor consultations increases when the patient is female, with increasing age, with low income, if covered by private health insurance, not covered by the government insurance for low income, high number of illness, high number of days of reduced activity, bad health score and with presence of chronic conditions. This indicates a poorer older woman with private insurance and higher number of illness is predicted to visit doctor more often.

- (f) For the last person in the dataset, compute the predicted probability distribution for their visits to the doctor, i.e., give the probability they visit 0, 1, 2, etc. times.

```
options(scipen = 99)
```

```
X = mod_Select$x
```

```
Y = mod_Select$y
```

```
hat_lambda_last = exp(t(X[nrow(dat),])%*%beta)
```

```
data.frame("Value" = 0:9, "Prob" = dpois(0:9,hat_lambda_last) %>% round(6) )
```

```
## Value Prob
```

```
## 1 0 0.858916
```

```
## 2 1 0.130628
```

```
## 3 2 0.009933
```

```
## 4 3 0.000504
```

```
## 5 4 0.000019
```

```
## 6 5 0.000001
```

```
## 7 6 0.000000
```

```
## 8 7 0.000000
```

```
## 9 8 0.000000
```

```
## 10 9 0.000000
```

- (g) Tabulate the frequencies of the number of doctor visits. Compute the expected frequencies of doctor visits under your most recent model. Compare the observed with the expected frequencies and comment on whether it is worth fitting a zero-inflated count model.

```
observed_freq <- with(dvisits, table(doctorco))
est <- matrix(nrow=dim(dvisits)[1], ncol=10)
for(i in 1:dim(dvisits)[1]){
  est[i,] <- dpois(0:9, fitted.values(mod_Select)[i])
}
expected_freq <- colMeans(est)*dim(dvisits)[1]
cbind.data.frame(observed_freq, expected_freq)
```

```
##      doctorco Freq expected_freq
## 1          0 4141  4013.6020569
## 2          1  782   928.3492327
## 3          2  174   168.0095991
## 4          3   30    45.4859546
## 5          4   24    18.9118479
## 6          5    9     8.8170066
## 7          6   12     4.0118473
## 8          7   12     1.7230627
## 9          8    5     0.6931622
## 10         9    1     0.2608052
```

From the two tables, the observed and expected frequencies are close enough and thus it does not seem worth fitting a zero inflated model.

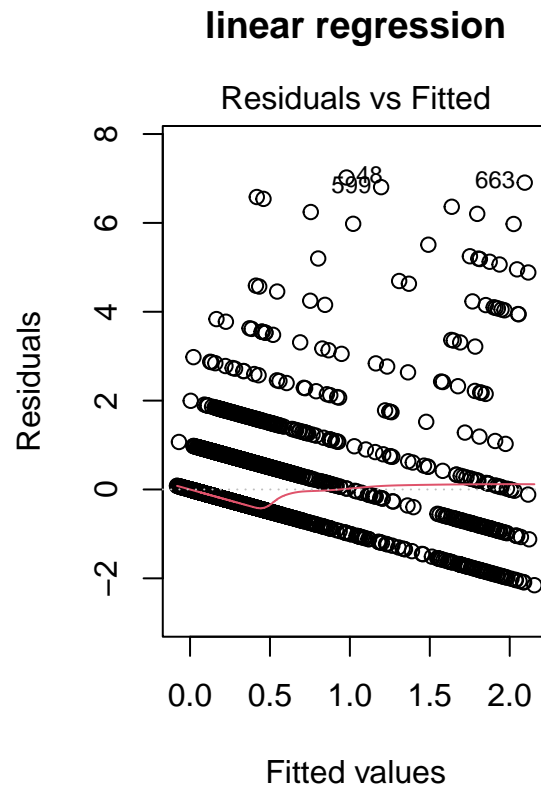
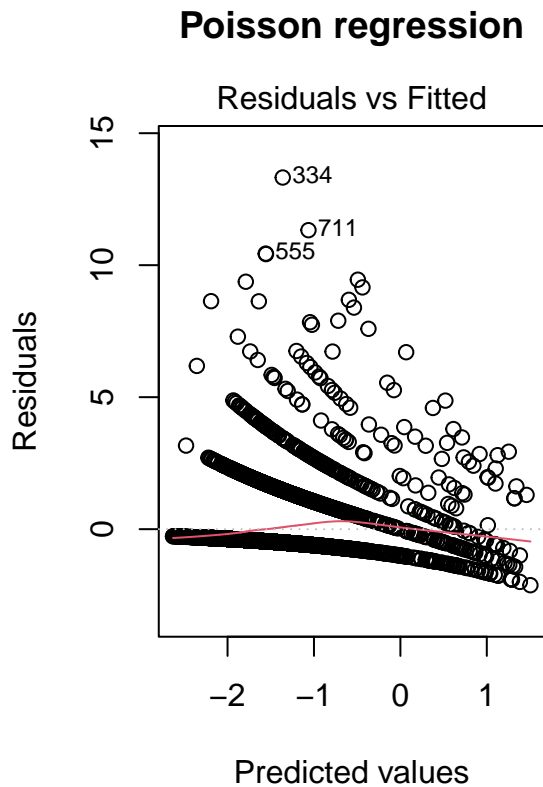
- (h) Fit a comparable (Gaussian) linear model and graphically compare the fits. Describe how they differ.

```
dvisitsmod_lm <- lm(Y ~ X)
summary(dvisitsmod_lm)
```

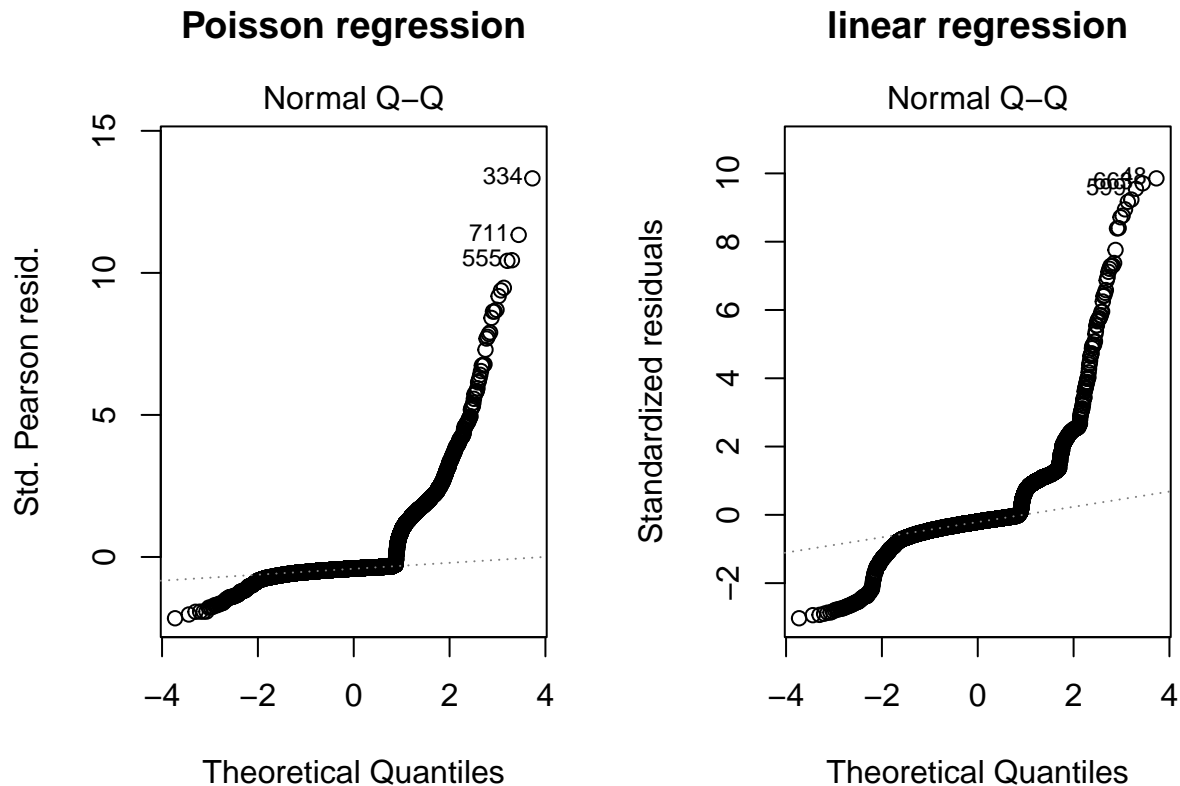
```
##
## Call:
## lm(formula = Y ~ X)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.1543 -0.2584 -0.1440 -0.0434  7.0211
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept)  0.036781   0.035794   1.028    0.304201
## X(Intercept)         NA          NA      NA         NA
## Xsex           0.035574   0.021505   1.654    0.098137 .
## Xage           0.180024   0.055912   3.220    0.001291 **
## Xincome       -0.061208   0.030522  -2.005    0.044975 *
## Xlevyplus      0.024041   0.021235   1.132    0.257626
## Xfreepoor     -0.111650   0.051532  -2.167    0.030308 *
## Xillness       0.060148   0.008332   7.219  0.000000000000602 ***
## Xactdays      0.103140   0.003656  28.213 < 0.000000000000002 ***
## Xhscore       0.017064   0.005180   3.294    0.000994 ***
## Xchcond1      0.005006   0.023709   0.211    0.832776
## Xchcond2      0.045319   0.035352   1.282    0.199921
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Residual standard error: 0.7138 on 5179 degrees of freedom
## Multiple R-squared:  0.2017, Adjusted R-squared:  0.2002
## F-statistic: 130.9 on 10 and 5179 DF,  p-value: < 0.00000000000000022
```

```
par(mfrow=c(1,2))
plot(mod_Select, which=1)
title(main = "Poisson regression\n")
plot(dvisitsmod_lm, which=1)
title(main = "linear regression\n")
```



```
par(mfrow=c(1,2))
plot(mod_Select, which=2)
title(main = "Poisson regression\n")
plot(dvisitsmod_lm, which=2)
title(main = "linear regression\n")
```



This seems to indicate how the Poisson regression is a little better because the line for the residuals plot is much more linear around zero, along with the Q-Q plot. In conclusion, we can see that even though we need to account for overdispersion in the Poisson model it is still a better fit than the linear model.

Problem 8 [20 points]: Analyze the CCSO data set with Days in Jail as the response variable. You are allowed to dichotomize the response into a binary variable. Restrict attention to other traffic offenses as done in class. Your analysis needs to consider the variables considered in class as well as repeat offenders, multiple offenses, released reason, and agency. The determination of repeat offenders and multiple offenses can be done via the jacket number variable. Report interesting significant and null findings and determine that your final model is appropriate. You are allowed to use other inferential techniques than GLM. If you do so, then you need to justify your choices.

Solution 8

Data Wrangling: We first select the necessary data, here we focus on OTHER TRAFFIC OFFENSES type crimes. Our response variable is the binary variable denoting whether the offender stayed in jail for at least one day or not. We also restrict ourselves with some specific Race, Sex, and Arrest Agencies.

```
CCSD <- fread("https://uofi.box.com/shared/static/9elozjsg99bgcb7gb546wlfr3r2gc9b7.csv")
```

The full data has 67764 observations and 35 features.

```
## data wrangling
CCSO_imp <- CCSO %>% rename(Days = "Days in Jail", Age = "Age at Arrest",
                           Date = "BOOKING DATE", Sex = "SEX", Race = "RACE",
                           Crime = "CRIME CODE", Agency = "ARREST AGENCY", JNum = "JACKET NUMBER") %>%
```

```

mutate(atleastone = ifelse(Days > 0,1,0)) %>%
filter(Crime == "OTHER TRAFFIC OFFENSES") %>%
filter(Race %in% c("Asian/Pacific Islander","Black","White","Hispanic")) %>%
filter(Sex %in% c("Female","Male")) %>%
## Agency (and possible interactions with Agency) is interesting if you want to go deeper filter(Agency
filter(Agency %in% c( "Illinois State Police", "Rantoul Police Department",
"University of Illinois Police Department", "Urbana Police Department",
"Champaign County Sherriff's Office","Champaign Police Department") ) %>%
dplyr::select(atleastone, Age, Sex, Date, Race, Agency, JNum, `RELEASED REASON`, `EMPLOYMENT STATUS`,
mutate(Race = fct_drop(Race), Sex = fct_drop(Sex), Agency = fct_drop(Agency))

CCSO_imp <- CCSO_imp[complete.cases(CCSO_imp), ]

```

Now, we will combine categories for the features RELEASED REASON, EMPLOYMENT STATUS, INCARCERATION REASON in order to reduce number of factors in these features.

```

CCSO_imp$`RELEASED REASON`[CCSO_imp$`RELEASED REASON` %in% c("Cash Bond Posted","Paid Fine, Court Costs
CCSO_imp$`RELEASED REASON`[CCSO_imp$`RELEASED REASON` %in% c("Credit Card Bond Posted")] = "CardBonds"
CCSO_imp$`RELEASED REASON`[CCSO_imp$`RELEASED REASON` %in% c("Placed on Probation

CCSO_imp$`RELEASED REASON`[CCSO_imp$`RELEASED REASON` %in% c("Release on Personal Recognizance
CCSO_imp$`RELEASED REASON`[CCSO_imp$`RELEASED REASON` %in% c("Served Sentence of Incarceration
CCSO_imp$`RELEASED REASON`[CCSO_imp$`RELEASED REASON` %in% c("Transfer to other county/state authorities
CCSO_imp$`RELEASED REASON`[!(CCSO_imp$`RELEASED REASON` %in% c("Bonds paid","Conditional Release","Rele

CCSO_imp$`EMPLOYMENT STATUS`[CCSO_imp$`EMPLOYMENT STATUS` %in% c("Employed - Full Time","Employed - Par
CCSO_imp$`EMPLOYMENT STATUS`[CCSO_imp$`EMPLOYMENT STATUS` %in% c("Student","Retired")] = "Student/Reti
CCSO_imp$`EMPLOYMENT STATUS`[CCSO_imp$`EMPLOYMENT STATUS` %in% c("Unemployed","Laid Off")] ="Unemployed

CCSO_imp$`INCARCERATION REASON`[CCSO_imp$`INCARCERATION REASON` %in% c("Arrest - Without Warrant","Arre
CCSO_imp$`INCARCERATION REASON`[CCSO_imp$`INCARCERATION REASON` %in% c("FTA - CITY WARRANT (OV)","FTA -
CCSO_imp$`INCARCERATION REASON`[CCSO_imp$`INCARCERATION REASON` %in% c("Sentenced - EHD","Sentenced")] :
CCSO_imp$`INCARCERATION REASON`[!(CCSO_imp$`INCARCERATION REASON` %in% c("Arrest","FTA","Sentenced"))] :

## Category Recoding
## 1. only keep arrest data; other categories are sparse, unknown, or have
## a disruptive effect on modeling (sentencing for example)
##
## 2. remove Conditional Release, Served/Acquitted, and Transferred from
## release reason. These levels seem to disrupt the model (have abnormally
## high propensity for incarcerations lasting one or more days)
CCSO_small <- CCSO_imp %>% rename(Inc_Reason = "INCARCERATION REASON",
Emp_Status = "EMPLOYMENT STATUS",
Rls_Reason = "RELEASED REASON") %>%
filter(Inc_Reason %in% "Arrest") %>%
filter(!Rls_Reason %in% c("Conditional Release", "Served/Acquitted",
"Transferred")) %>%
dplyr::select(-Inc_Reason)

```

Next, we add two derived features as directed in the problem: One denoting whether the offender did multiple offenses in the same day (called `multiple`), and the other denoting whether they are repeat offenders (called `repeatOffender`)

We show three rows from our final dataset. It has 5193 observations and 9 features

##	atleastone	Age	Sex	Race	Agency	Rls_Reason
## 1:	0	22	Male	White	Champaign Police Department	Other
## 2:	0	26	Male	White	Champaign County Sherriff's Office	Other
## 3:	0	32	Female	White	Champaign Police Department	Released
##	Emp_Status	multiple	repeat	Offender		
## 1:	Student/Retired		0			
## 2:	Employed		0			
## 3:	Employed		0			

Let us fit our full model here

Let us first test the goodness of the model against the null model and a model with no covariates (iid Bernoulli trials)

[illegible]

```
## Start:  AIC=2839.14
## atleastone ~ -1 + (Age + Sex + Race + Agency + Rls Reason + Emp Status +
```



```

##      multiple + repeatOffender)
##
##              Df Deviance    AIC
## - Age          1   2805.6 2837.6
## - Rls_Reason    1   2806.5 2838.5
## <none>          1   2805.1 2839.1
## - repeatOffender 1   2807.4 2839.4
## - multiple      1   2813.5 2845.5
## - Agency        5   2846.3 2870.3
## - Sex           2   2849.0 2879.0
## - Race          3   2892.9 2920.9
## - Emp_Status    3   2904.5 2932.5
##
## Step:  AIC=2837.64
## atleastone ~ Sex + Race + Agency + Rls_Reason + Emp_Status +
##      multiple + repeatOffender - 1
##
##              Df Deviance    AIC
## - Rls_Reason    1   2806.9 2836.9
## <none>          1   2805.6 2837.6
## - repeatOffender 1   2807.9 2837.9
## - multiple      1   2813.9 2843.9
## - Agency        5   2847.7 2869.7
## - Sex           2   2849.2 2877.2
## - Race          3   2894.1 2920.1
## - Emp_Status    3   2905.3 2931.3
##
## Step:  AIC=2836.93
## atleastone ~ Sex + Race + Agency + Emp_Status + multiple + repeatOffender -
##      1
##
##              Df Deviance    AIC
## <none>          1   2806.9 2836.9
## - repeatOffender 1   2809.2 2837.2
## - multiple      1   2815.5 2843.5
## - Agency        5   2849.7 2869.7
## - Sex           2   2849.6 2875.6
## - Race          3   2896.7 2920.7
## - Emp_Status    3   2907.7 2931.7
summary(mod_Select)

##
## Call:
## glm(formula = atleastone ~ Sex + Race + Agency + Emp_Status +
##      multiple + repeatOffender - 1, family = "binomial", data = CCSO_small,
##      x = "TRUE")
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.2735  -0.4534  -0.3169  -0.2628   2.7644
##
## Coefficients:
##
##                               Estimate Std. Error z value
## SexFemale                    -3.34684    0.66755  -5.014

```

```
## SexMale -2.88231 0.65704 -4.387
## RaceBlack 1.59003 0.52130 3.050
## RaceHispanic 0.41957 0.54283 0.773
## RaceWhite 0.69330 0.52209 1.328
## AgencyChampaign Police Department 0.14061 0.15919 0.883
## AgencyIllinois State Police 0.26231 0.15085 1.739
## AgencyRantoul Police Department 1.21401 0.19390 6.261
## AgencyUniversity of Illinois Police Department 0.00388 0.20902 0.019
## AgencyUrbana Police Department 0.47956 0.17086 2.807
## Emp_StatusEmployed -0.77713 0.39418 -1.972
## Emp_StatusStudent/Retired -0.75458 0.43082 -1.751
## Emp_StatusUnemployed 0.30138 0.39581 0.761
## multiple -0.36809 0.12637 -2.913
## repeatOffender 0.18779 0.12567 1.494
## Pr(>|z|)
## SexFemale 0.000000534185 ***
## SexMale 0.000011503639 ***
## RaceBlack 0.00229 **
## RaceHispanic 0.43957
## RaceWhite 0.18420
## AgencyChampaign Police Department 0.37707
## AgencyIllinois State Police 0.08207 .
## AgencyRantoul Police Department 0.000000000382 ***
## AgencyUniversity of Illinois Police Department 0.98519
## AgencyUrbana Police Department 0.00500 **
## Emp_StatusEmployed 0.04866 *
## Emp_StatusStudent/Retired 0.07986 .
## Emp_StatusUnemployed 0.44640
## multiple 0.00358 **
## repeatOffender 0.13511
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 7199.0 on 5193 degrees of freedom
## Residual deviance: 2806.9 on 5178 degrees of freedom
## AIC: 2836.9
##
## Number of Fisher Scoring iterations: 6
```

The Age and ‘release reason’ features are removed. Let us examine if this is a good fit. We can notice this model is significantly better than the null model and the full model is not significantly better than this model. We are satisfied with this selected model

```
#Testing against the saturated model
```

```
pchisq(mod_Select$deviance - mymod$deviance, df = mod_Select$df.residual - mymod$df.residual, lower = F
```

```
## [1] 0.4084187
```

```
#Testing against the null model
```

```
pchisq(mymod$null.deviance - mod_Select$deviance, df = mymod$df.null - mod_Select$df.residual,
lower = FALSE)
```

```
## [1] 0
```

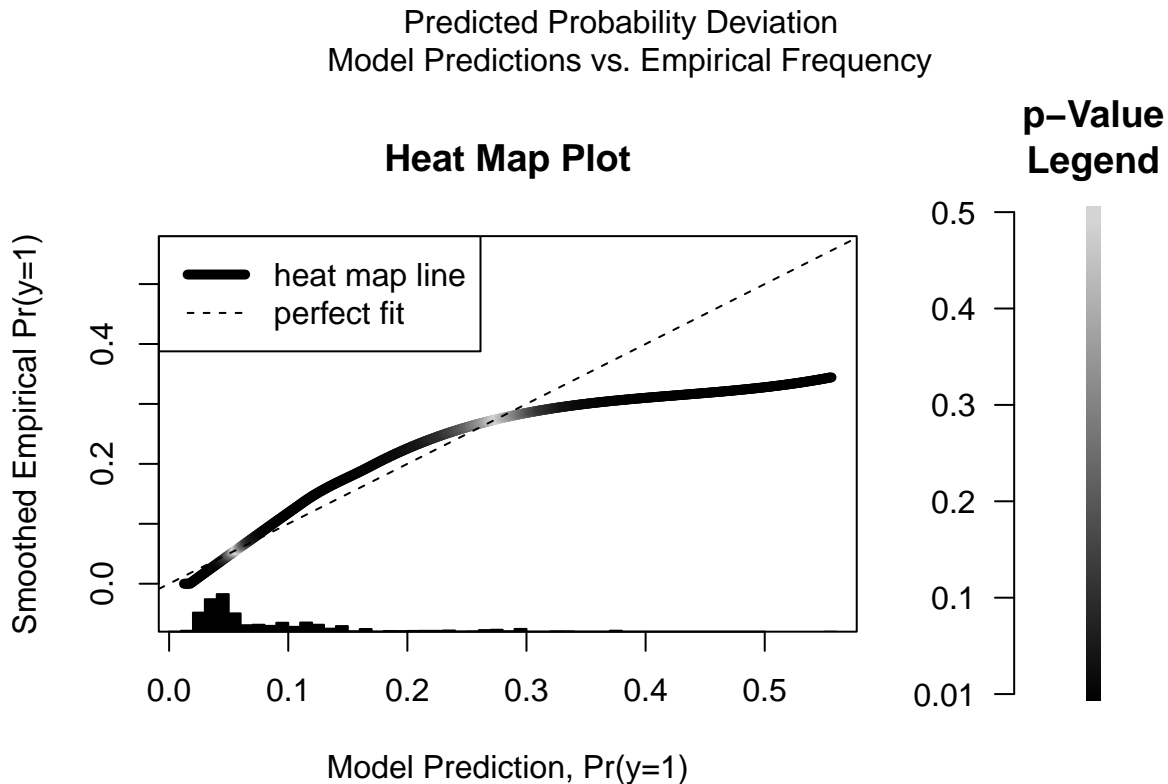
However, our graphical diagnostics tell a different story. It appears that there is some lack of model fit over

the bulk of the data and some individuals with predictive probability close to 1.

```
library(heatmapFit)
```

```
y <- CCSO_small$atleastone  
p1 <- predict(mod_Select, type = "response", se.fit = FALSE)  
heatmap.fit(y, p1)
```

```
##  
## Calculating optimal loess bandwidth...  
## aicc Chosen Span = 0.9504823  
##  
## Generating Bootstrap Predictions...  
## |
```



```
##  
##  
## *****  
## 74.48488% of Observations have one-tailed p-value <= 0.1  
## Expected Maximum = 20%  
## *****
```

The variable multiple offenses is inherently vague in nature. Some offenses can be more severe than others so this might disrupt our model. Let's remove everyone with multiple offenses and restrict attention to traffic offenses only. When we do this, the model is better fitting according to our visual diagnostics and it appears that some previous conclusions change (note that violations in fit appear where there is not much data and is extrapolated to where there is no data). Namely, race exhibits a much weaker association with propensity for incarcerations lasting one or more days, and repeat offenders seem to be punished more.

```

m1 <- glm(atleastone ~ Sex + Race + Agency + repeatOffender,
  family = "binomial", data = CCSO_small %>% filter(multiple == 0))

m2 <- glm(atleastone ~ Sex + Agency + repeatOffender,
  family = "binomial", data = CCSO_small %>% filter(multiple == 0))
anova(m2, m1, test = "LRT")

## Analysis of Deviance Table
##
## Model 1: atleastone ~ Sex + Agency + repeatOffender
## Model 2: atleastone ~ Sex + Race + Agency + repeatOffender
##   Resid. Df Resid. Dev Df Deviance    Pr(>Chi)
## 1      1775      1292.2
## 2      1772      1256.1  3    36.117 0.00000007074 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

heatmap.fit(CCSO_small %>% filter(multiple == 0) %>% pull(atleastone),
  predict(m1, type = "response", se.fit = FALSE))

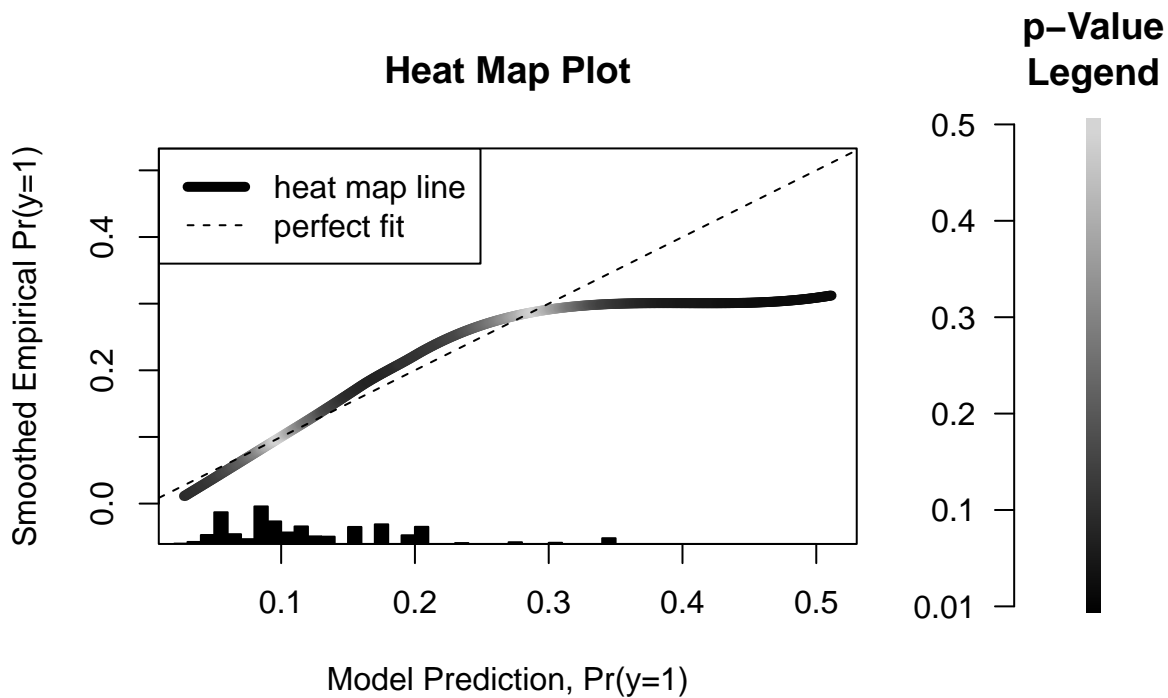
```

```

##
## Calculating optimal loess bandwidth...
## aicc Chosen Span = 0.9584928
##
## Generating Bootstrap Predictions...
## |

```

Predicted Probability Deviation Model Predictions vs. Empirical Frequency



```

##
##
## *****

```

```
## 15.5917% of Observations have one-tailed p-value <= 0.1
## Expected Maximum = 20%
## *****
```

```
## compare with saturated
```

```
pchisq(m1$deviance, m1$df.residual, lower = FALSE)
```

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

Summary table for final model

```
summary(m1)
```

```
##
```

```
## Call:
```

```
## glm(formula = atleastone ~ Sex + Race + Agency + repeatOffender,
##      family = "binomial", data = CCSO_small %>% filter(multiple ==
##      0))
##
```

```
## Deviance Residuals:
```

```
##      Min       1Q   Median       3Q      Max
## -1.1960  -0.5328  -0.4264  -0.3339   2.4723
##
```

```
## Coefficients:
```

```
##
##              Estimate Std. Error z value
## (Intercept)    -3.47204    0.76869  -4.517
## SexMale         0.54505    0.19623   2.778
## RaceBlack       1.23425    0.74135   1.665
## RaceHispanic   -0.08086    0.77987  -0.104
## RaceWhite       0.55369    0.74313   0.745
## AgencyChampaign Police Department 0.14864    0.22179   0.670
## AgencyIllinois State Police       0.24726    0.25868   0.956
## AgencyRantoul Police Department   1.03533    0.28236   3.667
## AgencyUniversity of Illinois Police Department 0.02118    0.30176   0.070
## AgencyUrbana Police Department    0.33929    0.23078   1.470
## repeatOffender 0.70111    0.26443   2.651
##
```

```
##              Pr(>|z|)
## (Intercept) 0.00000628 ***
## SexMale      0.005475 **
## RaceBlack    0.095936 .
## RaceHispanic 0.917424
## RaceWhite    0.456224
## AgencyChampaign Police Department 0.502732
## AgencyIllinois State Police       0.339133
## AgencyRantoul Police Department   0.000246 ***
## AgencyUniversity of Illinois Police Department 0.944039
## AgencyUrbana Police Department    0.141502
## repeatOffender 0.008017 **
## ---
```

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

```
##
```

```
## (Dispersion parameter for binomial family taken to be 1)
```

```
##
```

```
##      Null deviance: 1324.5  on 1782  degrees of freedom
```

```
## Residual deviance: 1256.1  on 1772  degrees of freedom
```

```
## AIC: 1278.1
```

##

Number of Fisher Scoring iterations: 5

Discussion Now, finally - if we look at the summary table - the features that are most strongly associated with propensity of incarcerations lasting one or more days:

- – Race seems to be a factor, and being black seems to be associated with longer incarcerations. However, the association between race and propensity of incarcerations lasting one or more days is not very strong when we consider important confounding factors.
- – One agency ‘Rantoul Police Department’ stands out among all agencies. It appears people arrested by this agency tends to have longer incarcerations.
- – Repeat offenses is strongly associated with longer incarcerations.
- – Males are strongly associated with longer incarcerations.

We should remember this is observational data and we can not make causal connections with confidence.