

Homework 02

yourname

Septemeber 16, 2018

Introduction

In homework 2 you will fit many regression models. You are welcome to explore beyond what the question is asking you.

Please come see us we are here to help.

Data analysis

Analysis of earnings and height data

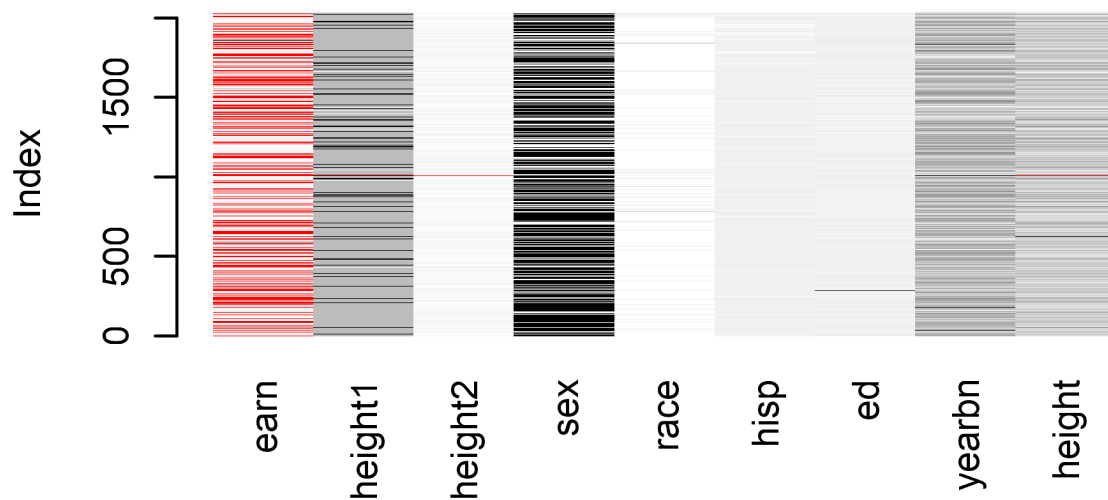
The folder **earnings** has data from the Work, Family, and Well-Being Survey (Ross, 1990). You can find the codebook at <http://www.stat.columbia.edu/~gelman/arm/examples/earnings/wfwcodebook.txt>

Pull out the data on earnings, sex, height, and weight.

1. In R, check the dataset and clean any unusually coded data.

Table 1: Quality of Data

	earn	height1	height2	sex	race	hisp	ed	yearbn	height
Number of Unique Value	135.00	4	14.0	2	5	3	19	74	24
Number of Missing Value	650.00	8	6.0	0	0	0	0	0	8
Number of 0	187.00	0	194.0	0	0	0	0	3	0
% of Missing Value	0.41	0	0.1	0	0	0	0	0	0

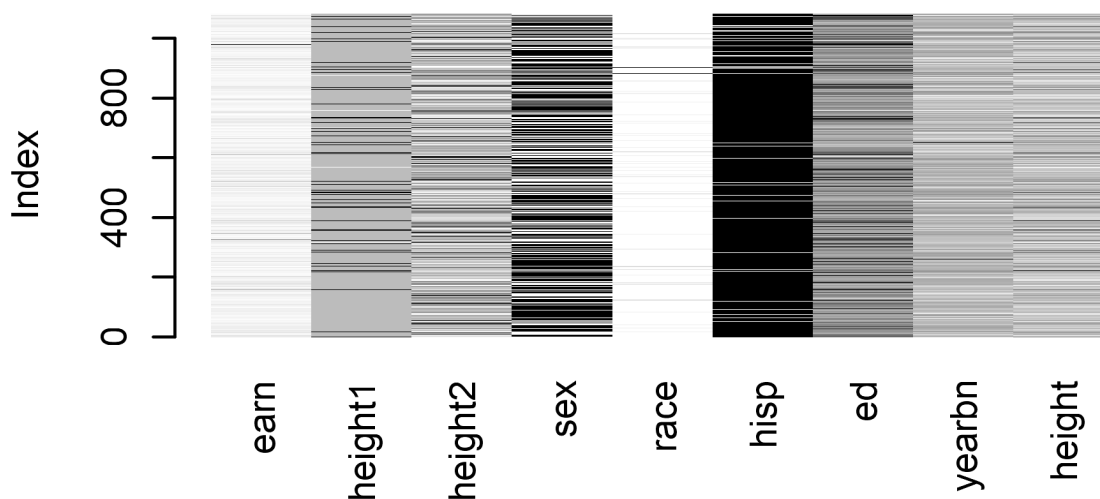


Above chart and figure show the distribution of missing value among variables, in the figure, the red means there is a missing value, and for the rest, the darker the color, the higher the value.

So according to the analysis above, the most of the missing value ,which is represented by ‘NA’ and ‘0’, contained in variable “earning”, about 41% of data doesn’t have variable ‘earning’, which is reasonable, beacause earning is kond of a private quetion, due to the extent of missing, we can hardly apply any imputaion, so simply delete it.

Table 2: Quality of Data

	earn	height1	height2	sex	race	hisp	ed	yearbn	height
Number of Unique Value	130	3	11	2	5	2	16	71	18
Number of Missing Value	0	0	0	0	0	0	0	0	0
Number of 0	0	0	0	0	0	0	0	0	0
% of Missing Value	0	0	0	0	0	0	0	0	0

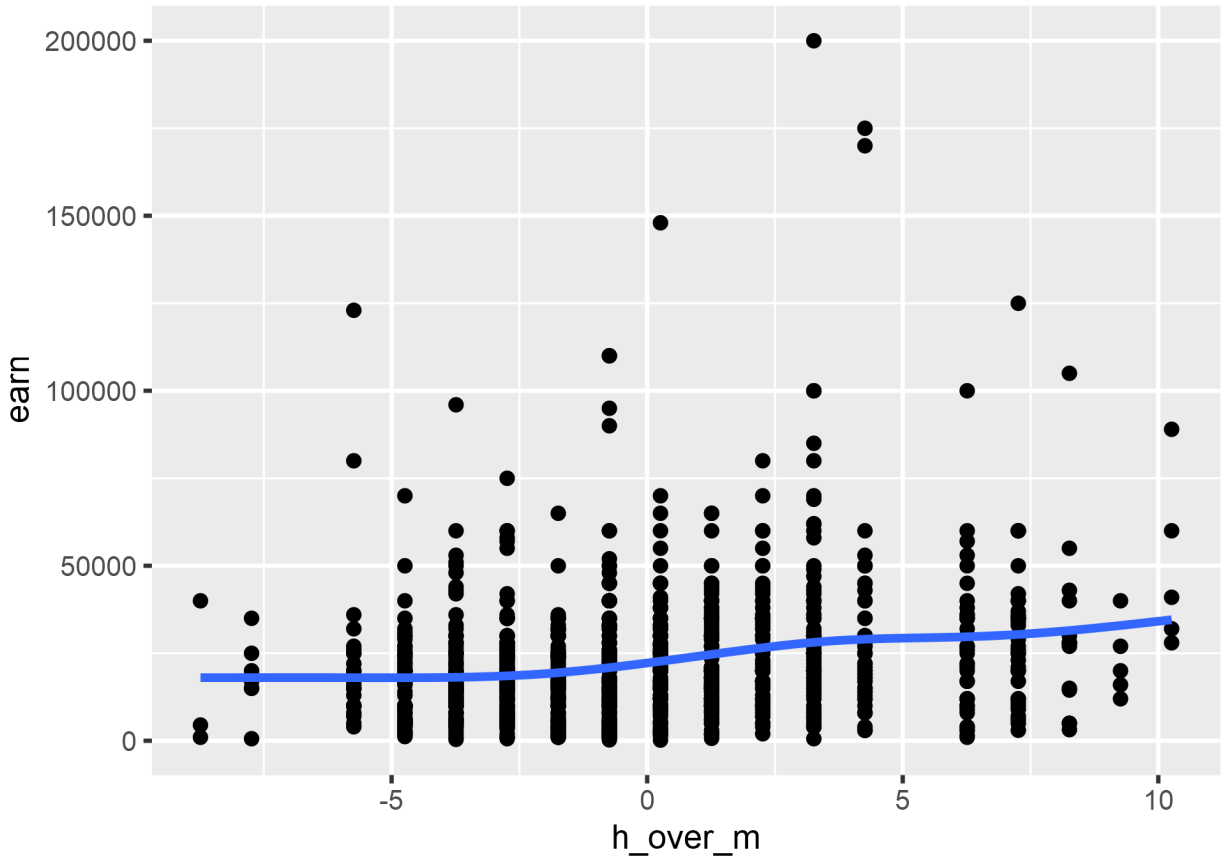


So now the data is clean and ready to be analysisd.

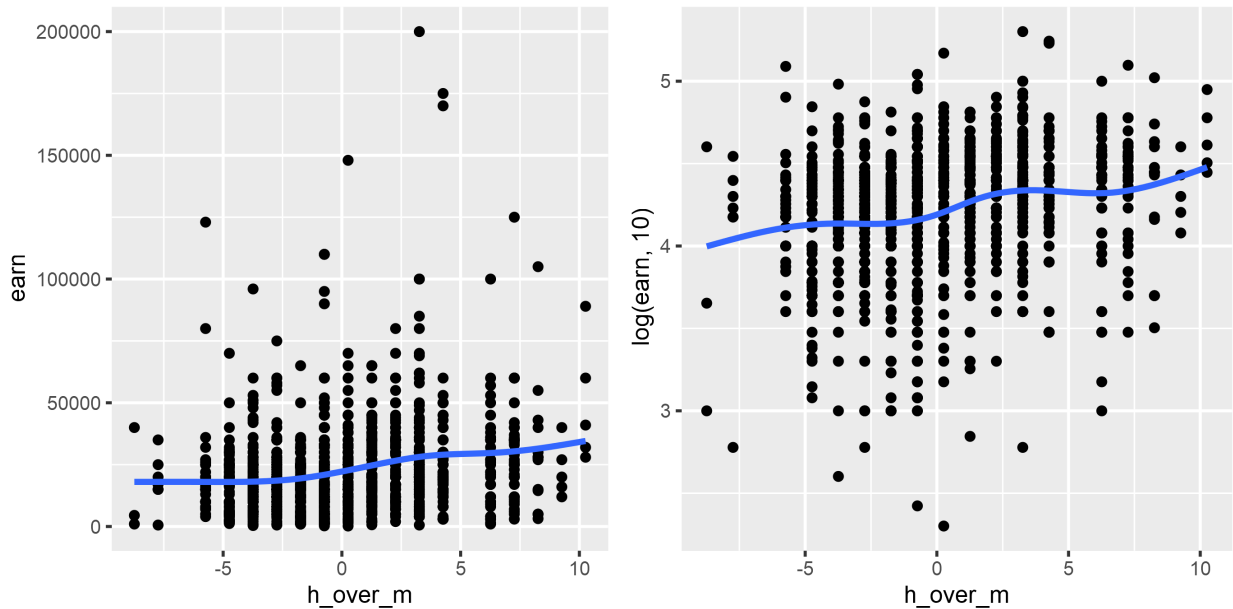
2.Fit a linear regression model predicting earnings from height. What transformation should you perform in order to interpret the intercept from this model

as average earnings for people with average height?

Look at the data, we can find three variables that related to height, “height 1”, “height 2” and “height”, after reading the original article, we know that real height should be “hight 1” feet and “height 2” inches, so $RealHeight = 12 * height_1 + height_2$, which is exactly the variable “height”.



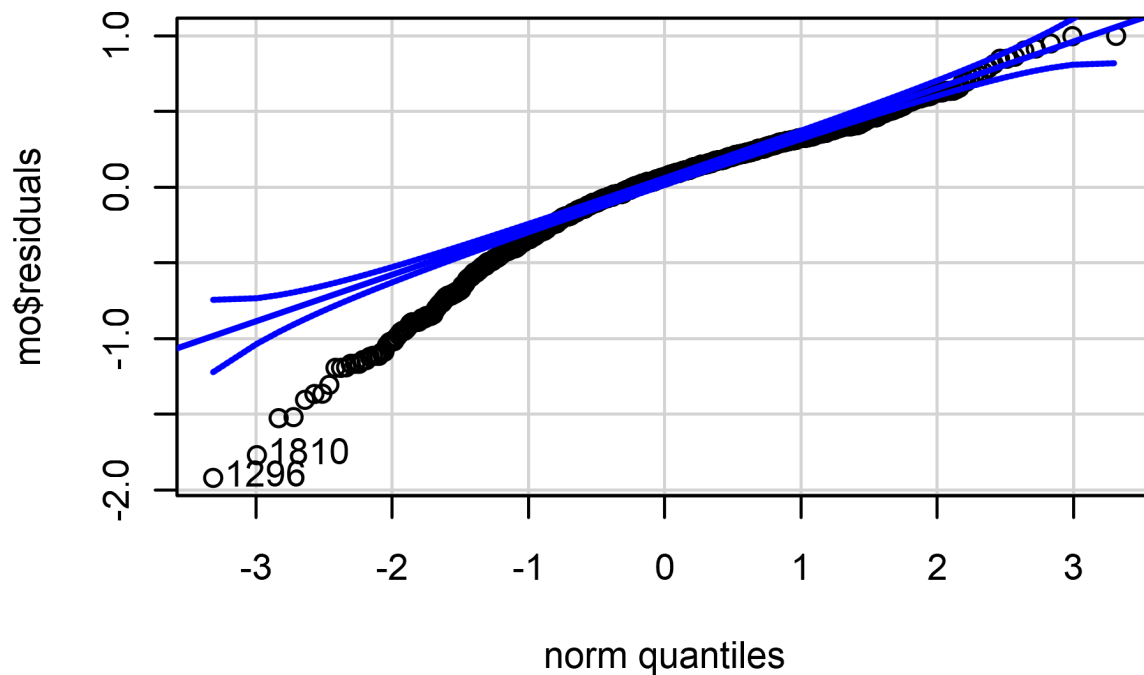
Look at the plot of height over average against earning, we can easily find some outliers, and a linear relationship may seem weird, so we apply some non-linear transformation, first we consider using $\log(\text{earn})$ to replace 'earn' to eliminate the effect of outliers. As for the predictor, we try to include the polynomial in our model.



From the above plot, we can tell that replacing 'earn' by $\log_{10} \text{Earn}$ does condense the data and alleviate

the effect of outliers a bit. Than apply regression analysis.

```
##
## Call:
## lm(formula = log(earn, 10) ~ h_over_m + I(h_over_m^2) + I(h_over_m^3),
##     data = hi)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.9194 -0.1696  0.0658  0.2459  0.9998
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   4.2130970   0.0153305  274.818 < 2e-16 ***
## h_over_m       0.0281404   0.0057486   4.895 1.13e-06 ***
## I(h_over_m^2)  0.0003175   0.0008552   0.371  0.710
## I(h_over_m^3) -0.0001484   0.0001499  -0.990  0.322
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.386 on 1077 degrees of freedom
## Multiple R-squared:  0.04665,    Adjusted R-squared:  0.044
## F-statistic: 17.57 on 3 and 1077 DF,  p-value: 3.872e-11
```

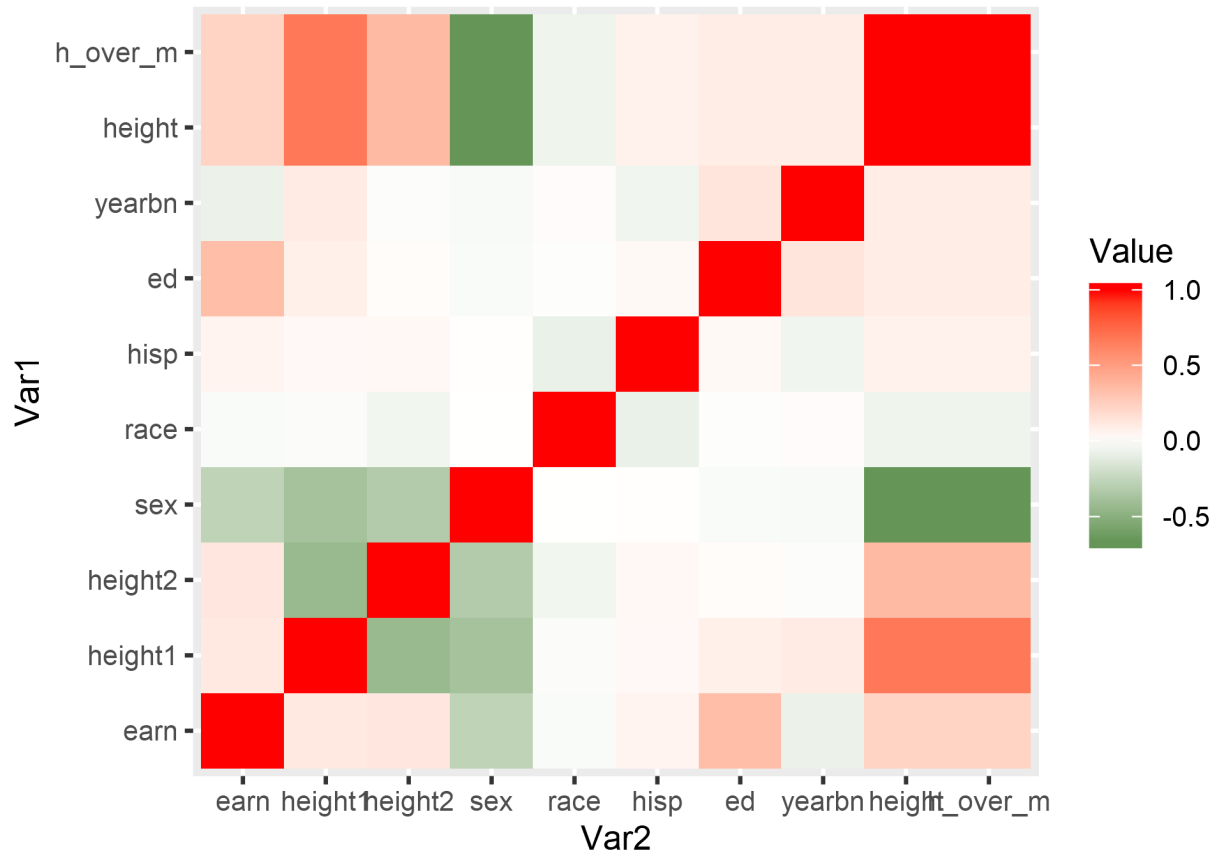


```
## 1296 1810
## 699 955
```

From result above, we can see that the performance of the model is poor, whose residuals are not normally distributed and even not close. I tried different combination of predictors, including heights and the power of height. The coefficient of model shows that only the coefficient of height is significant, Which means one unit increase(1 inch) will increase the earn by around 2.8%.

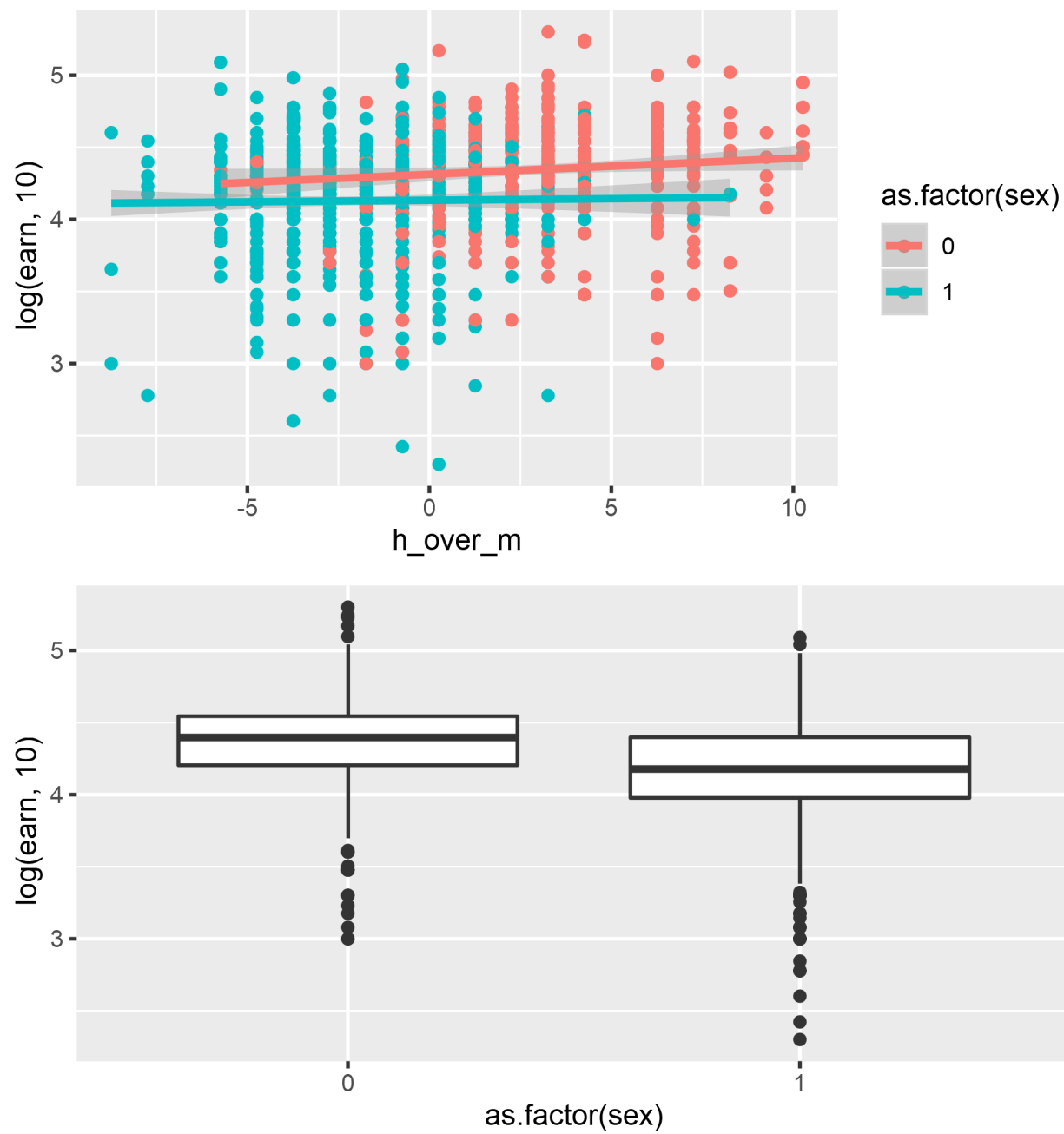
3. Fit some regression models with the goal of predicting earnings from some combination of sex, height, and weight. Be sure to try various transformations and interactions that might make sense. Choose your preferred model and justify.

To select proper predictors, we first do some analysis.

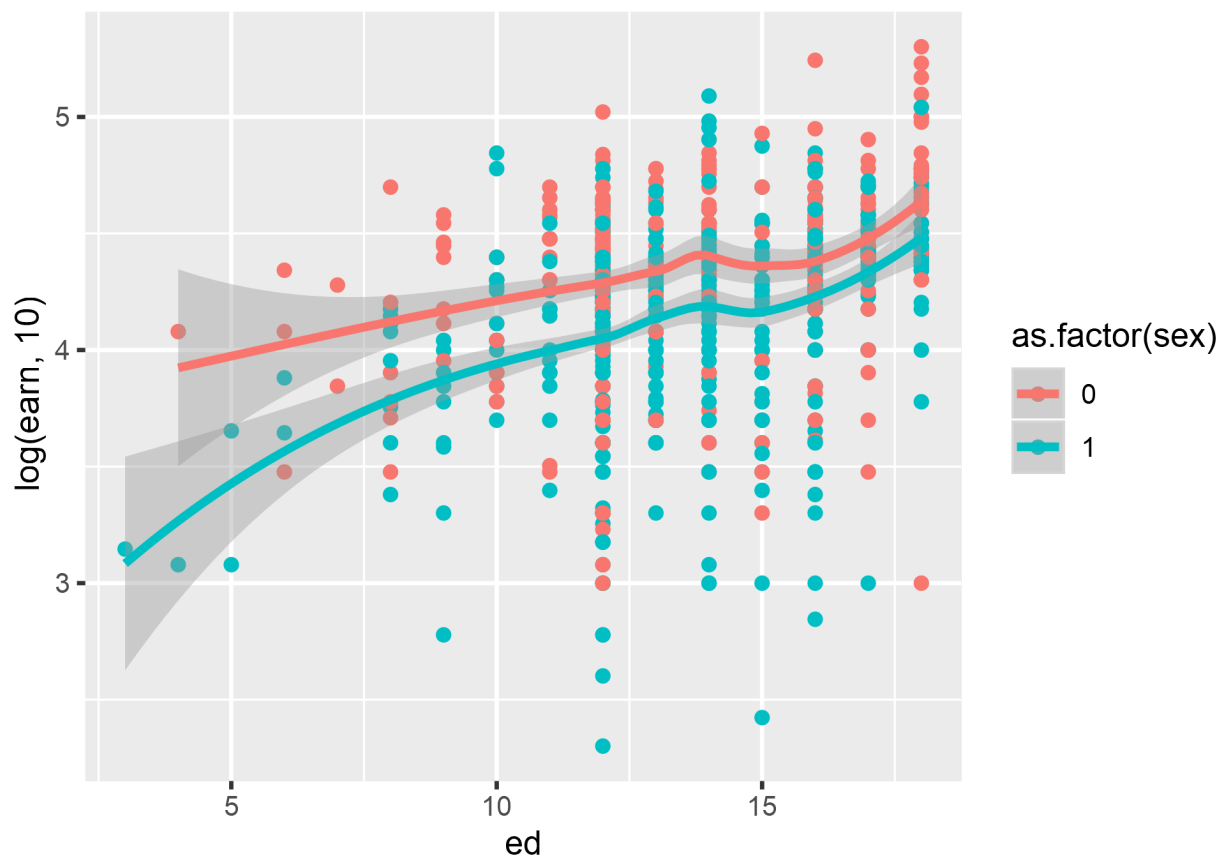


We can tell that sex is highly related to height, earn is related to earn, race is not obviously related to any other variables. So we may consider put sex, ed into our model.

Before we analyze sex, we have to transform sex into a binary variable, where 0 represent male, 1 represent female.

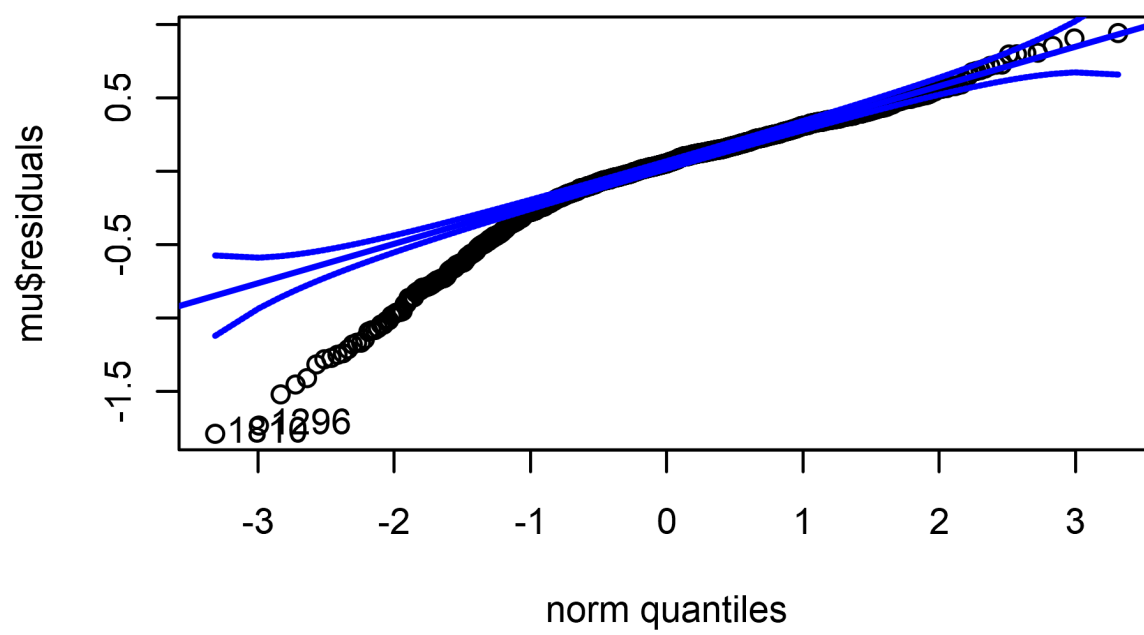
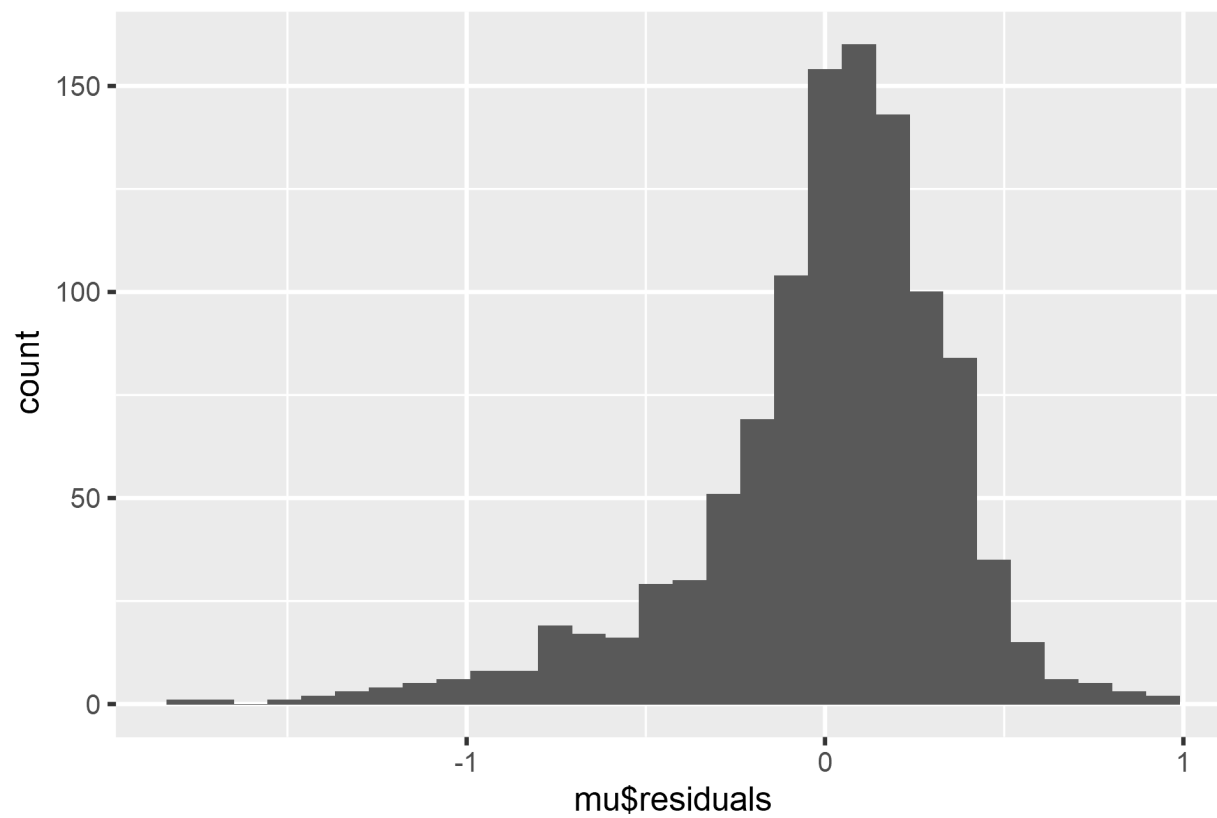


For sex, we can reasonably guess that sex does affect the average of earning and the way height affecting earn.



We can tell the education level is related to earn, rough judge from above figure, we can tell that higher education level, higher earning, and sex has no obvious effect on the slope of education level.

```
##
## Call:
## lm(formula = log(earn, 10) ~ h_over_m + sex + I(log(ed)) + sex:h_over_m +
##     ed:sex, data = hi)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.7910 -0.1399  0.0520  0.2228  0.9413
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.037750   0.213508  14.228 < 2e-16 ***
## h_over_m       0.008665   0.005729   1.513  0.130690
## sex           -0.480932   0.125001  -3.847  0.000126 ***
## I(log(ed))     0.494849   0.082338   6.010  2.54e-09 ***
## h_over_m:sex  -0.013955   0.008169  -1.708  0.087889 .
## sex:ed         0.020896   0.008966   2.330  0.019964 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3586 on 1075 degrees of freedom
## Multiple R-squared:  0.1786, Adjusted R-squared:  0.1748
## F-statistic: 46.75 on 5 and 1075 DF,  p-value: < 2.2e-16
```




```
## 1810 1296
## 955 699
```

After adding new variables to our model, R^2 increased a little, but the residuals of model is still not normally distributed.

4. Interpret all model coefficients.

5. Construct 95% confidence interval for all model coefficients and discuss what they mean.

Analysis of mortality rates and various environmental factors

The folder `pollution` contains mortality rates and various environmental factors from 60 U.S. metropolitan areas from McDonald, G.C. and Schwing, R.C. (1973) 'Instabilities of regression estimates relating air pollution to mortality', *Technometrics*, vol.15, 463-482.

Variables, in order:

- PREC Average annual precipitation in inches
- JANT Average January temperature in degrees F
- JULY Same for July
- OVR65 % of 1960 SMSA population aged 65 or older
- POPN Average household size
- EDUC Median school years completed by those over 22
- HOUS % of housing units which are sound & with all facilities
- DENS Population per sq. mile in urbanized areas, 1960
- NONW % non-white population in urbanized areas, 1960
- WWDRK % employed in white collar occupations
- POOR % of families with income < \$3000
- HC Relative hydrocarbon pollution potential
- NOX Same for nitric oxides
- SO@ Same for sulphur dioxide
- HUMID Annual average % relative humidity at 1pm
- MORT Total age-adjusted mortality rate per 100,000

For this exercise we shall model mortality rate given nitric oxides, sulfur dioxide, and hydrocarbons as inputs. This model is an extreme oversimplification as it combines all sources of mortality and does not adjust for crucial factors such as age and smoking. We use it to illustrate log transformations in regression.

1. Create a scatterplot of mortality rate versus level of nitric oxides. Do you think linear regression will fit these data well? Fit the regression and evaluate a residual plot from the regression.
2. Find an appropriate transformation that will result in data more appropriate for linear regression. Fit a regression to the transformed data and evaluate the new residual plot.
3. Interpret the slope coefficient from the model you chose in 2.
4. Construct 99% confidence interval for slope coefficient from the model you chose in 2 and interpret them.
5. Now fit a model predicting mortality rate using levels of nitric oxides, sulfur dioxide, and hydrocarbons as inputs. Use appropriate transformations when helpful. Plot the fitted regression model and interpret the coefficients.
6. Cross-validate: fit the model you chose above to the first half of the data and then predict for the second half. (You used all the data to construct the model in 4, so this is not really cross-validation, but it gives a sense of how the steps of cross-validation can be implemented.)

Study of teenage gambling in Britain

1. Fit a linear regression model with gamble as the response and the other variables as predictors and interpret the coefficients. Make sure you rename and transform the variables to improve the interpretability of your regression model.
2. Create a 95% confidence interval for each of the estimated coefficients and discuss how you would interpret this uncertainty.
3. Predict the amount that a male with average status, income and verbal score would gamble along with an appropriate 95% CI. Repeat the prediction for a male with maximal values of status, income and verbal score. Which CI is wider and why is this result expected?

School expenditure and test scores from USA in 1994-95

1. Fit a model with total sat score as the outcome and expend, ratio and salary as predictors. Make necessary transformation in order to improve the interpretability of the model. Interpret each of the coefficient.
2. Construct 98% CI for each coefficient and discuss what you see.
3. Now add takers to the model. Compare the fitted model to the previous model and discuss which of the model seem to explain the outcome better?

Conceptual exercises.

Special-purpose transformations:

For a study of congressional elections, you would like a measure of the relative amount of money raised by each of the two major-party candidates in each district. Suppose that you know the amount of money raised by each candidate; label these dollar values D_i and R_i . You would like to combine these into a single variable that can be included as an input variable into a model predicting vote share for the Democrats.

Discuss the advantages and disadvantages of the following measures:

- The simple difference, $D_i - R_i$
- The ratio, D_i/R_i
- The difference on the logarithmic scale, $\log D_i - \log R_i$
- The relative proportion, $D_i/(D_i + R_i)$.

Transformation

For observed pair of x and y , we fit a simple regression model

$$y = \alpha + \beta x + \epsilon$$

which results in estimates $\hat{\alpha} = 1$, $\hat{\beta} = 0.9$, $SE(\hat{\beta}) = 0.03$, $\hat{\sigma} = 2$ and $r = 0.3$.

1. Suppose that the explanatory variable values in a regression are transformed according to the $x^* = x - 10$ and that y is regressed on x^* . Without redoing the regression calculation in detail, find $\hat{\alpha}^*$, $\hat{\beta}^*$, $\hat{\sigma}^*$, and r^* . What happens to these quantities when $x^* = 10x$? When $x^* = 10(x - 1)$?
2. Now suppose that the response variable scores are transformed according to the formula $y^{**} = y + 10$ and that y^{**} is regressed on x . Without redoing the regression calculation in detail, find $\hat{\alpha}^{**}$, $\hat{\beta}^{**}$, $\hat{\sigma}^{**}$, and r^{**} . What happens to these quantities when $y^{**} = 5y$? When $y^{**} = 5(y + 2)$?
3. In general, how are the results of a simple regression analysis affected by linear transformations of y and x ?

4. Suppose that the explanatory variable values in a regression are transformed according to the $x^* = 10(x - 1)$ and that y is regressed on x^* . Without redoing the regression calculation in detail, find $SE(\hat{\beta}^*)$ and $t_0^* = \hat{\beta}^*/SE(\hat{\beta}^*)$.
5. Now suppose that the response variable scores are transformed according to the formula $y^{**} = 5(y + 2)$ and that y^{**} is regressed on x . Without redoing the regression calculation in detail, find $SE(\hat{\beta}^{**})$ and $t_0^{**} = \hat{\beta}^{**}/SE(\hat{\beta}^{**})$.
6. In general, how are the hypothesis tests and confidence intervals for β affected by linear transformations of y and x ?

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.