

# Homework 04

## Generalized Linear Models

*Kaiyu Yan*

*October 2, 2018*

## Data analysis

### Poisson regression:

The folder `risky.behavior` contains data from a randomized trial targeting couples at high risk of HIV infection. The intervention provided counseling sessions regarding practices that could reduce their likelihood of contracting HIV. Couples were randomized either to a control group, a group in which just the woman participated, or a group in which both members of the couple participated. One of the outcomes examined after three months was “number of unprotected sex acts”.

1. Model this outcome as a function of treatment assignment using a Poisson regression. Does the model fit well? Is there evidence of overdispersion?

```
# First round fupacts.
risky_behaviors$fupacts <- round(risky_behaviors$fupacts)
#We fit the model with constant term alone.
fit1 <- glm(fupacts ~ 1, data = risky_behaviors, family = poisson)
summary(fit1)

##
## Call:
## glm(formula = fupacts ~ 1, family = poisson, data = risky_behaviors)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -5.743  -5.743  -3.323   1.065  25.125
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  2.80266    0.01182   237.1  <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: 13299  on 433  degrees of freedom
## Residual deviance: 13299  on 433  degrees of freedom
## AIC: 14625
##
## Number of Fisher Scoring iterations: 6
#We fit the model by adding two indicators.
fit2 <- glm(fupacts ~ factor(women_alone)+factor(couples), data = risky_behaviors, family = poisson)
summary(fit2)

##
```

```
## Call:
## glm(formula = fupacts ~ factor(women_alone) + factor(couples),
##      family = poisson, data = risky_behaviors)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -6.6285  -4.9794  -3.2015   0.9847  27.1502
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      3.08960    0.01901  162.55 <2e-16 ***
## factor(women_alone)1 -0.57212    0.03023  -18.93 <2e-16 ***
## factor(couples)1    -0.32243    0.02737  -11.78 <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: 13299  on 433  degrees of freedom
## Residual deviance: 12925  on 431  degrees of freedom
## AIC: 14256
##
## Number of Fisher Scoring iterations: 6
"The fit looks better than null since the residual deviance reduced 374 from 13299 to 12925 "
```

```
## [1] "The fit looks better than null since the residual deviance reduced 374 from 13299 to 12925 "
```

```
#Check for overdispersion.
n1 <- nrow(risky_behaviors)
k1 <- length(fit2$coef)
yhat1 <- predict (fit2, type="response")
z1 <- (risky_behaviors$fupacts-yhat1)/sqrt(yhat1)
cat ("overdispersion ratio is ", sum(z1^2)/(n1-k1), "\n")
```

```
## overdispersion ratio is  44.13458
cat ("p-value of overdispersion test is ", pchisq (sum(z1^2), n1-k1), "\n")
```

```
## p-value of overdispersion test is  1
"In summary, the risky behavior data are overdispersed by a factor of 44.13, which is huge"
```

```
## [1] "In summary, the risky behavior data are overdispersed by a factor of 44.13, which is huge"
```

2. Next extend the model to include pre-treatment measures of the outcome and the additional pre-treatment variables included in the dataset. Does the model fit well? Is there evidence of overdispersion?

```
Risks <- risky_behaviors[risky_behaviors$bupacts >0, ]
fit3 <- glm(round(fupacts) ~ factor(women_alone)+factor(couples)+ factor(sex) + factor(bs_hiv),offset =
summary(fit3)
```

```
##
## Call:
## glm(formula = round(fupacts) ~ factor(women_alone) + factor(couples) +
##      factor(sex) + factor(bs_hiv), family = poisson, data = Risks,
##      offset = log(bupacts))
##
```

```
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -16.315   -3.165   -1.072    2.218   21.552
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -0.03222    0.02250  -1.432   0.152
## factor(women_alone)1  -0.55581    0.03043 -18.267 < 2e-16 ***
## factor(couples)1    -0.40263    0.02804 -14.362 < 2e-16 ***
## factor(sex)man      -0.11843    0.02372  -4.994 5.92e-07 ***
## factor(bs_hiv)positive -0.32512    0.03573  -9.099 < 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: 10577  on 419  degrees of freedom
## Residual deviance: 10032  on 415  degrees of freedom
## AIC: 11356
##
## Number of Fisher Scoring iterations: 6
```

```
"This fit is better than previous model"
```

```
## [1] "This fit is better than previous model"
```

```
n2 <- nrow(Risks)
k2 <- length(fit3$coef)
yhat2 <- predict (fit3, type="response")
z2 <- (Risks$fupacts-yhat2)/sqrt(yhat2)
cat ("overdispersion ratio is ", sum(z2^2)/(n2-k2), "\n")
```

```
## overdispersion ratio is  46.30971
```

```
cat ("p-value of overdispersion test is ", pchisq (sum(z2^2), n2-k2), "\n")
```

```
## p-value of overdispersion test is  1
```

```
"There still is overdispersed by a factor of 46.31"
```

```
## [1] "There still is overdispersed by a factor of 46.31"
```

3. Fit an overdispersed Poisson model. What do you conclude regarding effectiveness of the intervention?

```
fit4 <- glm(round(fupacts) ~ factor(women_alone)+factor(couples)+ factor(sex) + factor(bs_hiv),offset =
summary(fit4)
```

```
##
## Call:
## glm(formula = round(fupacts) ~ factor(women_alone) + factor(couples) +
##      factor(sex) + factor(bs_hiv), family = quasipoisson, data = Risks,
##      offset = log(bupacts))
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -16.315   -3.165   -1.072    2.218   21.552
##
## Coefficients:
```

```
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -0.03222    0.15314  -0.210  0.83349
## factor(women_alone)1 -0.55581    0.20706  -2.684  0.00756 **
## factor(couples)1    -0.40263    0.19078  -2.110  0.03542 *
## factor(sex)man      -0.11843    0.16139  -0.734  0.46346
## factor(bs_hiv)positive -0.32512    0.24316  -1.337  0.18193
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for quasipoisson family taken to be 46.30972)
##
## Null deviance: 10577  on 419  degrees of freedom
## Residual deviance: 10032  on 415  degrees of freedom
## AIC: NA
##
## Number of Fisher Scoring iterations: 6
```

"The model seems to suggest that the group in which just the woman participated is effective compared to"

```
## [1] "The model seems to suggest that the group in which just the woman participated is effective compared to"
```

4. These data include responses from both men and women from the participating couples. Does this give you any concern with regard to our modeling assumptions?

"Including men and women in the model is not consistent with the second treatment that only woman are a"

```
## [1] "Including men and women in the model is not consistent with the second treatment that only woman are a"
```

## Comparing logit and probit:

Take one of the data examples from Chapter 5. Fit these data using both logit and probit model. Check that the results are essentially the same (after scaling by factor of 1.6)

```
wells = read.table("http://www.stat.columbia.edu/~gelman/arm/examples/arsenic/wells.dat")
wells$log.arsenic = log(wells$arsenic)
#summary(wells)
logit = glm(switch ~ log(arsenic) + dist + educ, family=binomial(link="logit"), data=wells)
display(logit)

## glm(formula = switch ~ log(arsenic) + dist + educ, family = binomial(link = "logit"),
##      data = wells)
##               coef.est coef.se
## (Intercept)    0.32    0.08
## log(arsenic)    0.89    0.07
## dist          -0.01    0.00
## educ           0.04    0.01
## ---
##      n = 3020, k = 4
##      residual deviance = 3878.2, null deviance = 4118.1 (difference = 239.9)
probit = glm(switch ~ log(arsenic) + dist + educ, family=binomial(link="probit"), data=wells)
display(probit)

## glm(formula = switch ~ log(arsenic) + dist + educ, family = binomial(link = "probit"),
##      data = wells)
##               coef.est coef.se
```

```
## (Intercept)  0.19      0.05
## log(arsenic) 0.54      0.04
## dist        -0.01     0.00
## educ         0.03     0.01
## ---
## n = 3020, k = 4
## residual deviance = 3878.3, null deviance = 4118.1 (difference = 239.8)
"The coefficient of probit model are essentially the same after scaling by factor of 1.6"
## [1] "The coefficient of probit model are essentially the same after scaling by factor of 1.6"
```

## Comparing logit and probit:

construct a dataset where the logit and probit models give different estimates.

```
arsenic = runif(10,0.51,9.65)
dist = runif(10,0.387,339.53)
educ = sample(0:17,10,replace = T)

predict_data = data.frame(arsenic,dist,educ)
predict(logit,predict_data)
```

```
##          1          2          3          4          5          6
## -0.2143880  0.6557475  1.2826922  0.9402813 -2.1284207 -0.2949080
##          7          8          9         10
## -1.6592925 -1.4513781  1.3001958  1.1267057
```

```
predict(probit,predict_data)
```

```
##          1          2          3          4          5          6
## -0.1269346  0.4023063  0.7844282  0.5745490 -1.3043668 -0.1771053
##          7          8          9         10
## -1.0133589 -0.8852687  0.7976256  0.6930441
```

## Tobit model for mixed discrete/continuous data:

experimental data from the National Supported Work example are available in the folder `1alonde`. Use the treatment indicator and pre-treatment variables to predict post-treatment (1978) earnings using a tobit model. Interpret the model coefficients.

- sample: 1 = NSW; 2 = CPS; 3 = PSID.
- treat: 1 = experimental treatment group (NSW); 0 = comparison group (either from CPS or PSID) - Treatment took place in 1976/1977.
- age = age in years
- educ = years of schooling
- black: 1 if black; 0 otherwise.
- hisp: 1 if Hispanic; 0 otherwise.
- married: 1 if married; 0 otherwise.
- nodegree: 1 if no high school diploma; 0 otherwise.
- re74, re75, re78: real earnings in 1974, 1975 and 1978
- educ\_cat = 4 category education variable (1=<hs, 2=hs, 3=sm college, 4=college)

## Robust linear regression using the t model:

The csv file `congress` has the votes for the Democratic and Republican candidates in each U.S. congressional district in between 1896 and 1992, along with the parties' vote proportions and an indicator for whether the incumbent was running for reelection. For your analysis, just use the elections in 1986 and 1988 that were contested by both parties in both years.

1. Fit a linear regression (with the usual normal-distribution model for the errors) predicting 1988 Democratic vote share from the other variables and assess model fit.
2. Fit a t-regression model predicting 1988 Democratic vote share from the other variables and assess model fit; to fit this model in R you can use the `vglm()` function in the VGLM package or `tlm()` function in the hett package.
3. Which model do you prefer?

## Robust regression for binary data using the robit model:

Use the same data as the previous example with the goal instead of predicting for each district whether it was won by the Democratic or Republican candidate.

1. Fit a standard logistic or probit regression and assess model fit.
2. Fit a robit regression and assess model fit.
3. Which model do you prefer?

## Salmonella

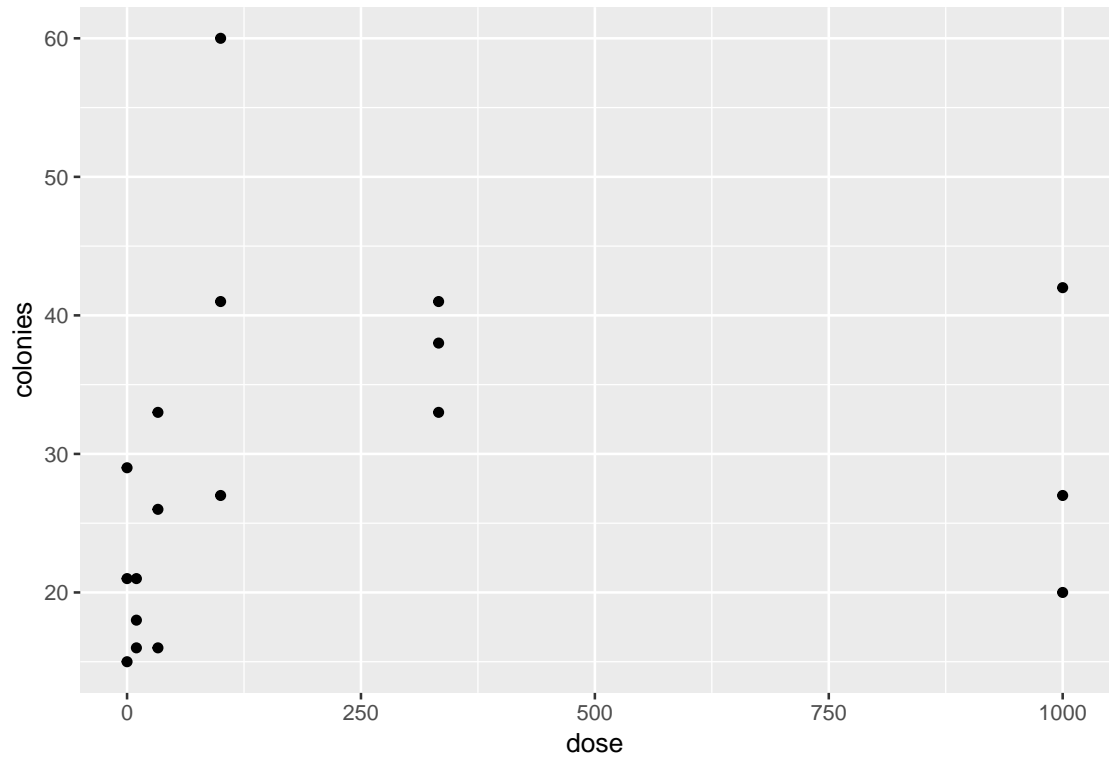
The `salmonella` data was collected in a salmonella reverse mutagenicity assay. The predictor is the dose level of quinoline and the response is the numbers of revertant colonies of TA98 salmonella observed on each of three replicate plates. Show that a Poisson GLM is inadequate and that some overdispersion must be allowed for. Do not forget to check out other reasons for a high deviance.

```
data(salmonella)
?salmoneilla
```

```
## starting httpd help server ... done
```

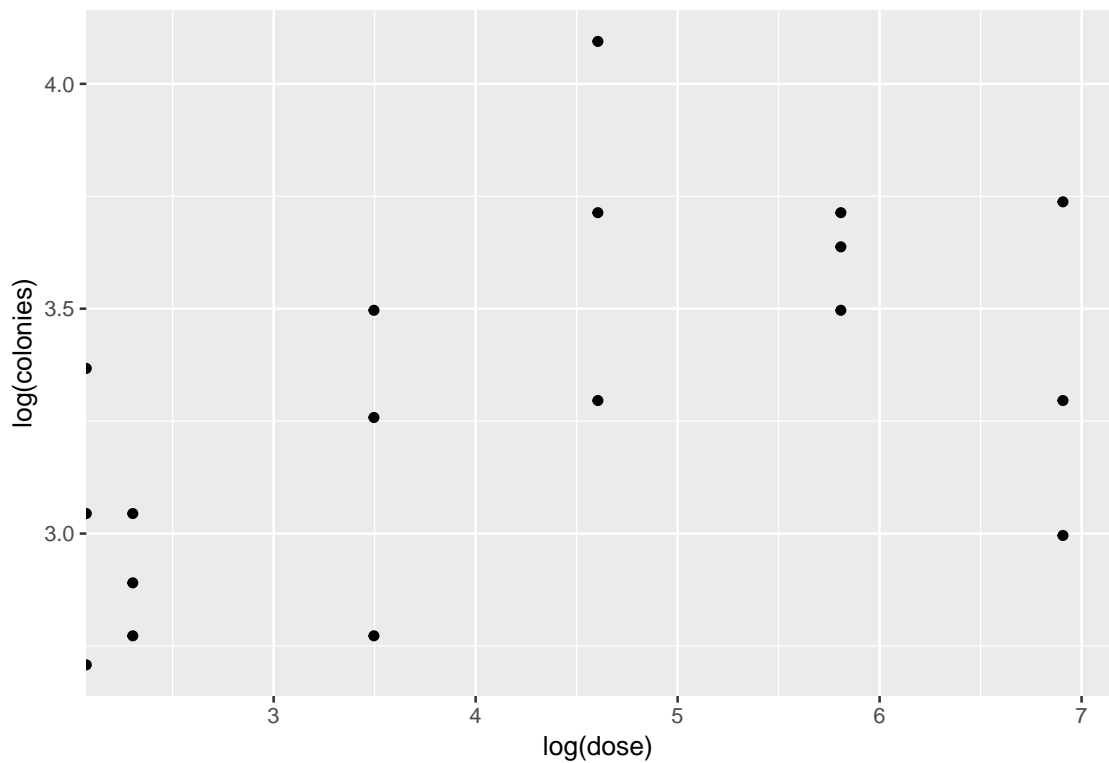
When you plot the data you see that the number of colonies as a function of dose is not monotonic especially around the dose of 1000.

```
ggplot(data = salmonella) + geom_point(aes(x = dose, y=colonies))
```



Since we are fitting log linear model we should look at the data on log scale. Also because the dose is not equally spaced on the raw scale it may be better to plot it on the log scale as well.

```
ggplot(data = salmonella) + geom_point(aes(x = log(dose), y=log(colonies)))
```

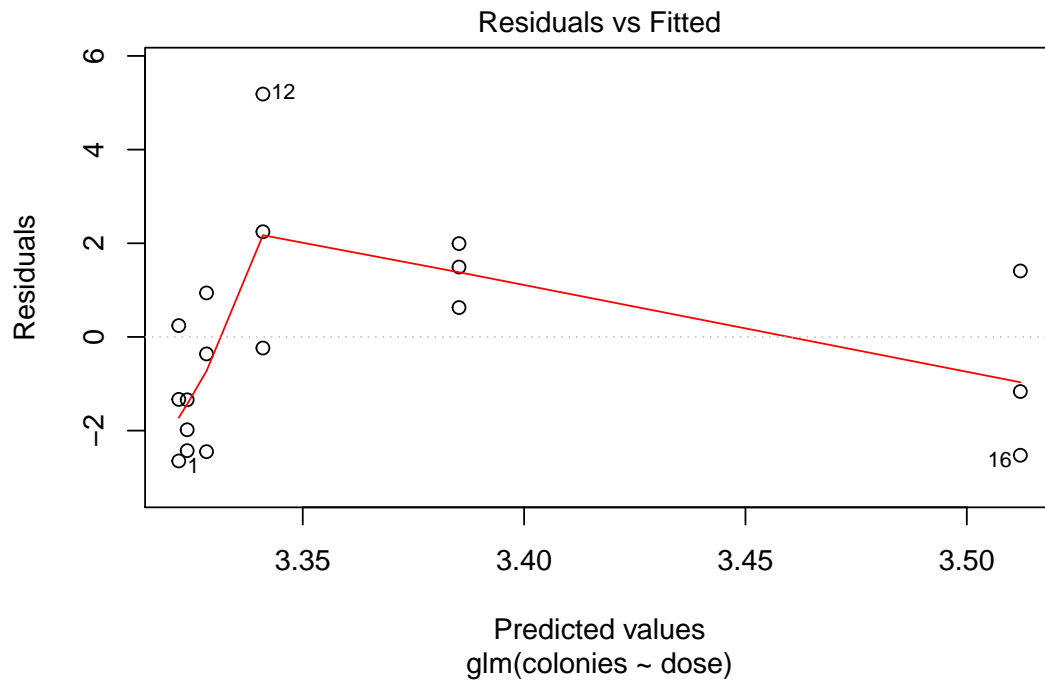


This shows that the trend is not monotonic. Hence when you fit the model and look at the residual you will see a trend.

```
salmonella_fit = glm(colonies ~ dose, data = salmonella, family = poisson(link = "log"))
summary(salmonella_fit)
```

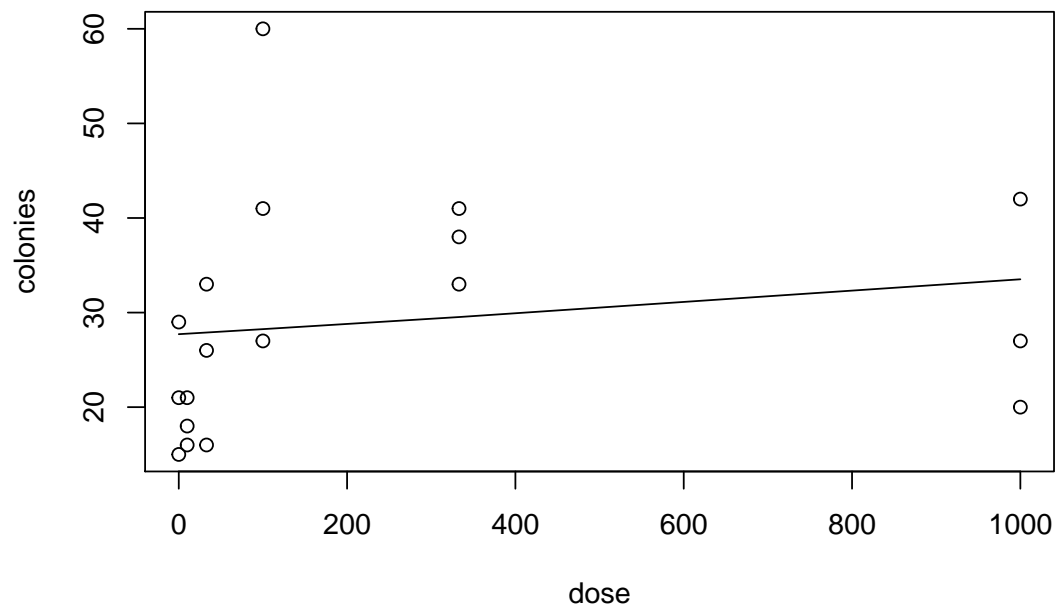
```
##
## Call:
## glm(formula = colonies ~ dose, family = poisson(link = "log"),
##      data = salmonella)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.6482  -1.8225  -0.2993   1.2917   5.1861
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  3.3219950  0.0540292  61.485   <2e-16 ***
## dose         0.0001901  0.0001172   1.622    0.105
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 78.358  on 17  degrees of freedom
## Residual deviance: 75.806  on 16  degrees of freedom
## AIC: 172.34
##
## Number of Fisher Scoring iterations: 4
plot(salmonella_fit, which = 1)
```





The lack of fit is also evident if we plot the fitted line onto the data.

```
plot(colonies ~ dose, data = salmonella)
lines(salmonella$dose, predict.glm(salmonella_fit, type="response"))
```



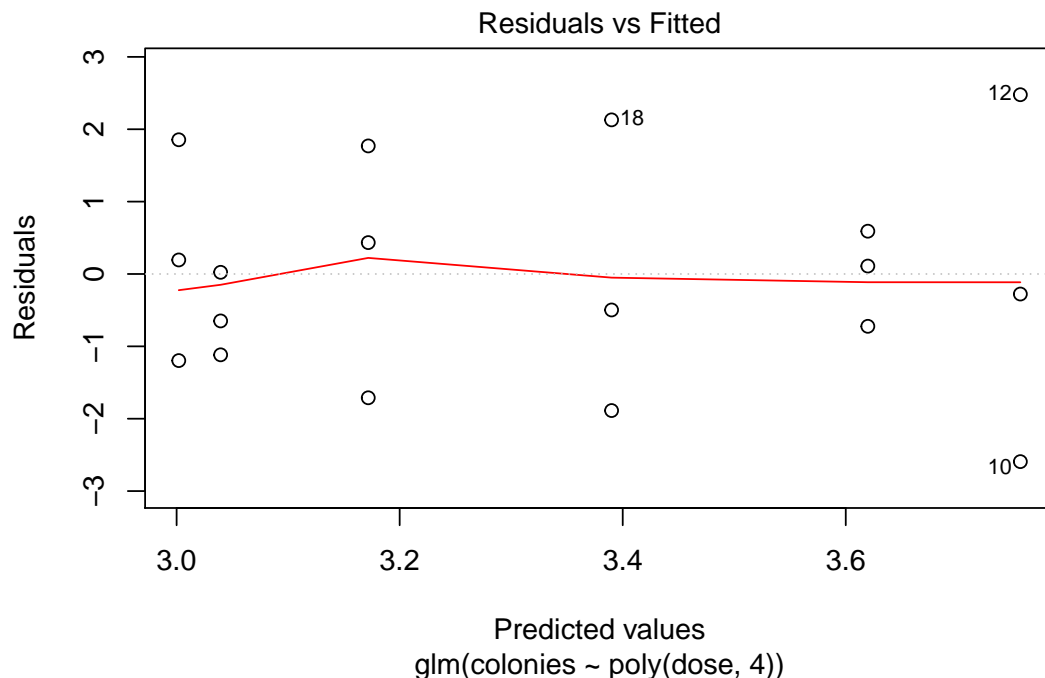
How do we address this problem? The serious problem to address is the nonlinear trend of dose rather than the overdispersion since the line is missing the points. Let's add a beny line with 4th order polynomial.

```
salmonella_fit2 = glm(colonies ~ poly(dose,4),data=salmonella,family=poisson(link="log"))
summary(salmonella_fit2)
```

```
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)    3.329932   0.045475  73.2262 < 2.2e-16
## poly(dose, 4)1  0.380047   0.190138   1.9988  0.04563
## poly(dose, 4)2 -0.853239   0.176572  -4.8322 1.350e-06
## poly(dose, 4)3  0.737453   0.172733   4.2693 1.961e-05
## poly(dose, 4)4  0.208570   0.203321   1.0258  0.30498
##
## n = 18 p = 5
## Deviance = 34.98914 Null Deviance = 78.35758 (Difference = 43.36844)
```

The resulting residual looks nice and if you plot it on the raw data. Whether the trend makes real contextual sense will need to be validated but for the given data it looks feasible.

```
plot(salmonella_fit2,which=1)
```



Dispite the fit, the overdispersion still exists so we'd be better off using the quasi Poisson model.

```
salmonella_fit3 = glm(colonies ~ poly(dose,4),data = salmonella,family=quasipoisson(link="log"))
summary(salmonella_fit3)
```

```
##
## Call:
## glm(formula = colonies ~ poly(dose, 4), family = quasipoisson(link = "log"),
##      data = salmonella)
##
## Deviance Residuals:
```

```
##      Min      1Q   Median      3Q      Max
## -2.5928 -1.0187 -0.1270  0.5518  2.4771
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    3.32993    0.07494  44.434 1.38e-15 ***
## poly(dose, 4)1  0.38005    0.31334   1.213  0.2468
## poly(dose, 4)2 -0.85324    0.29098  -2.932  0.0117 *
## poly(dose, 4)3  0.73745    0.28466   2.591  0.0224 *
## poly(dose, 4)4  0.20857    0.33506   0.622  0.5444
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for quasipoisson family taken to be 2.715769)
##
##      Null deviance: 78.358  on 17  degrees of freedom
## Residual deviance: 34.989  on 13  degrees of freedom
## AIC: NA
##
## Number of Fisher Scoring iterations: 4
```

## Ships

The `ships` dataset found in the `MASS` package gives the number of damage incidents and aggregate months of service for different types of ships broken down by year of construction and period of operation.

```
data(ships)
?ships
```

Develop a model for the rate of incidents, describing the effect of the important predictors.

```
fit_ship<- glm(incidents ~ ., family=poisson, data=ships)
summary(fit_ship)

##
## Call:
## glm(formula = incidents ~ ., family = poisson, data = ships)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -4.1013  -1.9648  -0.5380   0.9899   4.6212
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -5.706e+00  1.221e+00  -4.673 2.96e-06 ***
## typeB        8.135e-01  2.023e-01   4.021 5.79e-05 ***
## typeC       -1.205e+00  3.275e-01  -3.679 0.000234 ***
## typeD       -8.595e-01  2.875e-01  -2.989 0.002795 **
## typeE       -2.226e-01  2.348e-01  -0.948 0.343173
## year         4.519e-02  1.341e-02   3.370 0.000752 ***
## period       6.055e-02  8.945e-03   6.768 1.30e-11 ***
## service      5.970e-05  7.016e-06   8.509 < 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: 730.25   on 39   degrees of freedom
## Residual deviance: 174.00   on 32   degrees of freedom
## AIC: 287.86
##
## Number of Fisher Scoring iterations: 6
"in this fitted mode, the only predictor that not significant is type E."

## [1] "in this fitted mode, the only predictor that not significant is type E."
"The constant term gives the intercept of the regression, that means the incident is -5.706e+00 when the"

## [1] "The constant term gives the intercept of the regression, that means the incident is -5.706e+00 v
"The expected multiplicative increase of incidents is e^8.135e-01 difference of having a typeB"

## [1] "The expected multiplicative increase of incidents is e^8.135e-01 difference of having a typeB"
"The expected multiplicative increase of incidents is e^-1.205 difference of having a typeC"

## [1] "The expected multiplicative increase of incidents is e^-1.205 difference of having a typeC"
"The expected multiplicative increase of incidents is e^-8.595e-01 difference of having a typeD"

## [1] "The expected multiplicative increase of incidents is e^-8.595e-01 difference of having a typeD"
"The expected multiplicative increase of incidents is e^4.519e-02 difference of per year different"

## [1] "The expected multiplicative increase of incidents is e^4.519e-02 difference of per year differer
"The expected multiplicative increase of incidents is e^6.055e-02 difference per period different"

## [1] "The expected multiplicative increase of incidents is e^6.055e-02 difference per period differen
"The expected multiplicative increase of incidents is e^5.970e-05 difference of per service change"

## [1] "The expected multiplicative increase of incidents is e^5.970e-05 difference of per service chang
```

## Australian Health Survey

The `dvisits` data comes from the Australian Health Survey of 1977-78 and consist of 5190 single adults where young and old have been oversampled.

```
data(dvisits)
?dvisits
```

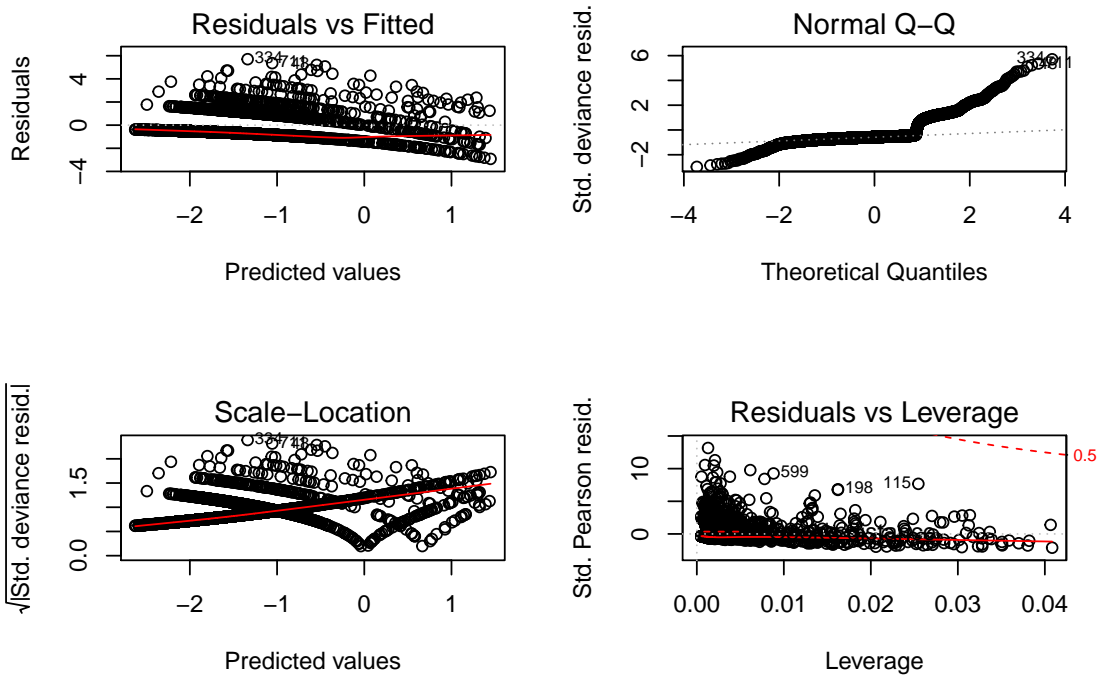
1. Build a Poisson regression model with `doctorco` as the response and `sex`, `age`, `agesq`, `income`, `levyplus`, `freepoor`, `freerepa`, `illness`, `actdays`, `hscore`, `chcond1` and `chcond2` as possible predictor variables. Considering the deviance of this model, does this model fit the data?

```
doctor_fit <- glm(doctorco ~ sex + age + agesq + income + levyplus + freepoor + freerepa + illness + ac
summary(doctor_fit)
```

```
##
## Call:
## glm(formula = doctorco ~ sex + age + agesq + income + levyplus +
##      freepoor + freerepa + illness + actdays + hscore + chcond1 +
```

```
##      chcond2, family = poisson, data = dvisits)
##
## 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
"Since the residual deviance and AIC are quite high, so it may not the best fit"

## [1] "Since the residual deviance and AIC are quite high, so it may not the best fit"
2. Plot the residuals and the fitted values-why are there lines of observations on the plot?
par(mfrow=c(2,2))
plot(doctor_fit)
```



"There are lines because the responses are discrete continuous value"

```
## [1] "There are lines because the responses are discrete continuous value"
```

3. What sort of person would be predicted to visit the doctor the most under your selected model?

"Predictors of age, income, hscore, actdays, and illness are significant, so it may the sort of person"

```
## [1] "Predictors of age, income, hscore, actdays, and illness are significant, so it may the sort of person"
```

4. 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.

```
predict(doctor_fit, dvisits[5190,], type="response")
```

```
##      5190
```

```
## 0.1533837
```

```
print(paste0("Probability of 0 doctor's visits: ", dpois(0, lambda = 0.153)))
```

```
## [1] "Probability of 0 doctor's visits: 0.858129721811394"
```

```
print(paste0("Probability of 1 doctor's visits: ", dpois(1, lambda = 0.153)))
```

```
## [1] "Probability of 1 doctor's visits: 0.131293847437143"
```

```
print(paste0("Probability of 2 doctor's visits: ", dpois(2, lambda = 0.153)))
```

```
## [1] "Probability of 2 doctor's visits: 0.0100439793289415"
```

```
print(paste0("Probability of 3 doctor's visits: ", dpois(3, lambda = 0.153)))
```

```
## [1] "Probability of 3 doctor's visits: 0.000512242945776013"
```

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

```
doctor_fit2<- lm(doctorco ~ sex + age + agesq + income + levyplus + freepoor + freerepa + illness + actdays + hscore + chcond1 + chcond2, data = dvisits)
summary(doctor_fit2)
```

```
##
## Call:
## lm(formula = doctorco ~ sex + age + agesq + income + levyplus +
##     freepoor + freerepa + illness + actdays + hscore + chcond1 +
##     chcond2, data = dvisits)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.1352 -0.2588 -0.1435 -0.0433  7.0327
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.027632   0.072220   0.383  0.70202
## sex          0.033811   0.021604   1.565  0.11764
## age          0.203201   0.410016   0.496  0.62020
## agesq       -0.062103   0.458716  -0.135  0.89231
## income      -0.057323   0.033089  -1.732  0.08326 .
## levyplus     0.035179   0.024882   1.414  0.15748
## freepoor    -0.103314   0.052471  -1.969  0.04901 *
## freerepa     0.033241   0.038157   0.871  0.38371
## illness      0.059946   0.008357   7.173 8.39e-13 ***
## actdays     0.103192   0.003657  28.216 < 2e-16 ***
## hscore       0.016976   0.005190   3.271  0.00108 **
## chcond1      0.004384   0.023740   0.185  0.85349
## chcond2      0.041617   0.035863   1.160  0.24592
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7139 on 5177 degrees of freedom
## Multiple R-squared:  0.2018, Adjusted R-squared:  0.2
## F-statistic: 109.1 on 12 and 5177 DF,  p-value: < 2.2e-16
```

```
predict(doctor_fit2, dvisits[5190,])
```

```
##      5190
## 0.1606531
```

```
"It appears that it isn't likely to be too different"
```

```
## [1] "It appears that it isn't likely to be too different"
```