

Problem Set 5

QTM 200: Applied Regression Analysis

Vanessa Wong

Instructions

- Please show your work! You may lose points by simply writing in the answer. If the problem requires you to execute commands in R, please include the code you used to get your answers. Please also include the .R file that contains your code. If you are not sure if work needs to be shown for a particular problem, please ask.
- Your homework should be submitted electronically on the course GitHub page in .pdf form.
- This problem set is due at the beginning of class on Wednesday, March 4, 2020. No late assignments will be accepted.
- Total available points for this homework is 100.

Using the `teengamb` dataset, fit a model with `gamble` as the response and the other variables as predictors.

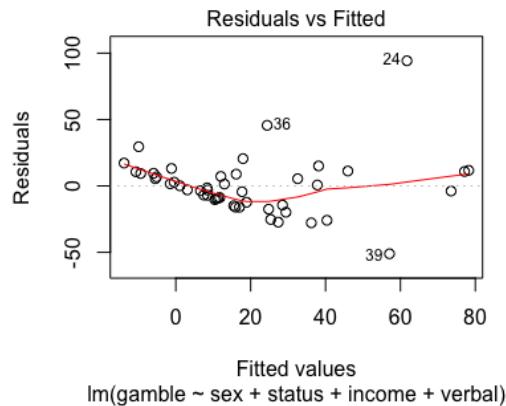
```
1 #####  
2  
3 # load data
```

Answer the following questions:

- (a) Check the constant variance assumption for the errors by plotting the residuals versus the fitted values.

```
1 # a. Residuals vs. Fitted Values  
2 par(mfrow=c(2,2)); plot(model1)  
3 plot(model1, which=1)
```

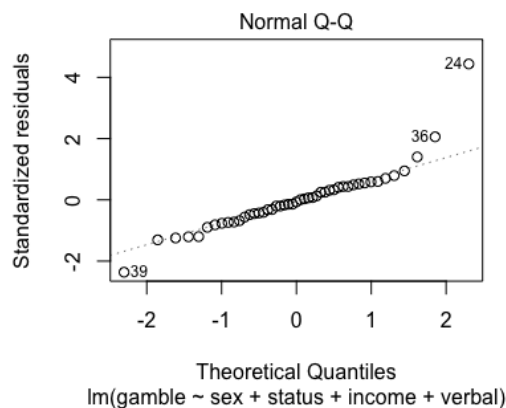
- After fitted value (\hat{y}) equals 45, variance seems to increase/residuals deviate more from 0. variance does not appear to be constant.



(b) Check the normality assumption with a Q-Q plot of the studentized residuals.

```
1 # b. Q-Q plot of studentized residuals
2 plot(modell, which=2)
```

- The normality assumption appears to be met because almost all points are very tightly clustered along the QQ line.

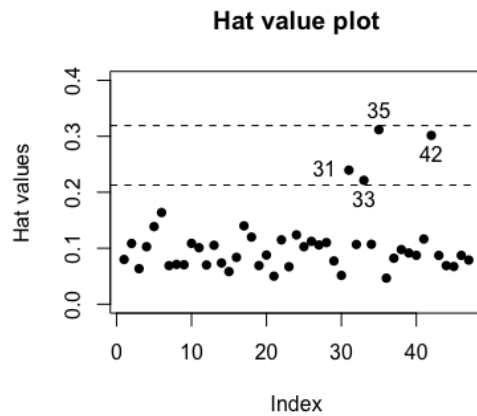


(c) Check for large leverage points by plotting the h values.

```
1 # c. Plot hat values
2 plot(hatvalues(modell), pch=16, cex=1, ylab="Hat values", main="Hat value
   plot", ylim=c(0,0.4))
3 abline(h=(2*5)/47, lty=2)
4 abline(h=(3*5)/47, lty=2)
5 identify(1:47, hatvalues(modell), row.names(gamble))
```

- Using the average hat value as a threshold (thresholds = $2 * \bar{h}$ and $3 * \bar{h}$), there are 4 high leverage points in this dataset: observations 31, 33, 35, and 42. These four points

have are potentially influential i.e. have the potential to greatly affect the regression model.



- (d) Check for outliers by running an `outlierTest`.

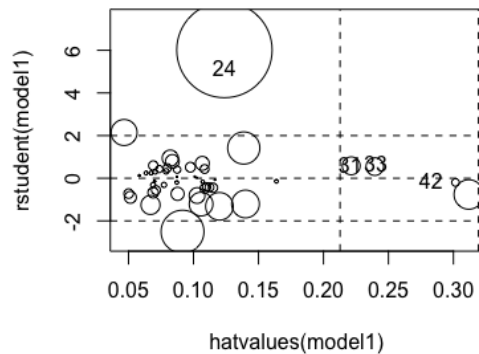
```
1 outlierTest(model1)
```

- According to outlier test, observation 24 has an adjusted (Bonferroni) p-value 0. therefore, observation 24 is an extreme residual in this model.

- (e) Check for influential points by creating a "Bubble plot" with the hat-values and studentized residuals.

```
1 plot(hatvalues(model1), rstudent(model1), ylim=c(-3,7.5), type="n")
2 cook <- sqrt(cooks.distance(model1))
3 points(hatvalues(model1), rstudent(model1), cex=10*cook/max(cook))
4 abline(h=c(-2,0,2), lty=2)
5 abline(v=c(2,3)*5/47, lty=2)
6 identify(hatvalues(model1), rstudent(model1), row.names(gamble))
```

- Observations 31, 33, 35, and 42 are influential points because all 4 are above at least one threshold for studentized residuals and hat values. This means all 4 points are both high leverage and are regression outliers, making them influential points (they greatly affect the regression model). Observation 24 was also identified because while it did not pass any of the hat value thresholds, it has an unusually large Cook's distance.



```

1 #####
2 # load libraries
3 # set wd
4 # clear global .envir
5 #####
6
7 # remove objects
8 rm(list=ls())
9 # detach all libraries
10 detachAllPackages <- function() {
11   basic.packages <- c("package:stats", "package:graphics", "package:grDevices",
12     "package:utils", "package:datasets", "package:methods", "package:base")
13   package.list <- search()[ifelse(unlist(gregexpr("package:", search()))==1,
14     TRUE, FALSE)]
15   package.list <- setdiff(package.list, basic.packages)
16   if (length(package.list)>0) for (package in package.list) detach(package,
17     character.only=TRUE)
18 }
19 detachAllPackages()
20
21 # load libraries
22 pkgTest <- function(pkg){
23   new.pkg <- pkg[!(pkg %in% installed.packages()[, "Package"])]
24   if (length(new.pkg))
25     install.packages(new.pkg, dependencies = TRUE)
26   sapply(pkg, require, character.only = TRUE)
27 }
28
29 # here is where you load any necessary packages
30 install.packages(car)
31 library(car)
32 library(boot)
33 # ex: stringr
34 # lapply(c("stringr"), pkgTest)
35
36 lapply(c("faraway"), pkgTest)

```

```

34
35 # set working directory
36 setwd("~/GitHub/QTM200Spring2020/problem_sets/PS5")
37
38
39 #####
40 # Problem 1
41 #####
42
43 # load data
44 gamble <- (data=teengamb)
45 # run regression on gamble with specified predictors
46 modell <- lm(gamble ~ sex + status + income + verbal, gamble)
47 summary(modell)
48
49 # a. Residuals vs. Fitted Values
50 par(mfrow=c(2,2)); plot(modell)
51 plot(modell, which=1)
52 # after fitted value of y hat equals ~ 45, variance seems to increase/
    residuals deviate more from 0. variance does not appear to be constant.
53
54 # b. Q-Q plot of studentized residuals
55 plot(modell, which=2)
56 # the normality assumption appears to be met because almost all points are
    very tightly clustered along the QQ line.
57
58 # c. Plot hat values
59 plot(hatvalues(modell), pch=16, cex=1, ylab="Hat values", main="Hat value plot
    ", ylim=c(0,0.4))
60 abline(h=(2*5)/47, lty=2)
61 abline(h=(3*5)/47, lty=2)
62 identify(1:47, hatvalues(modell), row.names(gamble))
63 # using the average hat value as a threshold (thresholds = 2 * hbar and 3 *
    hbar), there are 4 high leverage points in this dataset: observations 31,
    33, 35, and 42.
64 # these four points have are potentially influential i.e. have the potential
    to greatly affect the regression model.
65
66 # d. Outlier test
67 outlierTest(modell)
68 # according to outlier test, observation 24 has an adjusted (Bonferroni) p-
    value ~ 0. therefore, observation 24 is an extreme residual in this model.
69
70 # e. Bubble plot
71 plot(hatvalues(modell), rstudent(modell), ylim=c(-3,7.5), type="n")
72 cook <- sqrt(cooks.distance(modell))
73 points(hatvalues(modell), rstudent(modell), cex=10*cook/max(cook))
74 abline(h=c(-2,0,2), lty=2)
75 abline(v=c(2,3)*5/47, lty=2)
76 identify(hatvalues(modell), rstudent(modell), row.names(gamble))
77 # observations 31, 33, 35, and 42 are influential points because all 4 are

```

above at least one threshold for studentized residuals and hat values. This means all 4 points are both high leverage and are regression outliers, making them influential points (they greatly affect the regression model). Observation 24 was also identified because while it did not pass any of the hat value thresholds, it has an unusually large Cook's distance.