

Problem Set 5

QTM 200: Applied Regression Analysis

Due: March 4, 2020

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.

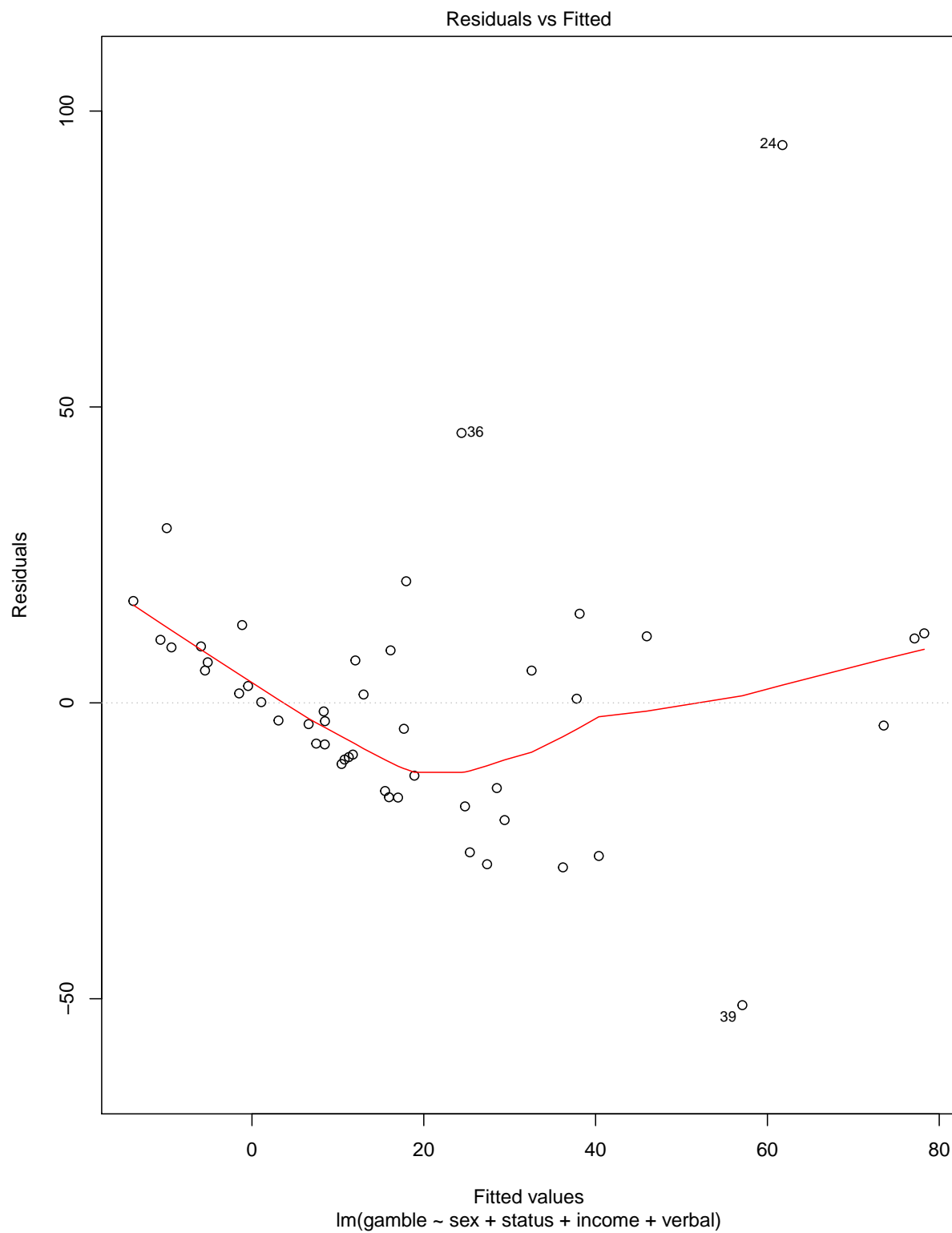
Using the **teengamb** dataset, fit a model with **gamble** as the response and the other variables as predictors.

```
1 # load data
2 gamble <- (data=teengamb)
3 # run regression on gamble with specified predictors
```

Answer the following questions:

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

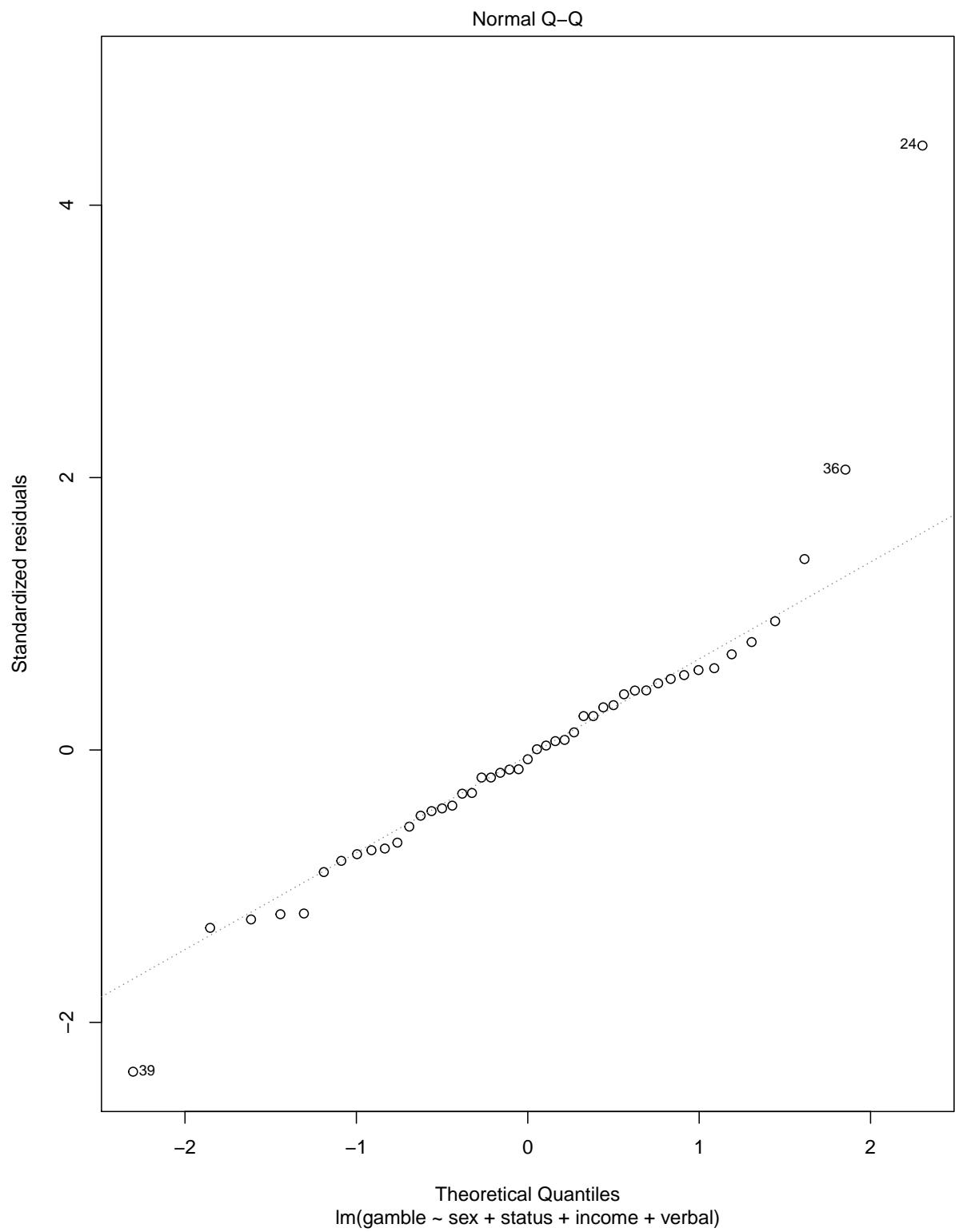
```
1 plot(model1)
```



Interpretation: After plotting the residuals versus the fitted values, we can see that the graph is not very linear and that the constant variance assumption may not necessarily hold for this example.

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

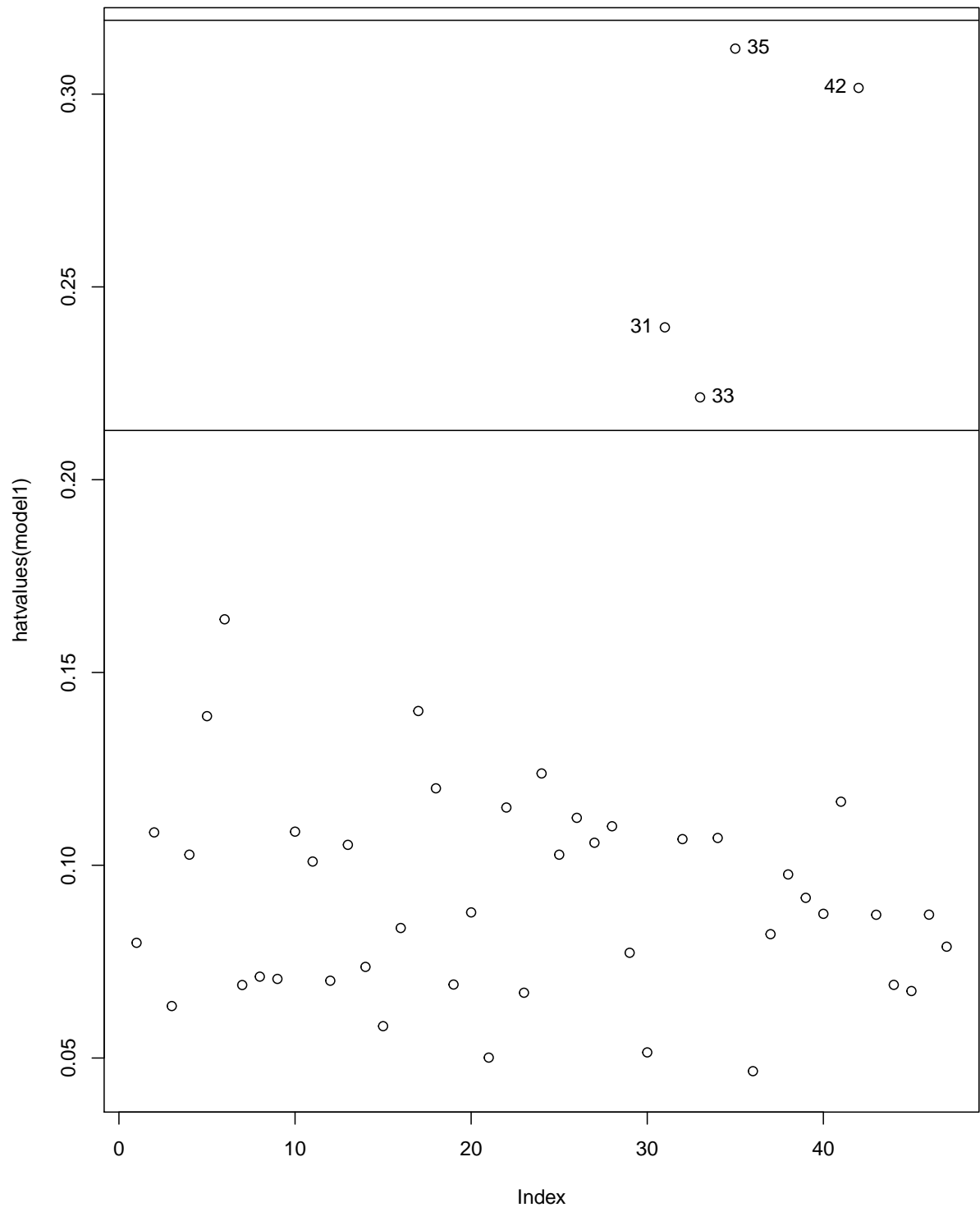
```
1 plot(model1)
```



Interpretation: After plotting a Q-Q plot of the studentized residuals, we can see that normality assumption is met because the vast majority of the observations fall on the line. Those observations that don't fall on the line happen to be at the tail ends where a less amount of data exists.

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

```
1 plot(hatvalues(model1))
2 abline(h=2*5/47)
3 abline(h=3*5/47)
4 identify(1:47, hatvalues(model1), row.names(gamble))
5 gamble
```



| | sex | status | income | verbal | gamble |
|----|-----|--------|--------|--------|--------|
| 31 | 0 | 18 | 12.00 | 2 | 88.00 |
| 33 | 0 | 38 | 15.00 | 7 | 90.00 |
| 35 | 0 | 28 | 1.50 | 1 | 14.10 |
| 42 | 0 | 61 | 15.00 | 9 | 69.70 |

Interpretation: Our analysis shows four large leverage points. These points are outlined on the graph above and display the characteristics tabulated above.

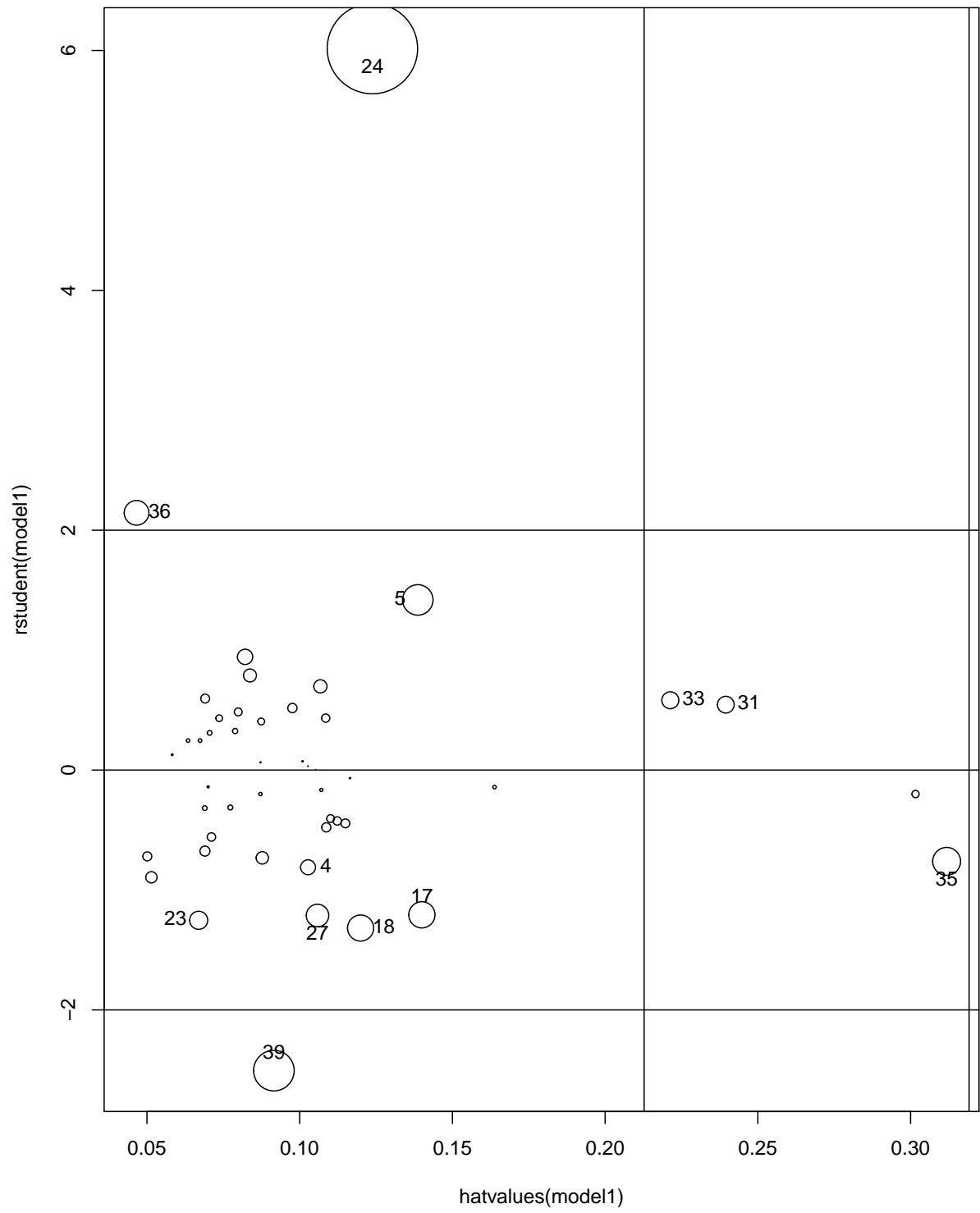
- (d) Check for outliers by running an `outlierTest`.

```
1 library(car)
2 outlierTest(model1)
```

| | <code>rstudent</code> | <code>unadjusted</code> | <code>p-value</code> | <code>Bonferroni</code> | <code>p</code> |
|----|-----------------------|-------------------------|----------------------|-------------------------|----------------|
| 24 | 6.016116 | | 4.1041e-07 | 1.9289e-05 | |

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

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



Interpretation: There are several influential points that were discovered after creating a "Bubble plot". They are outlined on the graph above.