

# MA678 homework 01

Ningze Zu

Septemeber 14, 2018

## Introduction

For homework 1 you will fit linear regression models and interpret them. You are welcome to transform the variables as needed. How to use `lm` should have been covered in your discussion session. Some of the code are written for you. Please remove `eval=FALSE` inside the knitr chunk options for the code to run.

This is not intended to be easy so please come see us to get help.

## Data analysis

### Pyth!

```
gelman_example_dir<-"http://www.stat.columbia.edu/~gelman/arm/examples/"
pyth <- read.table(paste0(gelman_example_dir,"pyth/exercise2.1.dat"),
  header=T, sep=" ")
```

The folder `pyth` contains outcome `y` and inputs `x1`, `x2` for 40 data points, with a further 20 points with the inputs but no observed outcome. Save the file to your working directory and read it into R using the `read.table()` function.

1. Use R to fit a linear regression model predicting `y` from `x1`, `x2`, using the first 40 data points in the file. Summarize the inferences and check the fit of your model.

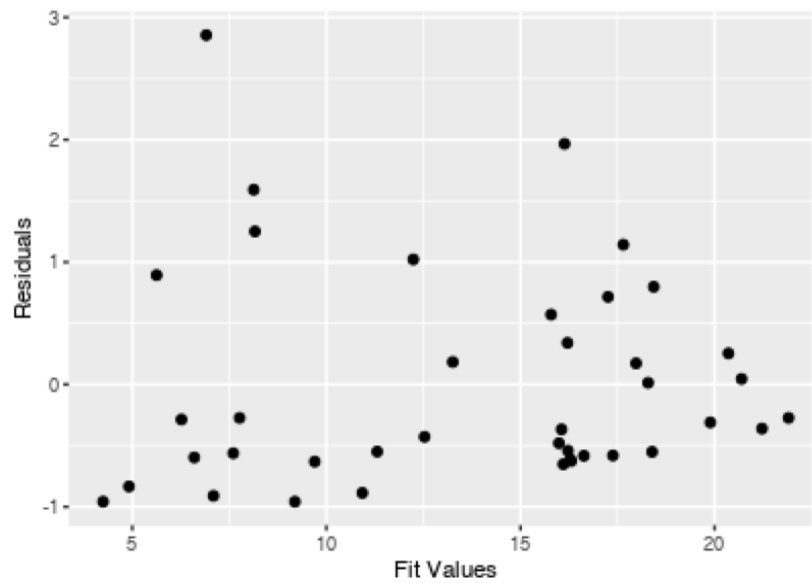
```
regout1 <- lm(y ~ x1 + x2, data =pyth )
summary(regout1)

##
## Call:
## lm(formula = y ~ x1 + x2, data = pyth)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.9585 -0.5865 -0.3356  0.3973  2.8548
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.31513    0.38769   3.392  0.00166 **
## x1           0.51481    0.04590  11.216 1.84e-13 ***
## x2           0.80692    0.02434  33.148 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.9 on 37 degrees of freedom
## (20 observations deleted due to missingness)
## Multiple R-squared:  0.9724, Adjusted R-squared:  0.9709
## F-statistic: 652.4 on 2 and 37 DF, p-value: < 2.2e-16
```

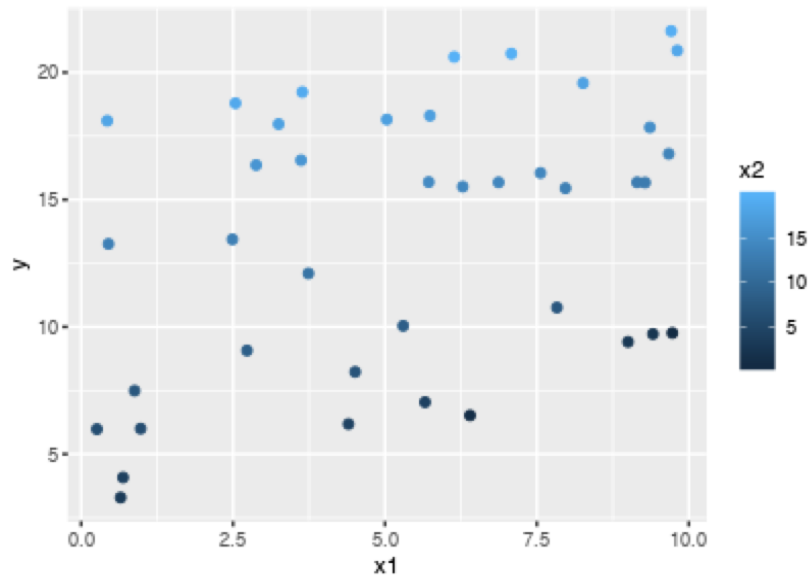
2. Display the estimated model graphically as in (GH) Figure 3.2.

```
fittedValues <- fitted(regout1)
Residuals    <- resid(regout1)

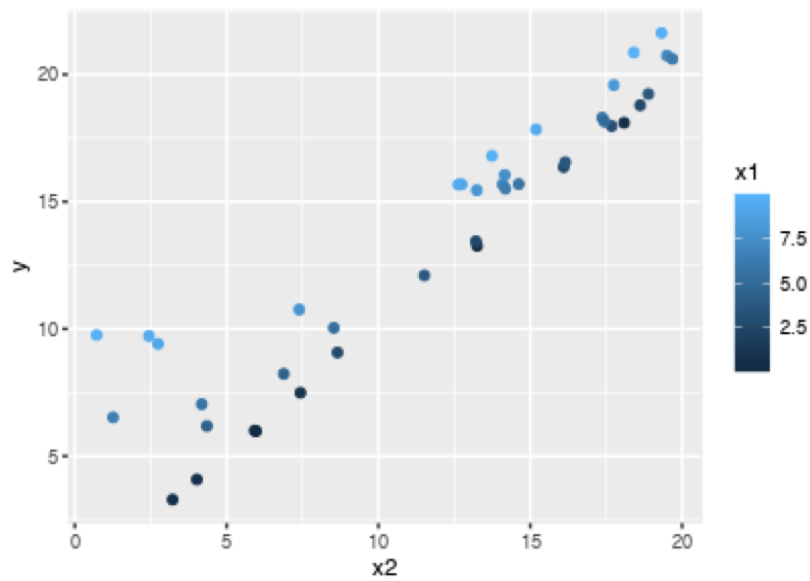
#ggplot2
library(ggplot2)
qplot(fittedValues, Residuals) + xlab("Fit Values") + ylab("Residuals")
```



```
#x1
qplot(x1,y,color=x2,data=regout1)
```

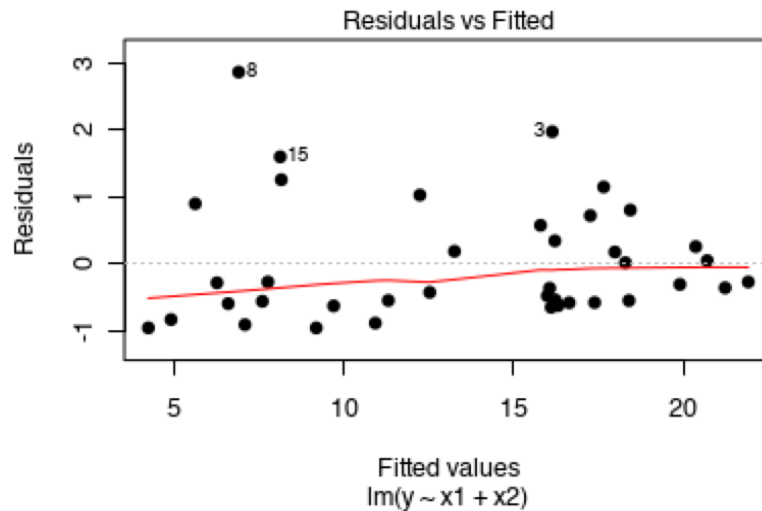


```
#x2
qplot(x2,y,color=x1,data=regout1)
```



3. Make a residual plot for this model. Do the assumptions appear to be met?

```
plot(regout1, pch=16, which=1)
```



```
"The Assumptions do not appear to be met"
```

```
## [1] "The Assumptions do not appear to be met"
```

4. Make predictions for the remaining 20 data points in the file. How confident do you feel about these predictions?

```
predictions <- predict(regout1, newdata = pyth[41:60,])
predictions
```

```
##      41      42      43      44      45      46      47
## 14.812484 19.142865  5.916816 10.530475 19.012485 13.398863  4.829144
##      48      49      50      51      52      53      54
##  9.145767  5.892489 12.338639 18.908561 16.064649  8.963122 14.972786
##      55      56      57      58      59      60
##  5.859744  7.374900  4.535267 15.133280  9.100899 16.084900
```

```
"I do not feel confident about these predictions since there are many residuals problem above."
```

```
## [1] "I do not feel confident about these predictions since there are many residuals problem above."
```

After doing this exercise, take a look at Gelman and Nolan (2002, section 9.4) to see where these data came from. (or ask Masanao)

### Earning and height

Suppose that, for a certain population, we can predict log earnings from log height as follows:

- A person who is 66 inches tall is predicted to have earnings of \$30,000.
- Every increase of 1% in height corresponds to a predicted increase of 0.8% in earnings.

- The earnings of approximately 95% of people fall within a factor of 1.1 of predicted values.
1. Give the equation of the regression line and the residual standard deviation of the regression.

Suppose the height is  $\hat{x}_i$  and the predictor earnings is  $\hat{y}_i$ .

$$\log(\hat{y}_i) = \beta_0 + \frac{0.008}{0.01} \log(\hat{x}_i)$$

$$\log(30000) = \beta_0 + \frac{0.008}{0.01} \log(66)$$

$$\beta_0 = 6.957229$$

Then the equation of the regression line is given by

$$\log(\hat{y}_i) = 6.957229 + \frac{0.008}{0.01} \log(\hat{x}_i)$$

```
sd = 0.1 * .50 / .95
sd
```

```
## [1] 0.05263158
```

2. Suppose the standard deviation of log heights is 5% in this population. What, then, is the  $R^2$  of the regression model described here?

```
sd.population = 0.05
R2 <- 1 - (sd^2 / sd.population^2)
R2
```

```
## [1] -0.1080332
```

### Beauty and student evaluation

The folder beauty contains data from Hamermesh and Parker (2005) on student evaluations of instructors' beauty and teaching quality for several courses at the University of Texas. The teaching evaluations were conducted at the end of the semester, and the beauty judgments were made later, by six students who had not attended the classes and were not aware of the course evaluations.

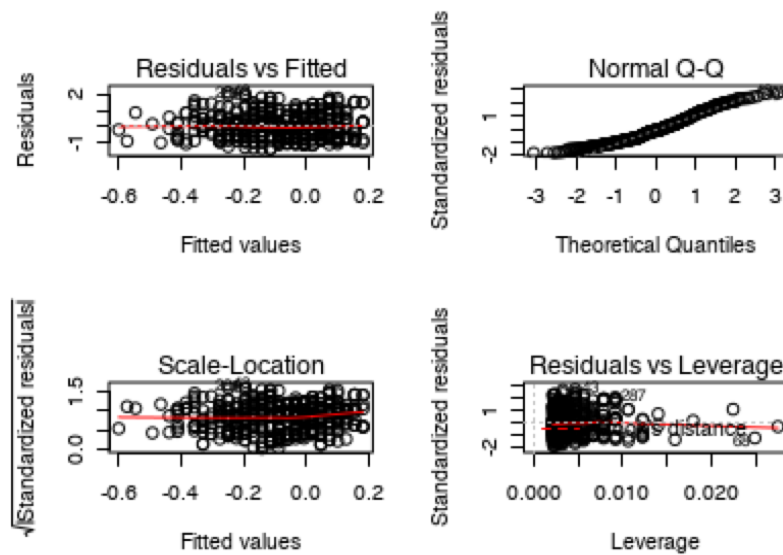
```
beauty.data <- read.table(paste0(gelman_example_dir, "beauty/ProfEvaltnsBeautyPublic.csv"), header=T, s
```

1. Run a regression using beauty (the variable btystdave) to predict course evaluations (courseevaluation), controlling for various other inputs. Display the fitted model graphically, and explaining the meaning of each of the coefficients, along with the residual standard deviation. Plot the residuals versus fitted values.

```
regout2 <- lm(btystdave ~ courseevaluation, data=beauty.data)
display(regout2)
```

```
## lm(formula = btystdave ~ courseevaluation, data = beauty.data)
##               coef.est coef.se
## (Intercept)    -1.16    0.26
## courseevaluation  0.27    0.07
## ---
## n = 463, k = 2
## residual sd = 0.78, R-Squared = 0.04
```

```
par(mfrow=c(2,2))
plot(regout2)
```

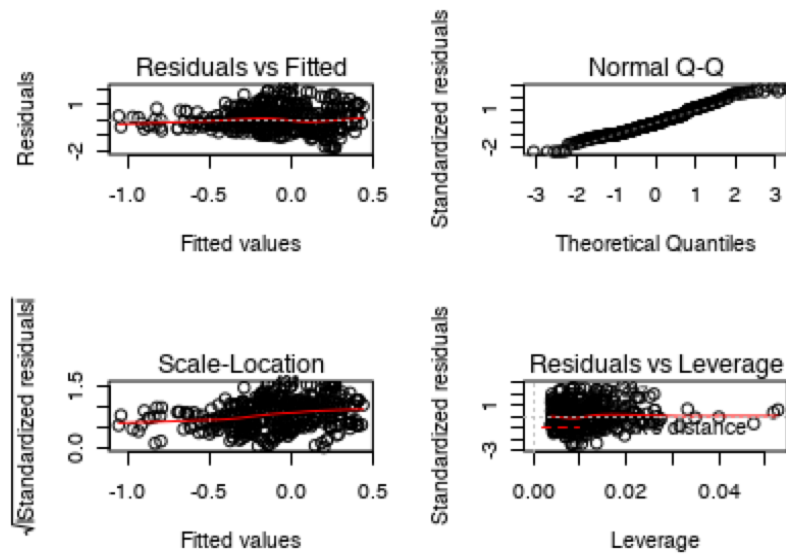


*#The coefficient of courseevaluation represents the slope of this variable.*

2. Fit some other models, including beauty and also other input variables. Consider at least one model with interactions. For each model, state what the predictors are, and what the inputs are, and explain the meaning of each of its coefficients.

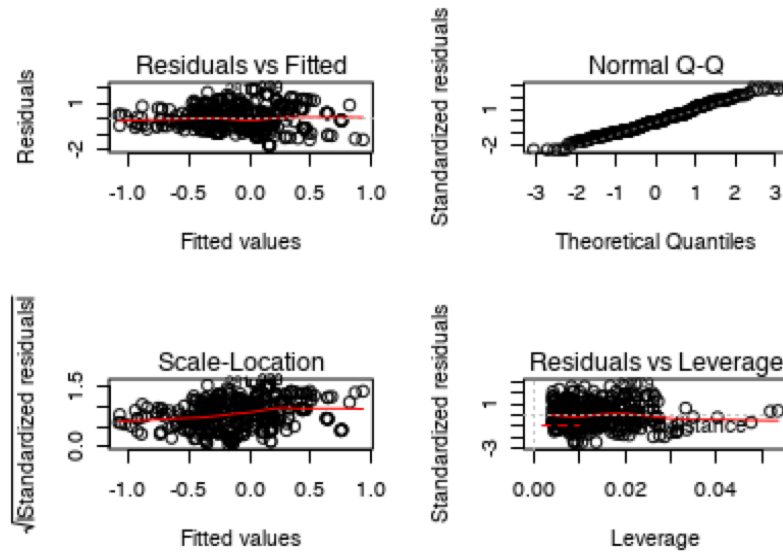
```
#
regout3 <- lm(btystdave ~ courseevaluation * female + age, data=beauty.data)
display(regout3)

## lm(formula = btystdave ~ courseevaluation * female + age, data = beauty.data)
##               coef.est coef.se
## (Intercept)    -0.51    0.39
## courseevaluation  0.36    0.08
## female          1.02    0.51
## age            -0.02    0.00
## courseevaluation:female -0.23    0.13
## ---
## n = 463, k = 5
## residual sd = 0.74, R-Squared = 0.13
par(mfrow=c(2,2))
plot(regout3)
```



```
#
regout4 <- lm(btystdave ~ courseevaluation * female + age + blkandwhite, data=beauty.data)
display(regout4)

## lm(formula = btystdave ~ courseevaluation * female + age + blkandwhite,
##    data = beauty.data)
##              coef.est coef.se
## (Intercept)    -0.34   0.37
## courseevaluation   0.32   0.08
## female           1.41   0.49
## age             -0.02   0.00
## blkandwhite       0.59   0.09
## courseevaluation:female -0.35   0.12
## ---
## n = 463, k = 6
## residual sd = 0.71, R-Squared = 0.20
par(mfrow=c(2,2))
plot(regout4)
```



In regout3 models, predictors are courseevaluation:female, courseevaluation, female, age.

The coefficient on the interaction term represents the difference in the slope, comparing females to males.

The coefficient of courseevaluation represents the slope of this variable when other variables are 0.

The coefficient of female represents the predicted difference for subjects that differ in sex when other variables are 0.

See also Felton, Mitchell, and Stinson (2003) for more on this topic link

## Conceptual exercises

On statistical significance.

Note: This is more like a demo to show you that you can get statistically significant result just by random chance. We haven't talked about the significance of the coefficient so we will follow Gelman and use the



approximate definition, which is if the estimate is more than 2 sd away from 0 or equivalently, if the z score is bigger than 2 as being “significant”.

( From Gelman 3.3 ) In this exercise you will simulate two variables that are statistically independent of each other to see what happens when we run a regression of one on the other.

1. First generate 1000 data points from a normal distribution with mean 0 and standard deviation 1 by typing in R. Generate another variable in the same way (call it var2).

```
var1 <- rnorm(1000,0,1)
var2 <- rnorm(1000,0,1)
m1 <- lm(var1~var2)
summary(m1)
```

Run a regression of one variable on the other. Is the slope coefficient statistically significant? [absolute value of the z-score(the estimated coefficient of var1 divided by its standard error) exceeds 2]

```
fit <- lm (var2 ~ var1)
z.scores <- coef(fit)[2]/se.coef(fit)[2]
z.scores
```

2. Now run a simulation repeating this process 100 times. This can be done using a loop. From each simulation, save the z-score (the estimated coefficient of var1 divided by its standard error). If the absolute value of the z-score exceeds 2, the estimate is statistically significant. Here is code to perform the simulation:

```
z.scores <- rep (NA, 100)
for (k in 1:100) {
  var1 <- rnorm (1000,0,1)
  var2 <- rnorm (1000,0,1)
  fit <- lm (var2 ~ var1)
  z.scores[k] <- coef(fit)[2]/se.coef(fit)[2]
}
sum(z.scores > 1.96)
```

How many of these 100 z-scores are statistically significant? What can you say about statistical significance of regression coefficient?

The regression coefficient significance at a 5% significance level.

### Fit regression removing the effect of other variables

Consider the general multiple-regression equation

$$Y = A + B_1X_1 + B_2X_2 + \cdots + B_kX_k + E$$

An alternative procedure for calculating the least-squares coefficient  $B_1$  is as follows:

1. Regress  $Y$  on  $X_2$  through  $X_k$ , obtaining residuals  $E_{Y|2,\dots,k}$ .
  2. Regress  $X_1$  on  $X_2$  through  $X_k$ , obtaining residuals  $E_{1|2,\dots,k}$ .
  3. Regress the residuals  $E_{Y|2,\dots,k}$  on the residuals  $E_{1|2,\dots,k}$ . The slope for this simple regression is the multiple-regression slope for  $X_1$  that is,  $B_1$ .
- (a) Apply this procedure to the multiple regression of prestige on education, income, and percentage of women in the Canadian occupational prestige data (<http://socserv.socsci.mcmaster.ca/jfox/Books/Applied-Regression-3E/datasets/Prestige.pdf>), confirming that the coefficient for education is properly recovered.

```

fox_data_dir<-"http://socserv.socsci.mcmaster.ca/jfox/Books/Applied-Regression-3E/datasets/"
Prestige<-read.table(paste0(fox_data_dir,"Prestige.txt"))
#1
lm1 <- lm(prestige~income+women,Prestige)
residuals1 <- resid(lm1)
meanY <- mean(residuals1)
#2
lm2 <- lm(education~income+women,Prestige)
residuals2 <- resid(lm2)
meanX1 <- mean(residuals2)
#3
lm3 <- lm(residuals1~residuals2)

```

- (b) The intercept for the simple regression in step 3 is 0. Why is this the case?
- (c) In light of this procedure, is it reasonable to describe  $B_1$  as the “effect of  $X_1$  on  $Y$  when the influence of  $X_2, \dots, X_k$  is removed from both  $X_1$  and  $Y$ ”?
- (d) The procedure in this problem reduces the multiple regression to a series of simple regressions ( in Step 3). Can you see any practical application for this procedure?

### Partial correlation

The partial correlation between  $X_1$  and  $Y$  “controlling for”  $X_2, \dots, X_k$  is defined as the simple correlation between the residuals  $E_{Y|2,\dots,k}$  and  $E_{1|2,\dots,k}$ , given in the previous exercise. The partial correlation is denoted  $r_{y1|2,\dots,k}$ .

1. Using the Canadian occupational prestige data, calculate the partial correlation between prestige and education, controlling for income and percentage women.

```

corr <- cor(residuals1,residuals2)
corr

```

```
## [1] 0.7362604
```

2. In light of the interpretation of a partial regression coefficient developed in the previous exercise, why is  $r_{y1|2,\dots,k} = 0$  if and only if  $B_1$  is 0?

### Mathematical exercises.

Prove that the least-squares fit in simple-regression analysis has the following properties:

1.  $\sum \hat{y}_i \hat{e}_i = 0$   
 $\sum \hat{y}_i \hat{e}_i = \sum \hat{e}_i(\beta_0 + \beta_1 X_i) = \sum \beta_0 \hat{e}_i + \sum \beta_1 \hat{e}_i X_i = \beta_0 \sum \hat{e}_i + \beta_1 \sum \hat{e}_i X_i = \beta_0 * 0 + \beta_1 * 0 = 0$
2.  $\sum (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = \sum \hat{e}_i(\hat{y}_i - \bar{y}) = 0$   
 $\sum (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = \sum \hat{e}_i(\hat{y}_i - \bar{y}) = \sum \hat{e}_i(\beta_0 + \beta_1 X_i - \bar{y}) = \sum (\beta_0 - \bar{y}) \hat{e}_i + \sum \beta_1 \hat{e}_i X_i = (\beta_0 - \bar{y}) \sum \hat{e}_i + \beta_1 \sum \hat{e}_i X_i = \sum \hat{e}_i(\hat{y}_i - \bar{y}) = 0$

Suppose that the means and standard deviations of  $\mathbf{y}$  and  $\mathbf{x}$  are the same:  $\bar{\mathbf{y}} = \bar{\mathbf{x}}$  and  $sd(\mathbf{y}) = sd(\mathbf{x})$ .

1. Show that, under these circumstances

$$\beta_{y|x} = \beta_{x|y} = r_{xy}$$

where  $\beta_{y|x}$  is the least-squares slope for the simple regression of  $y$  on  $x$ ,  $\beta_{x|y}$  is the least-squares slope for the simple regression of  $x$  on  $y$ , and  $r_{xy}$  is the correlation between the two variables. Show that the intercepts are also the same,  $\alpha_{y|x} = \alpha_{x|y}$ .

2. Why, if  $\alpha_{y|x} = \alpha_{x|y}$  and  $\beta_{y|x} = \beta_{x|y}$ , is the least squares line for the regression of  $y$  on  $x$  different from the line for the regression of  $x$  on  $y$  (when  $r_{xy} < 1$ )?
3. Imagine that educational researchers wish to assess the efficacy of a new program to improve the reading performance of children. To test the program, they recruit a group of children who are reading substantially below grade level; after a year in the program, the researchers observe that the children, on average, have improved their reading performance. Why is this a weak research design? How could it be improved?

### Feedback comments etc.

If you have any comments about the homework, or the class, please write your feedback here. We love to hear your opinions.