Olga Vitek

# Homework 7 - Solution

Each part of the problems 5 points

1. *[Methods qualifying exam, August 2009: use paper and pencil.]* Data were collected to study a type of damage caused by waves to the forward section of certain cargo-carrying vessels. The values of the following variables were recorded:

   type: ship type (coded as A to E);

   year: year of construction (1960, 1965, 1970, 1975);

   period: period of operation (1960-74, 75-79);

   service: aggregate months of service

   incidents: number of damage incidents

```
Call:
glm(formula = incidents ~ offset(log(service)) + type + year + period, family=poisson, data=wdships)

Deviance Residuals:
Min 1Q Median 3Q Max
-1.6768 -0.8293 -0.4370 0.5058 2.7912

Coefficients:
             Estimate Std. Error   z value     Pr(>|z|)
(Intercept) -6.40590    0.21744   -29.460      < 2e-16 ***
typeB       -0.54334    0.17759    -3.060      0.00222 **
typeC       -0.68740    0.32904    -2.089      0.03670 *
typeD       -0.07596    0.29058    -0.261      0.79377
typeE        0.32558    0.23588     1.380      0.16750
year65       0.69714    0.14964     4.659    3.18e-06 ***
year70       0.81843    0.16977     4.821    1.43e-06 ***
year75       0.45343    0.23317     1.945      0.05182 .
period75     0.38447    0.11827     3.251      0.00115 **

(Dispersion parameter for poisson family taken to be 1)

Null deviance:     146.328 on 33 degrees of freedom
Residual deviance: 38.695 on 25 degrees of freedom
AIC: 154.56
```

   (a) Describe the model used in the above analysis (in terms of the distribution for the response variable, and the relation of independent variable to the response).

   **Answer:**

   The distribution for the response is

   $$Y_i \overset{ind}{\sim} Poisson(\lambda_i)$$

   $$\log\lambda_i = \log(\texttt{Service}_i) + \sum_{j=2}^{5} \beta_{ij}^t \texttt{type}_{ij} + \sum_{k=2}^{4} \beta_{ik}^y \texttt{year}_{ik} + \beta_{i2}^p \texttt{period}_{i2}$$

where parameter for baseline $(\beta_A^t \texttt{type}_A, \beta_{60}^y \texttt{year}_{60}, \beta_{60-74}^p \texttt{period}_{60-74})$

(b) Give the best reason to use `offset(log(service))` in the above data analysis.

**Answer:**

Service measures the aggregate number of months each ship has been in service. This is a measure of the amount of exposure each ship has had to wave damage and hence represents a different interval in time for each response (the number of damage incidents). It is accounted for as an offset in the Poisson regression context.

(c) Does the model fit the data well? Please explain.

**Answer:**

$$H_0 : \log \lambda_i = \log(\texttt{Service}_i) + X'\beta$$
$$H_a : \log \lambda_i \neq \log(\texttt{Service}_i) + X'\beta$$

residual deviance= $38.695 > \chi^2_{25,0.95} = 37.65$ Therefore, reject $H_0$ and we conclude that the model does not fit well.

(d) Should we consider allowing the dispersion parameter to vary? Please explain.

**Answer:**
$\hat{\phi} = \frac{38.695}{25} = 1.5478 > 1$ So, the dispersion parameter should be allowed to vary.

(e) We define `year60 <- (year==1960)`, and rerun the above model with `year60` replacing `year`. The residual deviance is reported as 42.329. Please construct a hypothesis test to state which of these two models is better.

**Answer:**

$$H_0 : \text{new model with "year60"}$$
$$H_a : \text{stated model with "year"}$$

test statistic= $42.329 - 38.693 = 3.634 < \chi^2_{2,0.95} = 5.99$ Therefore, fail to reject $H_0$ and we conclude that replacing "year" with "year60" leads to a better-fitting model.

2. *[Methods qualifying exam, January 2007: use paper and pencil.]*

The following table contains the counts of leukemia in five counties of California from 1996 to 2000 as well as the county population in 2000.

| County | Count ($n_i$) | Population ($x_i$) |
|---|---|---|
| Marin | 22 | 247289 |
| Contra Costa | 146 | 948816 |
| Alameda | 226 | 1443741 |
| San Francisco | 47 | 776733 |
| San Mateo | 52 | 70716 |

Assume the count $N_i$ with observed value $n_i$ follows a Poisson distribution with parameter $\lambda_i$ times the county population $x_i, i = 1, \ldots, 5$. Derive the Pearson and deviance goodness-of-fit tests for the null hypothesis that the $\lambda_i$ are all the same.

**Answer:**

Here we assume $n_i \sim Poisson(\theta x_i)$. Then, the MLE of $\theta$ is

$$\hat{\theta} = \frac{\sum_{i=1}^{5} n_i}{\sum_{i=1}^{5} x_i} = 0.0001414$$

Thus, the predicted values of $n_i(\hat{n}_i)$ are 34.96, 134.13, 204.10, 109.81, 10.00. The Pearson $\chi^2$ statistic is

$$X^2 = \sum_{i=1}^{5} \frac{(n_i - \hat{n}_i)^2}{\hat{n}_i} = 220.6 > \chi^2_{4, o.95} = 9.49$$

Therefore reject $H_0$ and we conclude that at least one county has a different $\lambda_i$ from the others.

and the deviance goodness of fit is

$$G^2 = 2 \sum_{i=1}^{5} n_i \log \frac{n_i}{\hat{n}_i} = 142.14$$

Therefore, reject $H_0$ and we have the same conclusion.

3. *[Methods qualifying exam, January 2005: use paper and pencil.]* During nesting, each female horseshoe crab has a single male crab in the nest but other males reside nearby. These other males are known as satellite crabs. A study was done to investigate whether the size of the female influences the number of satellite crabs.

The investigators measured each females carapace width and grouped them by 8 width intervals (the midpoints of these intervals are given in the table). The total number of female crabs is given as `cases` and the total number of satellite crabs is given by `satell`. The sample mean and variance within each interval is also given.

| Width | Cases | Satell | Mean | Variance |
|---|---|---|---|---|
| 22.69 | 14 | 14 | 1.00 | 2.77 |
| 23.84 | 14 | 20 | 1.43 | 8.88 |
| 24.77 | 28 | 67 | 2.39 | 6.54 |
| 25.84 | 39 | 105 | 2.69 | 11.38 |
| 26.79 | 22 | 63 | 2.86 | 6.88 |
| 27.74 | 24 | 93 | 3.87 | 8.81 |
| 28.67 | 18 | 71 | 3.94 | 16.88 |
| 30.41 | 14 | 72 | 5.14 | 8.29 |

3

A log-linear model was fit to the data as follows:

```
> g <- glm(satell=width+offset(log(cases)), family=poisson)
> summary(g)
```

with output

```
Coefficients:
Estimate    Std. Error    t-value
Intercept  -3.5354702 0.57602076 -6.137748
width       0.1727192 0.02123161  8.135004
```

The null deviance is 72.37717 on 7 degrees of freedom. The residual deviance is 6.516796 on 6 degrees of freedom.

(a) Can you tell if this model fits the data?

**Answer:**

Yes, the model fits the data since residual deviance is 6.51. Based on $\chi_6^2$ distribution, it is small enough.

(b) Explain why an offset of this form was used in this model.

**Answer:**

The goal of the study is to find the relationship between the number of satellites and the width of the female crab. However the study measured locations where there were several female crabs, and the difference in the number of satellites per location can be due to the difference in the number of females. Therefore we add an offset to account for the difference in the number of females.

(c) Suppose the model `satell=width` was fit to the ungrouped data (i.e. without an offset). From the information given, state a problem that might arise.

**Answer:**

This model considers the number of cases instead of the ratio as the response variable. It will give a large predicted value for number of `satell` of crabs with a particular measurement of `width`. Result can not be explained according to the size influences the number of satellite crabs.

4. The dataset `arrests.txt` on the course website reports the number of individuals stopped during police checks in 75 precincts in New York city.

- The units $i$ are precincts and ethnic groups ($i = 1, \ldots, n = 3 \times 75$).
- The outcome of interest is `stops`, i.e. the number of stops of members of that ethnic group and precincts
- `past.arrests` is the number of past arrests of people from that ethnic group and precinct in the previous year
- `precinct` is the id of the precinct

- `eth` is the id of the ethnic group (1=african americans, 2=hispanic, 3=white)
- `crime` is the indicator of the type of the crime

One would like to know whether there is a difference in the average number of police stops by `precinct` and by `ethnicity`, while treating `past.arrests` as a measure of "exposure".

(a) Specify the model that is appropriate for this problem in mathematical notation, and state the assumption. Obtain the model fit. (*Hint:* use `offset=` option in `glm`).

**Answer:**
The response of interest is measured in counts. Therefore an appropriate model is based on Poisson distribution with an offset:

$$log[E(y)] = log(\lambda_i) = log(\mu_i\theta_i) = log(\mu_i) + log(\theta_i)$$

$$= log(\texttt{past.arrests}_i) + \beta_0 + \sum_{j=2}^{3}\beta_{ij}^{\texttt{e}}\texttt{ethnicity}_{ij} + \sum_{j=2}^{75}\beta_{ik}^{\texttt{p}}\texttt{precinct}_{ik}$$

The assumptions include

 i. $Y_i \sim Poisson(\lambda_i)$, where $\lambda_i = \mu_i\theta_i$.
 ii. $\mu_i$: the exposure, which is `past.arrests` in this case. It is multiplicative.
 iii. $\theta_i$: the rate, following Poisson regression model.
 iv. Observations are independent of each other.
 v. At each level, $E(Y_i) = Var(Y_i)$.

And the model fit is below,

### R code and output

```
mfull <- glm(formula = stops ~  factor(eth) + factor(precinct),
        family=poisson, data=stop, offset=log(past.arrests))
summary(mfull)

Coefficients:
                  Estimate Std. Error z value Pr(>|z|)
(Intercept)       -1.378863   0.051019 -27.026  < 2e-16 ***
factor(eth)2       0.010183   0.006802   1.497 0.134377
factor(eth)3      -0.419023   0.009435 -44.412  < 2e-16 ***
factor(precinct)2 -0.149049   0.074030  -2.013 0.044078 *
factor(precinct)3  0.559956   0.056758   9.866  < 2e-16 ***
factor(precinct)4  1.210638   0.057549  21.037  < 2e-16 ***
...
factor(precinct)74 1.151469   0.058023  19.845  < 2e-16 ***
factor(precinct)75 1.571238   0.075732  20.747  < 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: 183986  on 899  degrees of freedom
Residual deviance: 141292  on 823  degrees of freedom
AIC: 146785
```

(b) Conduct the null and the residual deviance tests. State the null and the alternative hypotheses, and your conclusions.

**Answer:**
For the null deviance test, we test

$$H_0 : log[E(Y_i)] = \beta_0 \text{ vs } H_a : \text{at least one of } \beta_{2j}^e \text{ or } \beta_{2k}^p \neq 0 \ \ j = 2, 3, k = 2, 3, ...75.$$

Decision rule: Reject $H_0$ hypothesis if
$G^2 =$ `Null deviance` - `Residual deviance` $> \chi^2(1 - \alpha,$ `numberExplanatoryVariables`$)$.
In this example, $G^2 = 183986 - 141292 = 42694$ with 76 degrees of freedom. The p-value of the test statistic in approximately 0, and we reject $H_0$.

Therefore, we have strong evidence to reject the null hypothesis and conclude that at least one of the predictors are significant.

For the residual deviance test, we test

$$H_0 : log[E(y)] = log(\texttt{past.arrests}_i) + \beta_0 + \sum_{j=2}^{3} \beta_{ij}^e \texttt{ethnicity}_{ij} + \sum_{j=2}^{75} \beta_{ik}^p \texttt{precinct}_{ik}$$

$$H_a : log[E(y)] \neq log(\texttt{past.arrests}_i) + \beta_0 + \sum_{j=2}^{3} \beta_{ij}^e \texttt{ethnicity}_{ij} + \sum_{j=2}^{75} \beta_{ik}^p \texttt{precinct}_{ik}$$

Decision rule: Reject $H_0$ if
$G^2 =$ `Residual deviance` $> \chi^2(1 - \alpha,$ `residualDf`$)$.
In this example: $G^2 = 141292$ based on 823 degrees of freedom. The p-value of the test is 0 and we reject $H_0$.

Therefore, we conclude that the proposed model does not fit the data well.
The deficiency can be due to the incomplete description of the structure of the data. E.g. it can be de to missing important predictors, such as crime type. This can also be due to deviations from the assumption of independent observations, which can arise when precincts have a latent clustering structure. This can be accounted for by modeling overdispersion.

(c) What are your conclusions regarding the influence of precincts and of ethnic groups on the number of stops?

**Answer:**
To test the influence of `precinct`, we compare

$$\text{Reduced Model} \quad : \quad log[\lambda_i] = log(\texttt{past.arrests}_i) + \beta_0 + \sum_{j=2}^{3} \beta_{ij}^e \texttt{ethnicity}_{ij}$$

$$\text{Full Model} \quad : \quad log[\lambda_i] = log(\texttt{past.arrests}_i) + \beta_0 + \sum_{j=2}^{3} \beta_{ij}^e \texttt{ethnicity}_{ij} + \sum_{j=2}^{75} \beta_{ik}^p \texttt{precinct}_{ik}$$

6

The decision rule is: Reject $H_0$ hypothesis if

$$G^2 = \texttt{ResidualDeviance}(\texttt{Reduced}) - \texttt{ResidualDeviance}(\texttt{Full}) > \chi^2(1 - \alpha, p - q)$$

In this example: $G^2 = 42011$, and the p-value based on the $\chi^2$ distribution, with 74 degrees of freedom is $0 < 0.05$.

Therefore, we reject the null hypothesis and conclude that at the `precint` a significant predictor of the number of police stops.

To test the influence of `ethnicity`,

$$\text{Reduced Model} \quad : \quad log[\lambda_i] = log(\texttt{past.arrests}_i) + \beta_0 + \sum_{j=2}^{75} \beta_{ik}^{\texttt{P}}\texttt{precinct}_{ik}$$

$$\text{Full Model} \quad : \quad log[\lambda_i] = log(\texttt{past.arrests}_i) + \beta_0 + \sum_{j=2}^{3} \beta_{ij}^{\texttt{e}}\texttt{ethnicity}_{ij} + \sum_{j=2}^{75} \beta_{ik}^{\texttt{P}}\texttt{precinct}_{ik}$$

The decision rule is: Reject $H_0$ hypothesis if

$$G^2 = \texttt{ResidualDeviance}(\texttt{Reduced}) - \texttt{ResidualDeviance}(\texttt{Full}) > \chi^2(1 - \alpha, p - q).$$

In this example: $G^2 = 2469.4$ and P-value based on the $\chi^2$ distribution, with degree of freedom of 2 is $0 < 0.05$.

Therefore, we reject the null hypothesis and conclude that at the ethnicity is a significant predictor of the number of police stops.

## R code and output

```
## without precint
>influ.prec <- glm(stops ~ factor(eth), family=poisson, data = stop,
              offset=log(past.arrests))
>anova(influ.prec, mfull)

Model 1: stops ~ factor(eth)
Model 2: stops ~ factor(eth) + factor(precinct)
  Resid. Df Resid. Dev Df Deviance
1       897     183303
2       823     141292 74    42011


## without ethnicity
>influ.eth <- glm(stops ~ factor(precinct), family=poisson, data = stop,
              offset=log(past.arrests))
>anova(influ.eth, mfull)
Model 1: stops ~ factor(precinct)
Model 2: stops ~ factor(eth) + factor(precinct)
  Resid. Df Resid. Dev Df Deviance
1       825     143762
2       823     141292  2   2469.4
```

(d) For each precinct and ethnicity, the dataset contains separate entries for crime type. Since crime type is not used as a predictor, these entries are viewed as independent replicates. Will there be

a change in (i) the values of null and residual deviance, (ii) the estimates of the parameters, and (iii) comparisons between nested models if we combine all the entries from a same precinct and ethnicity (i.e. add the number of stops and past arrests across all the crime types)? Explain the reasons for your answers, and show numeric output to support your conclusions.

**Answer:**

i. null and residual deviance are different.

|  | before combining, (a) | after combining, (d) |
|---|---|---|
| null deviance | 183986 | 46120.9 |
| df | 899 | 224 |
| residual deviance | 141292 | 3427.1 |
| df | 823 | 148 |

ii. parameter estimates are the same.

|  | before combining, (a) | after combining, (d) |
|---|---|---|
| (Intercept) | 1-1.378863 | -1.378863 |
| factor(eth)2 | 0.010183 | 0.010183 |
| factor(eth)3 | -0.419023 | -0.419023 |
| factor(precint)2 | -0.149049 | -0.149049 |
| factor(precint)3 | 0.559956 | 0.559956 |
| $\vdots$ | $\vdots$ | $\vdots$ |
| factor(precint)74 | 1.151469 | 1.151469 |
| factor(precint)75 | 1.571238 | 1.571238 |

iii. comparison between nested models is the same.

|  | before combining, (a) | after combining, (d) |
|---|---|---|
| $G^2$ | 42694 with 76 df | 42693.8 with 76 df |
| p-value | $\approx 0$ | $\approx 0$ |

Since this only changes constants in the log-likelihood equation, it only makes change in (i).

## R code and output

```
>stops.sum <- as.vector(t( tapply(stop$stops, list(stop$precinct,stop$eth), sum) ))
>past.arrests.sum <- as.vector(t( tapply(stop$past.arrests, list(stop$precinct,stop$eth), sum) ))
>XX <- data.frame(unique(stop[,c("pop", "precinct", "eth")]),
    stops=stops.sum, past.arrests=past.arrests.sum)

>fit.sum <- glm(stops~ factor(eth) + factor(precinct), family=poisson, data=XX, offset=log(past.arrests))
>summary(fit.sum)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-11.1396  -3.0893  -0.1934   2.0990  10.4185

Coefficients:
                 Estimate Std. Error z value Pr(>|z|)
(Intercept)      -1.378863   0.051019 -27.026  < 2e-16 ***
factor(eth)2      0.010183   0.006802   1.497 0.134377
factor(eth)3     -0.419023   0.009435 -44.412  < 2e-16 ***
factor(precinct)2 -0.149049   0.074030  -2.013 0.044078 *
factor(precinct)3  0.559956   0.056758   9.866  < 2e-16 ***
```

```
...
factor(precinct)73  0.991018   0.053585  18.494  < 2e-16 ***
factor(precinct)74  1.151469   0.058023  19.845  < 2e-16 ***
factor(precinct)75  1.571238   0.075731  20.747  < 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: 46120.9  on 224  degrees of freedom
Residual deviance:  3427.1  on 148  degrees of freedom
AIC: 5287.8
```

(e) Answer the same questions as in (d), while removing the offset from the model. Explain the reasons for your answers, and show numeric output to support your conclusions.

**Answer:**

i. null and residual deviance are different.

|                   | before combining | after combining |
|-------------------|-----------------:|----------------:|
| null deviance     |           182217 |          123333 |
| df                |              899 |             224 |
| residual deviance |           113372 |           54487 |
| df                |              823 |             148 |

ii. parameter estimates are the same except for the intercept, which will reflect different baselines.

|                    | before combining, (a) | after combining, (d) |
|--------------------|----------------------:|---------------------:|
| (Intercept)        |              3.934514 |             5.320809 |
| factor(eth)2       |             -0.447714 |            -0.447714 |
| factor(eth)3       |             -1.414281 |            -1.414281 |
| factor(precint)2   |             -0.103919 |            -0.103919 |
| factor(precint)3   |              1.426389 |             1.426389 |
| ⋮                  |                     ⋮ |                    ⋮ |
| factor(precint)74  |              1.237433 |             1.237433 |
| factor(precint)75  |             -0.178692 |            -0.178692 |

iii. comparison between nested models is the same. While the absolute values of the deviances are different, the difference in residual deviances between nested model remains the same.

|          | before combining, (a) | after combining, (d) |
|----------|----------------------:|---------------------:|
| $G^2$    |        68846 with 76 df |       68846 with 76 df |
| p-value  |              $\approx 0$ |            $\approx 0$ |

## R code and output

```
> mfull2 <- glm(formula = stops ~  factor(eth) + factor(precinct),
         family=poisson, data=stop)
> summary(mfull2)

Deviance Residuals:
    Min      1Q   Median      3Q      Max
-26.137  -8.894   -3.454   3.095   54.180
```

```
Coefficients:
                   Estimate Std. Error  z value Pr(>|z|)
(Intercept)        3.934514   0.051031   77.101  < 2e-16 ***
factor(eth)2      -0.447714   0.006061  -73.872  < 2e-16 ***
factor(eth)3      -1.414281   0.008558 -165.263  < 2e-16 ***
factor(precinct)2 -0.103919   0.074022   -1.404 0.160352
factor(precinct)3  1.426389   0.056756   25.132  < 2e-16 ***
...
factor(precinct)74 1.237433   0.057888   21.376  < 2e-16 ***
factor(precinct)75 -0.178692  0.075514   -2.366 0.017965 *
---
Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1   1


(Dispersion parameter for poisson family taken to be 1)


    Null deviance: 182217  on 899  degrees of freedom
Residual deviance: 113372  on 823  degrees of freedom
AIC: 118864


> fit.sum2 <- glm(stops~ factor(eth) + factor(precinct), family=poisson, data=XX)
> summary(fit.sum2)


Deviance Residuals:
    Min      1Q   Median      3Q      Max
-36.603  -12.500   -2.412   8.935   38.919


Coefficients:
                   Estimate Std. Error  z value Pr(>|z|)
(Intercept)        5.320809   0.051031  104.267  < 2e-16 ***
factor(eth)2      -0.447714   0.006061  -73.872  < 2e-16 ***
factor(eth)3      -1.414281   0.008558 -165.263  < 2e-16 ***
factor(precinct)2 -0.103919   0.074022   -1.404 0.160352
factor(precinct)3  1.426389   0.056756   25.132  < 2e-16 ***
...
factor(precinct)74 1.237433   0.057888   21.376  < 2e-16 ***
factor(precinct)75 -0.178692  0.075517   -2.366 0.017969 *
---
Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1   1


(Dispersion parameter for poisson family taken to be 1)


    Null deviance: 123333  on 224  degrees of freedom
Residual deviance:  54487  on 148  degrees of freedom
AIC: 56348
```

(f) For the following questions, consider the original dataset and the model with offset specified in (a). What is the estimated expected number of stops by police for white individuals from the first precinct?
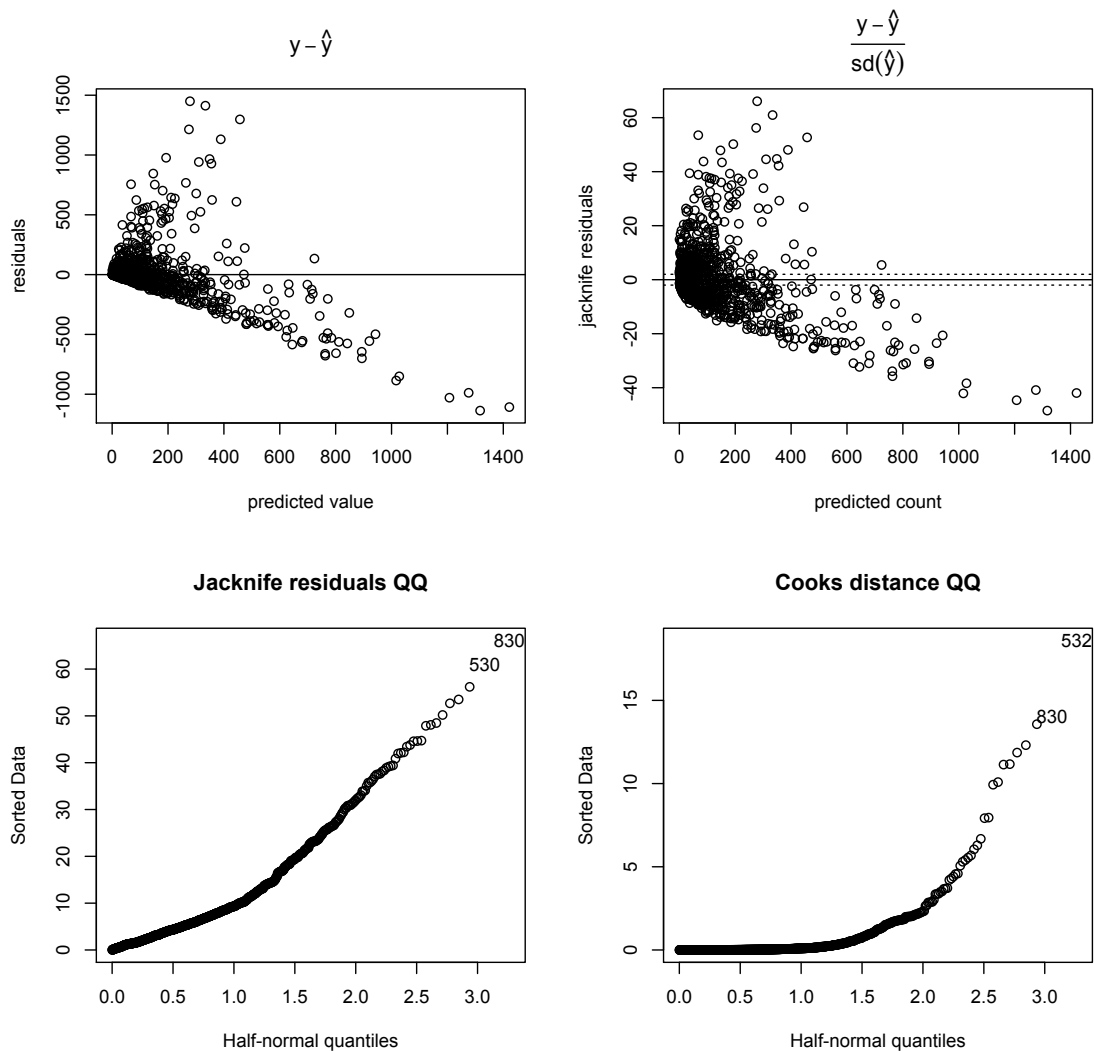
**Answer:**
Given $X_{23} = 1, X_{22} = 0, and$ all the $X_{3k} = 0$, the expected rate of stops given the model is $Exp\{-1.379 - 0.419\} = 0.1656$. To obtain the offset, we combine the number of past arrests over all types of crimes for white individuals from the first precinct, i.e. the offset $=$ 135+16+107+123=381, (where 135, 16, 107 and 123 are the numbers of past arrests in the 4 replicates with these values of the covariates). Therefore, the expected number of stops is $381 \times 0.1656 = 63.0936$

(g) Using diagnostics plots, explore the quality of model specification, and discuss the presence of influential observations or of outliers.

**Answer:**

Since the predictors are categorical, the quality of model specification focuses on residual variation, presence of outliers and on overdispersion. The residual plot on the response scale displays non constant variance (as expected in the case of the Poisson model) of residuals. However it also shows systematic over- or under-prediction of observations with a large predicted mean. The pattern persists when standardizing the residuals. The pattern does not change if deviance residuals are used instead. This points to the fact that other important predictors should be included.

The QQ plots of residuals and of Cooks distance show potentially outlying observations. However these observations should be examined after the deficiencies of functional form are addressed.



(h) Test for overdispersion and state your conclusions. What could be the reasons for overdispersion in this problem?

**Answer:**
The estimate of overdispersion is obtained using

$$\hat{\phi} = \frac{X^2}{\texttt{ResidualDF}}$$

where $X^2$ is the generalized Pearson $X^2$ statistic. In the case of Poisson model, it can be calculated as

$$X^2 = \sum_i (\frac{y_i - \hat{y}_i}{\sqrt{\hat{y}_i}})^2$$

(or equivalently, as a sum of squared Pearson residuals)

In this example, $\hat{\phi} = 260.9587$, the P-value of $\chi^2$ at 260.9587 with df as 823 is 1. Therefore, we reject the null hypothesis of no overdispersion.

There are serval potential reasons for the existance of overdispersion.

   i. Poisson distribution is not actually appropriate for this model. The variance of the observations do not equal to the mean.

   ii. The observations may be non-independent, and cluster according to a latent variable.

## R code and output

```
>yhat <- predict (mfull, type="response")
>z <- (stop$stops-yhat)/sqrt(yhat)

>cat ("overdispersion ratio is ", sum(z^2)/(823), "\n")
>cat ("p-value of overdispersion test is ", pchisq (sum(z^2), 823), "\n")

overdispersion ratio is   260.9587
p-value of overdispersion test is   1
```

(i) What would be a recommended correction for overdispersion in this case? Can you correct for overdispersion without fitting a new model?

**Answer:**
A possible way to correct overdispersion is to multiply all the standard errors of the parameter estimates by $\sqrt{\hat{\phi}} = 16.61$. After the multiplication, the estimated variances are given in the following list.

- The parameter estimates: they are unchanged,
- Estimated (co)variances: multiplied by $\sqrt{\hat{\phi}} = 16.61$.
- Residuals deviance: scaled (=divided) by $\sqrt{\hat{\phi}} = 16.61$.