

# Data 621 Final Project

*Omar Pineda, Jeff Littlejohn, Sergio Ortega Cruz, Chester Poon, Simon Ustoyev*

*Due Date December 20th*

## Contents

Data 621 Final Project . . . . .	2
Abstract . . . . .	2
Keywords . . . . .	2
Introduction . . . . .	2
Literature Review . . . . .	2
Methodology . . . . .	4
Experimentation and Results: . . . . .	4
The research question . . . . .	4
Understanding the Data Set . . . . .	4
Checking Multicollinearity . . . . .	8
Pairwise Scatterplots . . . . .	8
Correlation Matrices . . . . .	9
Variance Inflation Factor . . . . .	9
Visualizing the Data . . . . .	10
Plotting . . . . .	10
Adding Noise to the Scatterplot . . . . .	11
Multiple Linear Regression . . . . .	12
Building the Model . . . . .	12
Criteria for Comparing all Kinds of Models . . . . .	12
For Comparing Nested Models . . . . .	12
Creating the model . . . . .	12
Stepwise Selection . . . . .	12
All subsets . . . . .	13
Creating a New Model . . . . .	14
Checking Model Assumptions . . . . .	14
Outliers . . . . .	18
Indicator Variables . . . . .	18
The Regression Equation . . . . .	19
Interpreting Regression Coefficients . . . . .	19
Confidence Intervals for Regression Coefficients . . . . .	20
T-Tests . . . . .	20
Making Predictions . . . . .	20
Confidence Intervals . . . . .	21
Adding Interaction Terms . . . . .	22
Creating a New Model . . . . .	22
Checking Model Assumptions . . . . .	22
The Regression Equation . . . . .	25
The T-test Interpretation . . . . .	26
F-tests . . . . .	26
Partial F-tests . . . . .	27
Discussion and Conclusions: . . . . .	28
APPENDIX . . . . .	28
References . . . . .	28
Code . . . . .	28

# Data 621 Final Project

## Abstract

## Keywords

Demographic, Education, Attainment, Scores, Equality.

## Introduction

Despite the national narrative that America affords its inhabitants an unprecedented land of opportunity, intergenerational social mobility, defined as the likelihood that a child born to parents in the bottom fifth of the income distribution reaches the top fifth, is higher in many other advanced countries. Fewer than eight percent of Americans born in the bottom 20% of the income distribution reach the top 20%, whereas more than 13 percent of Canadians do. However, as Harvard Economist Raj Chetty has demonstrated, significant differences in upward mobility rates exist across the United States. “In this country, of all countries, a person’s zip code shouldn’t decide their destiny,” President Barack Obama said in 2015. We wish to understand which demographic and geographic factors are associated with social mobility. We will use the CollegeDistance data set that is included in the AER R package. It contains survey data from 1100 USA residents who were high school seniors in 1980 with follow-up data regarding education attainment in 1986. While educational attainment is not a perfect proxy for income, we will use it as our outcome variable, our metric of interest. We will start the construction of the regression model with the 13 other variables, including gender, ethnicity, a composite test score, school and community demographic factors, and family income, among others. After exploratory data analysis that includes checking multicollinearity, we will endeavor to construct a multiple linear or multiple logistic regression model that predicts educational attainment based on the independent variables. While our data captures educational attainment from more than 30 years ago, we think the construction of a model that predicts such a metric is nonetheless a worthwhile endeavor. An interesting follow-up study may be to repeat such a survey in the present day, construct a model, and then analyze differences.

## Literature Review

We reviewed three papers for our review that look at social mobility where each study had a slightly different approach at analyzing the different independent variables that affect intergenerational social and economic outcome. Broadly, the first paper reviewed observed economic outcomes children in families that moved from a poorer to a better income neighborhood. The second paper researched the interplay of a child’s gender from poorer neighborhoods versus higher income neighborhoods and those children’s adult economic outcome. While these two papers studies were focused on communities in the United States, our third paper reviewed looked at differences in Denmark that draw heavily from the theories of sociologist Pierre Bourdieu on the concept of “cultural capital”, which is the theory that individuals can possess a form of “capital” quantified by the amount of knowledge and behavior (culture) that promote social mobility in a socially stratified society (Møllegaard & Jæger 2015). The first paper we looked at, The Impacts of Neighborhoods on Intergenerational Mobility I: Childhood Exposure Effects found that families that moved from a lower to a better income area had children that had greater social mobility than those that stayed. There was a positive improvement in the child’s adult income, which rose at a rate of about 4% that was proportional the number of years spent growing up in the higher income area (Chetty & Hendren 2016). The dataset composed of federal income tax records from 1996 through 2012 and focused on families with children born between 1980 and 1988 and moved across neighborhoods between 1997 and 2010. This paper reconciles conflicting papers by introducing the number of years growing up in the better neighborhood. Interestingly, a symmetrical finding was found for families moving from higher to lower income neighborhoods in which the children of those families had adult incomes that reduced 4% per years lived in the lower income neighborhood (Chetty & Hendren 2016). The key variables are fairly different from the ones we are using from the dataset in AER.

For this paper, those variables are parent’s income, parent location, child’s adult earned income, teenage birth, marriage, educational attainment, and the child’s employment at adulthood. While the primary independent variable is the child’s adult earned income, ours is educational attainment. Our analysis is also more focused on key variables that are characteristic of a neighborhood, whereas this study is more focused on those associated with the family unit. Our dataset also does not contain any interventional events with each family as this study has. A way to take our study further would be to complete the dataset in AER with family income as well as the child’s adult earned income. The second paper we reviewed, *Childhood Environment and Gender Gaps in Adulthood* takes a deeper look into the effects of gender and social mobility. Like the first paper we reviewed, the data was drawn from tax records for families with children born in the 1980s. The major finding of this study found that the traditional gender gap in employment and income earnings are reversed in poor families, particularly in high-poverty neighborhoods. Boys from these families are less likely to work in adulthood and the issue is compounded with single-parent families. The authors write: Low-income boys who grow up in high-poverty, high-minority areas work less than girls. These areas also have higher rates of crime, consistent with a model in which boys with lower latent earnings potential who grow up in environments of concentrated poverty switch from the formal labor market to crime or other illicit activities (Chetty et al. 2016).

Though the paper does a comprehensive analysis on the impact of environment based on economic class and the interplay of gender, a deeper analysis that would include race as a more prominent factor would strengthen this paper significantly. This is especially true if we approach the analysis of these factors (race, class, gender, etc.) as observations that are not mutually exclusive. The last two papers we reviewed included clear and concrete independent variables for their analyses, the third paper we reviewed is different in its approach in many ways. The data quantifies three different forms of capital as measured in three major categories: economic, social and cultural. (Møllegaard & Jæger 2015). Economic capital is self-explanatory in that it’s a measurement of income and assets. Cultural capital is measured in less concrete terms such as years of educational attainment of the parents and grandparents, subscriptions to news outlets, and other things. Social capital is interesting in that it is best understood as a kind of network analysis manifesting as professional networks, friendships, familial connections, etc. The paper also addresses the intergenerational effects on educational attainment as far as two-generations, whereas most studies review a single generation prior. The findings of the paper found that in Denmark, educational attainment or more academically driven educational attainment depended more on cultural capital than of social/economic capital (Møllegaard & Jæger 2015). While these measures are fascinating and the results of the study show that cultural capital measures are certainly valid and useful independent variables, it’s difficult to apply an appropriate weight to some of the cultural capital examples. Generally, applying the appropriate weight of the impact of various aspects of cultural capital are difficult to accomplish and to justify. Further research is needed to properly assess the impact of specific variables associated with cultural capital to apply a stronger model in predicting educational outcomes. For our project, measuring social and cultural capital appropriately will be unfeasible. Overall, a common measure of all reviewed papers focuses on parental income or economic means to assess the probability of social upward mobility of the child’s adulthood income. This “intergenerational mobility” is the common theme as a measure of social mobility. For our project, while we are not focused on income as a measure, our project is focused primarily on educational attainment to serve as a proxy for social mobility. We are also looking to identify the strongest predictors for upward social mobility as measured in educational attainment.

## Methodology

### Experimentation and Results:

#### The research question

Are the numbers of years of education beyond highschool related to the relationship of elements such as exam scores, tuition paid, distance they had to travel to get to school and/or if their highschool was or not situated in and urban/rural environment?

#### Understanding the Data Set

The first data set we will use for this R Guide is CollegeDistance from the R package AER. This is a cross-sectional data set from a survey conducted by the Department of Education in 1980, with a follow-up in 1986, containing 14 variables relating to the characteristics of the students surveyed for this data set, their families and the area in which they live. We can see the names of all of the variables in the data set by using the command `names(dataset)` and see how the observations for each of these variables look by using the command `head(dataset)`

```
## Warning: package 'AER' was built under R version 3.5.3
## Loading required package: car
## Warning: package 'car' was built under R version 3.5.3
## Loading required package: carData
## Loading required package: lmtest
## Loading required package: zoo
##
## Attaching package: 'zoo'
## The following objects are masked from 'package:base':
##
##      as.Date, as.Date.numeric
## Loading required package: sandwich
## Loading required package: survival
## Warning: package 'survival' was built under R version 3.5.2
## Warning: package 'reshape2' was built under R version 3.5.3
## Warning: package 'dplyr' was built under R version 3.5.3
##
## Attaching package: 'dplyr'
## The following object is masked from 'package:car':
##
##      recode
## The following objects are masked from 'package:stats':
##
##      filter, lag
```

```
## The following objects are masked from 'package:base':
##
## intersect, setdiff, setequal, union
## Warning: package 'psych' was built under R version 3.5.3
##
## Attaching package: 'psych'
## The following objects are masked from 'package:ggplot2':
##
## %+%, alpha
## The following object is masked from 'package:car':
##
## logit
## Warning: package 'skmr' was built under R version 3.5.3
## gender ethnicity score fcollege mcollege home urban unemp wage distance
## 1 male other 39.15 yes no yes yes 6.2 8.09 0.2
## 2 female other 48.87 no no yes yes 6.2 8.09 0.2
## 3 male other 48.74 no no yes yes 6.2 8.09 0.2
## 4 male afam 40.40 no no yes yes 6.2 8.09 0.2
## 5 female other 40.48 no no no yes 5.6 8.09 0.4
## 6 male other 54.71 no no yes yes 5.6 8.09 0.4
## tuition education income region
## 1 0.88915 12 high other
## 2 0.88915 12 low other
## 3 0.88915 12 low other
## 4 0.88915 12 low other
## 5 0.88915 13 low other
## 6 0.88915 12 low other
## [1] "gender" "ethnicity" "score" "fcollege" "mcollege" "home"
## [7] "urban" "unemp" "wage" "distance" "tuition" "education"
## [13] "income" "region"
```

Though the data set contains 14 different variables, in this R Guide, we will only use score, the achievement test score obtained during the student's senior year of high school, urban, whether the student's high school is located in an urban area, distance, the distance the student lives from a 4-year college (in 10's of miles), tuition, the average 4-year college tuition in the student's state (in 1000's of dollars), and education, the number of years of education attained 6 years after high school graduation, in any of the models we create. Of these chosen variables, we can see that score, distance, tuition and education are quantitative and urban is categorical.

```
## vars n mean sd median trimmed mad min max range skew
## gender* 1 4739 1.55 0.50 2.00 1.56 0.00 1.00 2.00 1.00 -0.20
## ethnicity* 2 4739 1.55 0.79 1.00 1.43 0.00 1.00 3.00 2.00 0.99
## score 3 4739 50.89 8.70 51.19 50.91 10.24 28.95 72.81 43.86 -0.03
## fcollege* 4 4739 1.21 0.41 1.00 1.14 0.00 1.00 2.00 1.00 1.44
## mcollege* 5 4739 1.14 0.34 1.00 1.05 0.00 1.00 2.00 1.00 2.11
## home* 6 4739 1.82 0.38 2.00 1.90 0.00 1.00 2.00 1.00 -1.67
## urban* 7 4739 1.23 0.42 1.00 1.17 0.00 1.00 2.00 1.00 1.26
## unemp 8 4739 7.60 2.76 7.10 7.32 2.22 1.40 24.90 23.50 1.56
## wage 9 4739 9.50 1.34 9.68 9.47 1.17 6.59 12.96 6.37 0.09
## distance 10 4739 1.80 2.30 1.00 1.36 1.04 0.00 20.00 20.00 3.00
## tuition 11 4739 0.81 0.34 0.82 0.82 0.45 0.26 1.40 1.15 -0.15
```

```
## education    12 4739 13.81 1.79 13.00 13.66 1.48 12.00 18.00 6.00 0.44
## income*     13 4739 1.29 0.45 1.00 1.24 0.00 1.00 2.00 1.00 0.94
## region*     14 4739 1.20 0.40 1.00 1.12 0.00 1.00 2.00 1.00 1.51
##             kurtosis  se
## gender*      -1.96 0.01
## ethnicity*    -0.69 0.01
## score         -0.88 0.13
## fcollege*     0.07 0.01
## mcollege*     2.44 0.01
## home*         0.78 0.01
## urban*        -0.40 0.01
## unemp         5.40 0.04
## wage          -0.26 0.02
## distance      13.03 0.03
## tuition       -1.05 0.00
## education     -1.23 0.03
## income*       -1.12 0.01
## region*       0.27 0.01
```

Table 1: Data summary

Name	CollegeDistance
Number of rows	4739
Number of columns	14
Column type frequency:	
factor	8
numeric	6
Group variables	None

#### Variable type: factor

skim_variable	n_missing	complete_rate	ordered	n_unique	top_counts
gender	0	1	FALSE	2	fem: 2600, mal: 2139
ethnicity	0	1	FALSE	3	oth: 3050, his: 903, afa: 786
fcollege	0	1	FALSE	2	no: 3753, yes: 986
mcollege	0	1	FALSE	2	no: 4088, yes: 651
home	0	1	FALSE	2	yes: 3887, no: 852
urban	0	1	FALSE	2	no: 3635, yes: 1104
income	0	1	FALSE	2	low: 3374, hig: 1365
region	0	1	FALSE	2	oth: 3796, wes: 943

#### Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100	hist
score	0	1	50.89	8.70	28.95	43.92	51.19	57.77	72.81	
unemp	0	1	7.60	2.76	1.40	5.90	7.10	8.90	24.90	
wage	0	1	9.50	1.34	6.59	8.85	9.68	10.15	12.96	
distance	0	1	1.80	2.30	0.00	0.40	1.00	2.50	20.00	
tuition	0	1	0.81	0.34	0.26	0.48	0.82	1.13	1.40	

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100	hist
education	0	1	13.81	1.79	12.00	12.00	13.00	16.00	18.00	

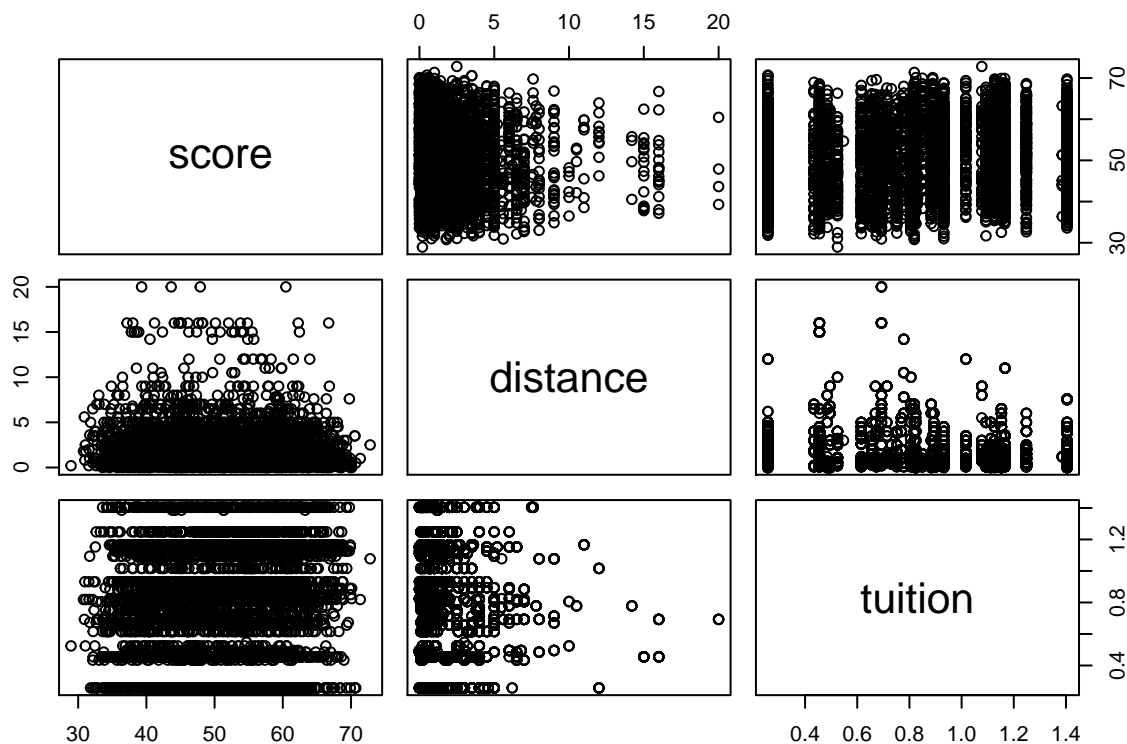
## Checking Multicollinearity

One of the very first things we want to do before making any model is check for multicollinearity between our quantitative predictor variables. Multicollinearity is a problem because it will inflate the standard error of a model as well as make the parameter estimates inconsistent.

So, to look at multicollinearity, let's first look at the pairwise scatterplots for each of our quantitative predictor variables.

### Pairwise Scatterplots

We want to look at pairwise scatterplots of our predictor variables so that we can determine visually if any pairs of predictors appear highly correlated. To do this, we can use the command `plot(variable, variable, ..., variable)` using the variables of interest as our arguments, or more simply, specify the columns to use within your data frame like this:

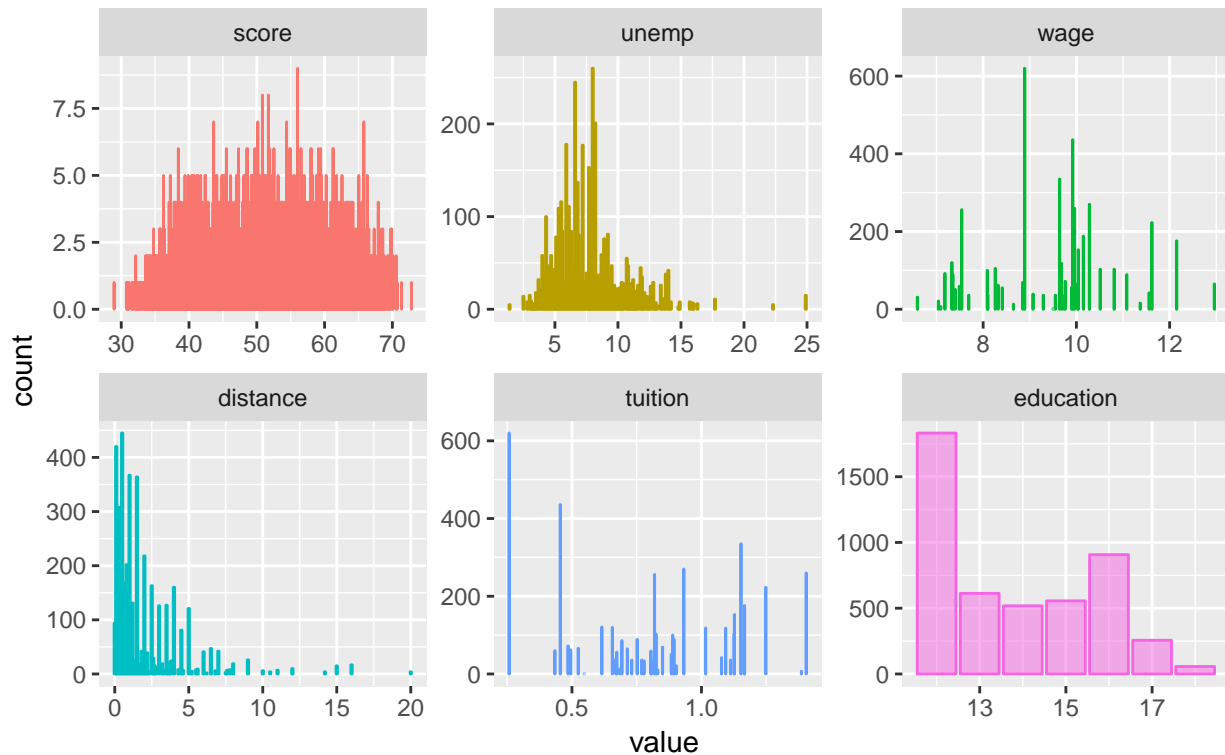


Looking at our plot, it does not appear that any of our quantitative predictor variables are highly correlated, or have a strong linear relationship with one another.

```
## No id variables; using all as measure variables
```



## Distribution of all Continuous Variables



## Correlation Matrices

Just to be certain there is no strong correlation, let's calculate the correlation matrix of the three variables using `cor()` which takes the variables of interest as its arguments, similar to `plot()`

```
##           score  distance  tuition
## score      1.00000000 -0.06797927  0.1298585
## distance  -0.06797927  1.00000000 -0.1009806
## tuition    0.12985848 -0.10098058  1.0000000
```

We see that the pairwise correlations between our quantitative predictor variables are very low. We only consider multicollinearity to be a problem when the absolute value of the correlation is greater than .9, so we would say that there is no multicollinearity issue in this model.

## Variance Inflation Factor

I'll use a third way to detect multicollinearity in a model. We consider the variance inflation factor (VIF) of a variable by first making a model using that variable as the response and all other predictor variables in the original model as the predictors. For example, to calculate the VIF of achievement test score we need to create this linear model:

Then, we need to plug the R-squared value of that model into our VIF equation,  $VIF = \frac{1}{1-r^2}$

```
## [1] 1.020309
```

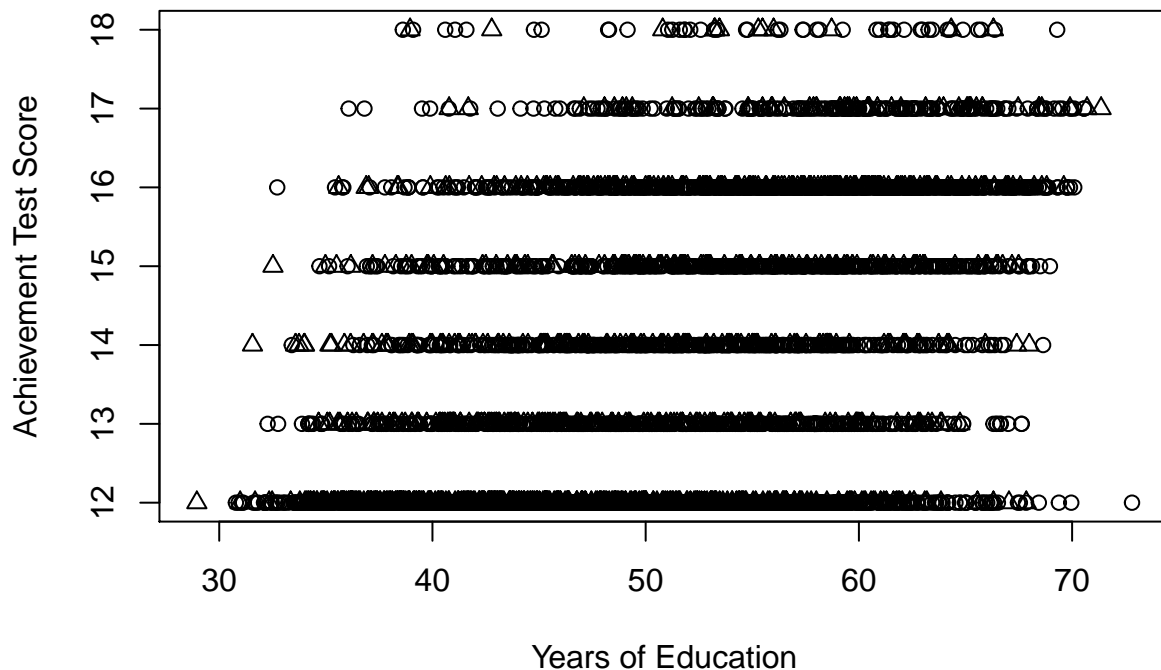
So the VIF of achievement test score for this model is around 1.02. This is great news since we only consider multicollinearity to be an issue when  $VIF > 10$ . We could then continue checking VIF by following the same process for our distance and tuition variables.

## Visualizing the Data

Now that we know there is no evidence of multicollinearity between our predictor variables, we can move on to visualizing some of the data. For simple linear regression, we easily built a scatterplot for exploratory data analysis since we only had two variables; however, in many multiple linear regression situations, the variables we are using cannot be simultaneously represented two-dimensionally so exploring the data visually is far more difficult. Nevertheless, using an example, I will demonstrate that we can represent a fairly uncomplicated multiple linear regression situation two-dimensionally. For this, we will need a quantitative response variable and two predictor variables, one quantitative and one categorical.

## Plotting

We can plot the student's achievement test score, the quantitative response, versus the number of years of education they have attained, the quantitative predictor, varying the shape of their plotted point by whether or not their high school is located in an urban area, the categorical predictor, because this satisfies our variable requirements. We'll use the `plot()` command, to make the scatterplot, adding one unusual argument, `pch = variable`, which varies the shape of the points by the value of a certain variable. The variable must be numeric.

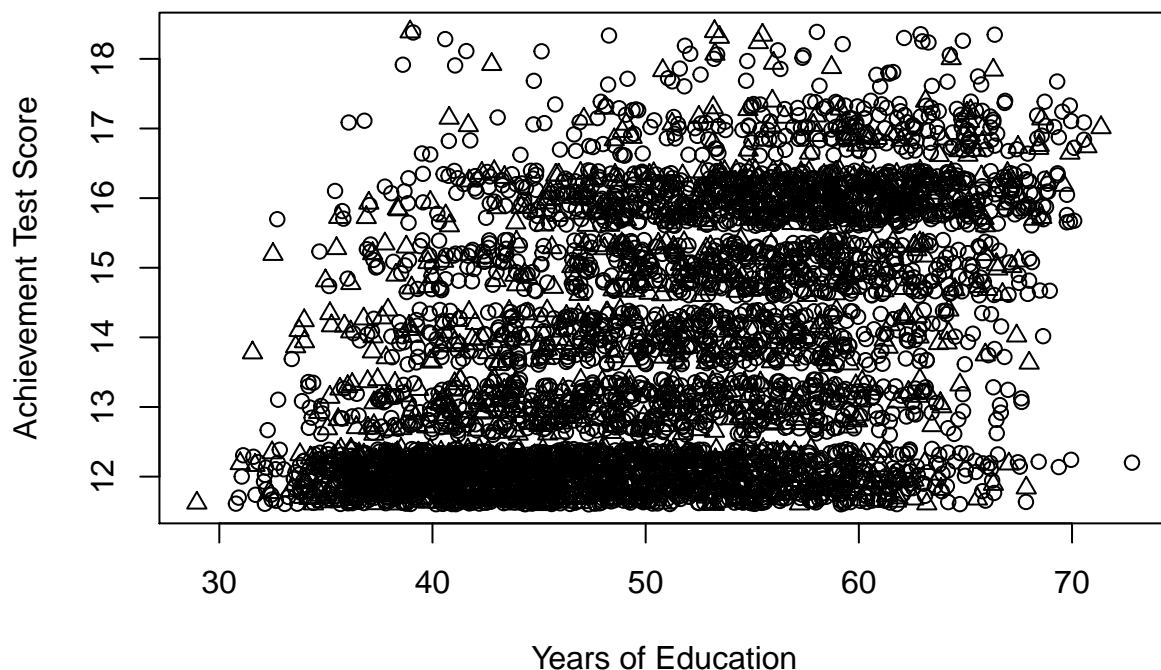


The circles indicate that the student represented went to a high school in an urban area. The triangles indicate that the student represented went to a high school in a non-urban area. From this plot we can't

really see whether there is any difference between the two groups because there are a lot of overlapping data points. We'll try to fix this by adding noise to, or slightly increasing the variation of, the scatterplot so that we can see any patterns in the data more clearly.

### Adding Noise to the Scatterplot

We can add more variation to the scatterplot by using the command `jitter(variable, factor = number)` on our predictor and response variables, where `factor` is the argument specifying how much noise to add to the plot, and then using these altered variables as our new predictor and response like this:



Using the `jitter()` command has unveiled some of the hidden data points. We can now more distinctly see two types of data points, circles and triangles. However, there doesn't appear to be a distinct separation of the data points. This might lead us to believe that there is not a significant difference between the educational attainment of students at urban and non-urban schools, but we can confirm whether or not this is true using a test later on.

## Multiple Linear Regression

### Building the Model

#### Criteria for Comparing all Kinds of Models

There are many types of criteria that we can use to determine whether or not one model of a set of data is better than another model of that same set of data. First, when comparing any two models we can use the Akaike Information Criteria (AIC) or Mallow's Cp. When looking at AIC, we want to weigh the number of parameters in a model and the AIC value together. For example, if we have a model with an AIC value of X and 4 parameters, in order to prefer a model with more than 4 parameters, we would want the AIC value for that model to be at least 10 units less than X. Generally, we can think of lower AIC values corresponding to better models. Then, if we instead use Mallow's Cp, we see that a smaller value means the model fits better, so we would choose the model with the lowest Mallow's Cp. One important note on these comparison tools is that they have no units so they can only compare models made from the same data set, not models from varied data sets.

#### For Comparing Nested Models

If we are comparing nested models, i.e. one model's predictors are a subset of the other model's predictors, we have additional comparison options. For comparing nested models, we can use the p-values from t-tests or partial F-tests, only choosing the more complex model if the p-value of the t-test or F-test is below .05. We can also use adjusted R-squared as a comparison tool; however, I would discourage anyone from using adjusted R-squared as their comparison criteria because the other methods of comparison are more rigorous.

### Creating the model

#### Stepwise Selection

Now, on to the topic of actually building the model. First, we can discuss the idea of building a model in a stepwise fashion. There are two ways of creating models stepwise, either forward selection or backwards elimination. In forward selection, you start with the simplest desired model and add predictors to the model until your chosen criteria indicate that adding more predictors to the model would actually worsen the model's abilities. In backwards elimination, you start with the most complex model acceptable and remove predictors from the model until your chosen criteria indicate that removing more predictors from the model would worsen the model's abilities.

By their nature, these models will always be nested so any of the comparison methods previously mentioned can be used.

As an example, I'll be illustrating the backwards elimination method to create a model that predicts educational attainment 6 years after graduation using all or a subset of the variables discussed earlier and AIC as my model comparison criteria. We'll need to use the command

```
stepAIC(starting.model, scope = list(upper = complex.model, lower = simple.model), direction = "backward")
```

in the R package MASS to do this. `starting.model` is the model that we will begin the stepwise selection with, `complex.model` is the most complex model we would want to have, so in the case of backward elimination this would be the same as `starting.model` and `simple.model` is the most simplistic model you would want to obtain. I will specify the simplest model as a model using only the intercept to model educational attainment 6 years after graduation.

```
## Warning: package 'MASS' was built under R version 3.5.3
```

```
##
```

```
## Attaching package: 'MASS'
```

```
## The following object is masked from 'package:dplyr':
##
##      select
##
## Start:  AIC=4339.19
## education ~ score + urban + distance + tuition
##
##           Df Sum of Sq  RSS    AIC
## - urban      1      0.5 11815 4337.4
## <none>                11815 4339.2
## - tuition     1     10.8 11826 4341.5
## - distance     1     53.3 11868 4358.5
## - score        1    3178.2 14993 5466.2
##
## Step:  AIC=4337.39
## education ~ score + distance + tuition
##
##           Df Sum of Sq  RSS    AIC
## <none>                11815 4337.4
## - tuition     1       11 11826 4339.8
## - distance     1       62 11877 4360.2
## - score        1    3205 15020 5472.8
##
## Call:
## lm(formula = education ~ score + distance + tuition, data = CollegeDistance)
##
## Coefficients:
## (Intercept)      score    distance    tuition
##    9.15680    0.09547   -0.05013   -0.14369
```

We can see from the output that our final model actually contains all of the chosen variables, except for whether or not a student's high school is located in an urban area, as the optimal way of predicting that student's educational attainment 6 years after graduation. One thing to note about this process is that, although the two models' AIC differ by less than 10, the chosen model is the model with fewer predictor variables because of the necessary balance between accuracy and complexity that AIC uses.

We could easily use a forward selection method by specifying `starting.model` as the simplest model we would deem acceptable and changing direction to "forward". However, using a forward selection method might produce a different model than the backwards selection method because each process takes a different set of paths to reach the optimal model.

## All subsets

Lastly, we have the all subsets method. This method requires that you calculate all of the AIC values, or Mallows' Cp's, for all possible models that can be created from any subset of the variables you are considering for the model. This method is very time consuming and so, it is used infrequently. I will note that these models will not all be nested so we can only use AIC or Mallows' Cp as comparison methods.

Besides the automated method of model building, using `stepAIC()` for stepwise selection, we can also simply create models manually. Although these models may not be the best possible models with respect to any one of the numerous criteria we discussed, they do work well for demonstration purposes. Thus, throughout the rest of this R guide, the models used for examples will be created manually.

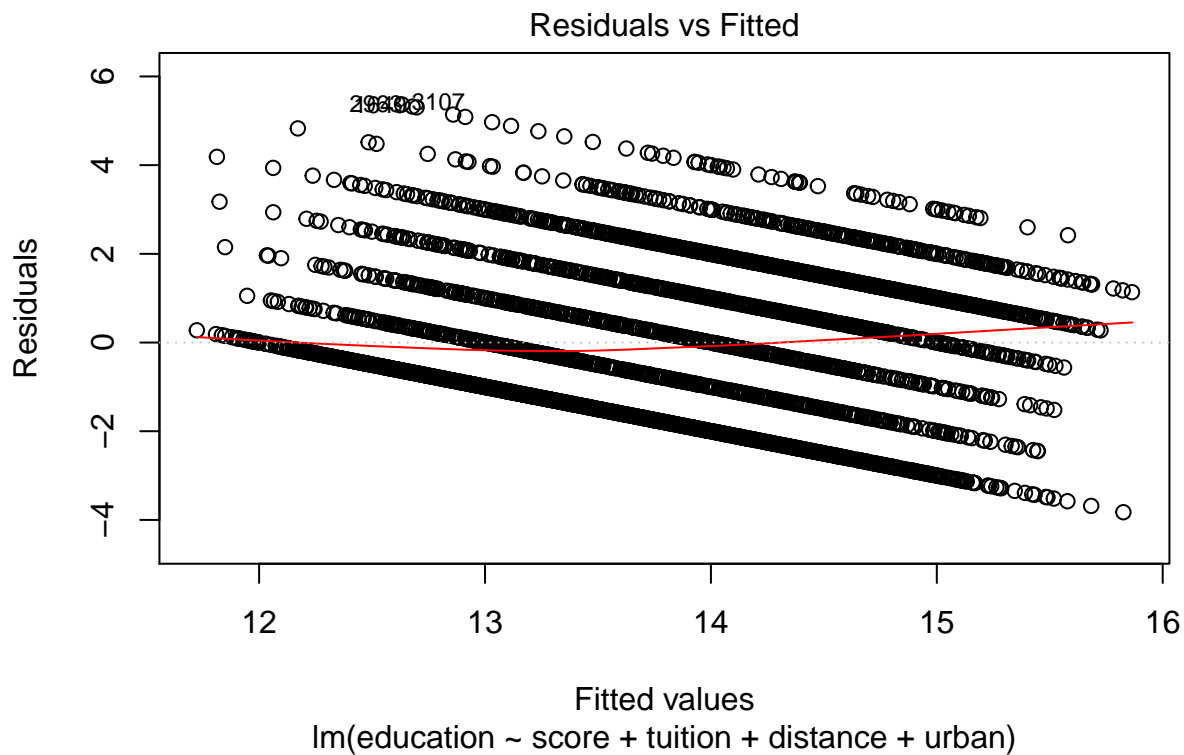
## Creating a New Model

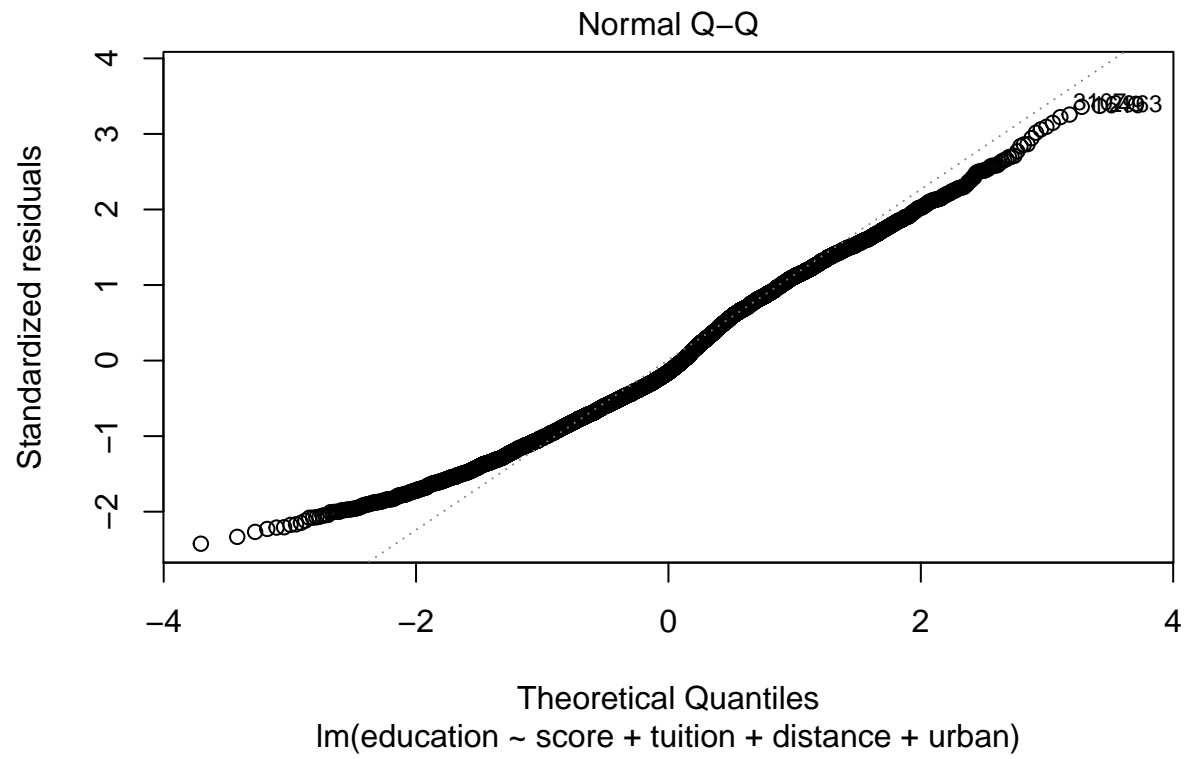
We've looked at our data set, gone through a short example of how one might visualize a relatively simple multiple linear regression situation and discussed some sophisticated techniques of model building as well as how to compare models. Now, we can create a multiple linear regression model manually. For this model, I will use achievement test score, average tuition for a four-year college in the student's state, distance from a four year college and whether or not the student's high school was located in an urban area to predict the number of years of education the student had attained six years after graduation.

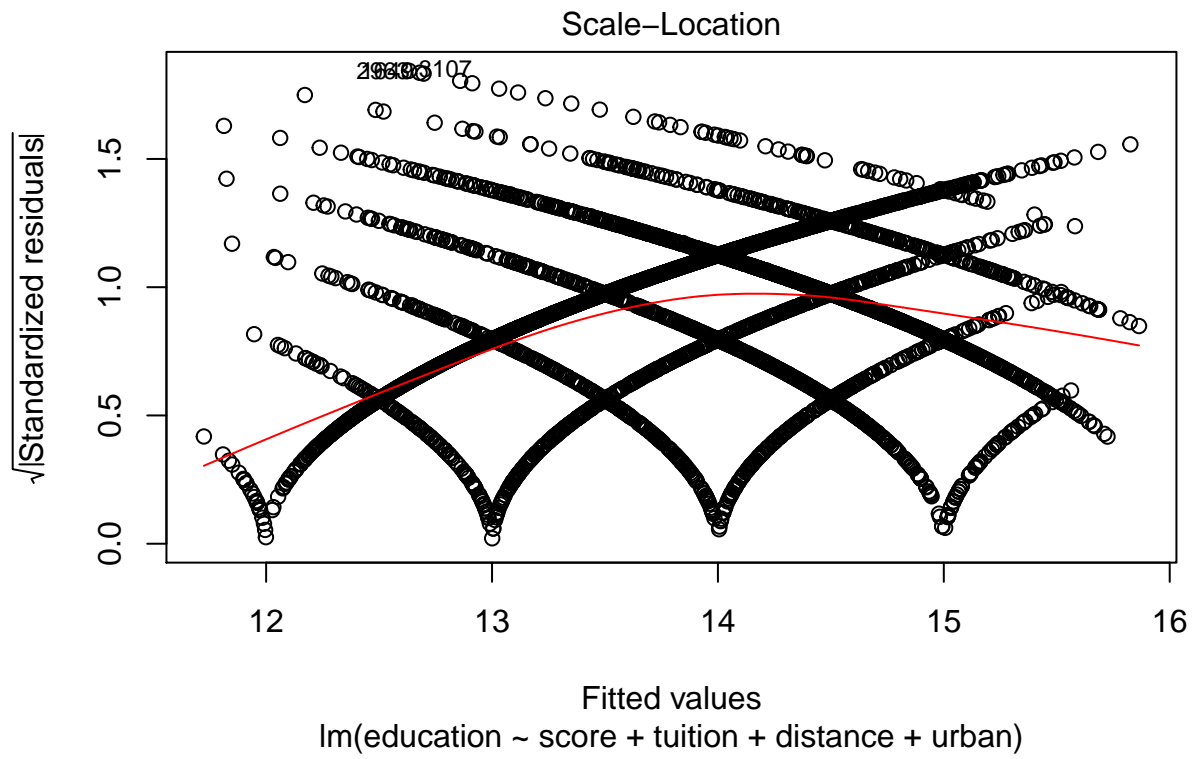
We can create the model similarly to how we created the simple linear regression model, only adding more predictors. We use `lm(response ~ predictor + predictor + ... + predictor, data = dataset)` to create our model.

## Checking Model Assumptions

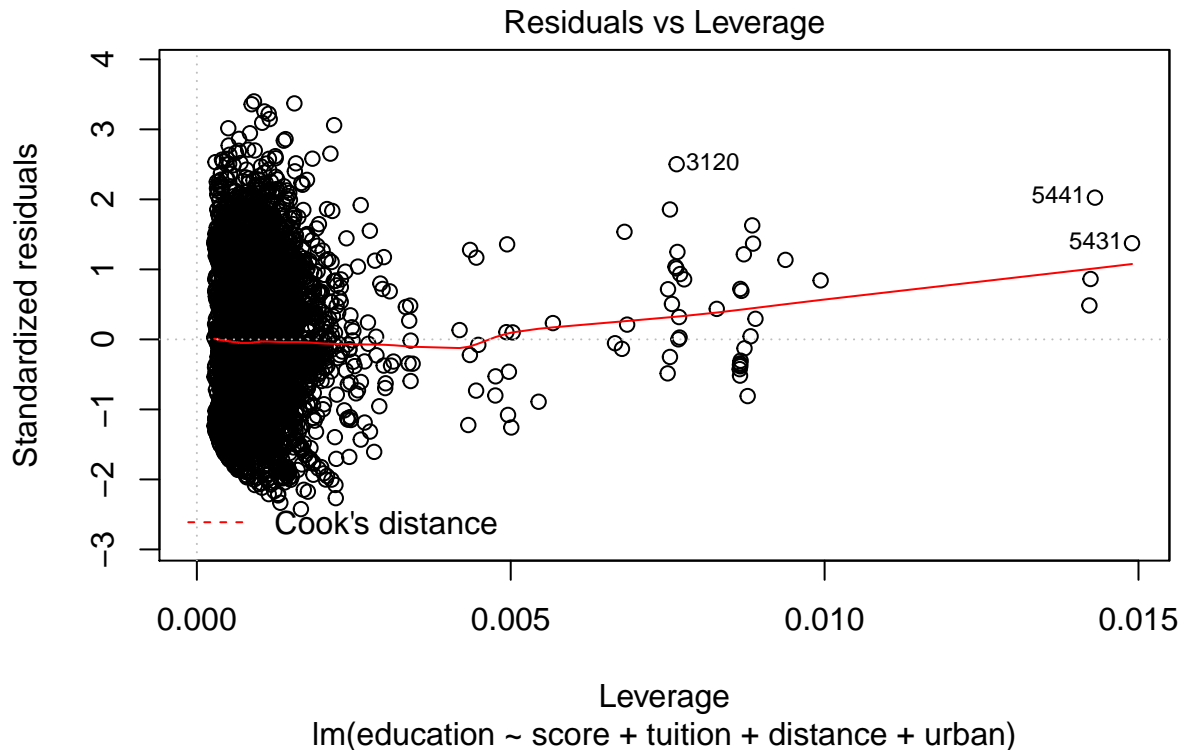
At this point, we need to check the assumptions of our model to be sure that the model provides accurate results. In my R guide on simple linear regression, we created the diagnostic plots for our model individually, but there is a function in R that allows us to create all of these plots at once. We create the plots using the function `plot(model)`.











The first plot is the residual plot, a comparison of the residuals of our model against the fitted values produced by our model, and is the most important plot because it can tell us about trends in our residuals, evidence of heteroskedasticity and possible outliers. The plot for this model indicates that our model is systematically underpredicting the lower values of educational attainment and systematically overpredicting the higher values of educational attainment. Although the residuals do not seem to be evenly spread around 0 for all fitted values, the range of the residuals at each fitted value appears to be roughly the same, so we can conclude there is no evidence of heteroskedasticity. Finally, this plot indicates that there are likely no outliers because there are no points on the plot well-separated from the rest.

The next plot is the QQ-plot. Though most of the points seem to fall on the line which indicates that our residuals come from a normal distribution, there are some points that stray from the line in the lower and upper quantiles of the plot. It is possible that these points do not come from a normal distribution, but most of our points seem to come from a normal distribution so there is not a lot to worry about here.

The third plot created is the scale-location plot. This plot is similar to the residual plot, but uses the square root of the standardized residuals instead of the residuals themselves. This makes trends in residuals more evident and, from this plot, we can see that there is likely a U-shaped trend in our residuals.

## Outliers

So, although our data did not appear to have any outliers, it's still important to know what to do if we do have outliers. Unfortunately, this is not an easy question to answer.

### The Definition

First, we should define what an outlier is. Generally, we define an outlier as any point well-separated from the rest. More specifically, we can define an outlier as any point with a studentized residual of greater than 2 or less than negative 2. This criteria indicates whether a point is an outlier in the y direction. Then, a data point outside the Cook's distance boundaries of a leverage plot would indicate that that point is an outlier in the x direction.

What can we do?

So, what can we do if we find that we have an outlier?

There is no simple answer to this question. However, the first thing you'll want to check is whether the data was recorded correctly. If you have access to the original data collection materials, check to see if the data was inputted incorrectly. If you do not have access to this information, but the recorded value is not contextually possible, you may remove this observation.

On the other hand, if the observation was recorded correctly, or if the observation is possible in the context of the study, you should think about why the observation is an outlier. One possible explanation for the outlier is that there is a predictor missing from the model that could explain the outlier. To really understand your outlier, you'll need to look at the specific context in which you are working and, hopefully, a reason for the outlier can be determined.

Finally, we see the leverage plot. This plot graphs the standardized residuals against their leverage. It also includes the Cook's distance boundaries. Any point outside of those boundaries would be an outlier in the x direction. Since we cannot even see the boundaries on our plot, we can conclude that we have no outliers.

## Indicator Variables

```
##
## Call:
## lm(formula = education ~ score + tuition + distance + urban,
##     data = CollegeDistance)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.826 -1.181 -0.246  1.219  5.367
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  9.141015   0.148905  61.388 < 2e-16 ***
## score        0.095596   0.002679  35.686 < 2e-16 ***
## tuition     -0.142627   0.068517  -2.082  0.0374 *
## distance    -0.048723   0.010539  -4.623 3.88e-06 ***
## urbanyes     0.025619   0.057090   0.449  0.6536
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.58 on 4734 degrees of freedom
## Multiple R-squared:  0.221, Adjusted R-squared:  0.2203
## F-statistic: 335.7 on 4 and 4734 DF, p-value: < 2.2e-16
```

We should first take note of how our categorical variable *urban* is written in our summary of the model. In the model, this variable is treated as an indicator variable meaning that it can take one of two values, either 0 or 1. In the summary, this variable is labeled as *urbanyes*. This means that the variable is given the value of 1 if a student's school is in an urban area and a 0 if a student's school is not in an urban area. Since this variable has only 2 categories, yes and no, it only needs one indicator function. However, if a categorical variable has more than 2 categories, it will need more indicator functions. To be exact, the variable would need the number of categories minus one indicator functions.

## The Regression Equation

From our summary, we can also get the multiple linear regression equations:

$$\widehat{\text{EducationalAttainment}} = 9.141 + .0956 \cdot \text{AchievementScore} - .1426 \cdot \text{Tuition} - .0487 \cdot \text{Distance} + .0256 \cdot I(\text{urban} = \text{yes})$$

Since we have an indicator variable, *urban*, in our equation, this equation can actually be broken down into two separate equations, one for if the student's school is located in an urban area and another for if the student's school is in a non-urban area. For students going to school in urban areas:

$$\widehat{\text{EducationalAttainment}} = (9.141 + .0256) + .0956 \cdot \text{AchievementScore} - .1426 \cdot \text{Tuition} - .0487 \cdot \text{Distance}$$

For students going to school in non-urban areas:

$$\widehat{\text{EducationalAttainment}} = 9.141 + .0956 \cdot \text{AchievementScore} - .1426 \cdot \text{Tuition} - .0487 \cdot \text{CollegeDistance}$$

Notice that the only change in the two equations is the intercept.

## Interpreting Regression Coefficients

Returning again to our summary, we can interpret the coefficients for each of our variables in a way similar to how they were interpreted for simple linear regression. For example, for every increase in achievement test score by 1 point, the predicted years of education of a student increases by .0956 years when that student's state's average tuition for a 4-year college, distance from a 4-year college and location of school (urban or non-urban) all stay the same. You might notice that the only change in the interpretation that we made from the interpretation for simple linear regression is that you must state that all of the other predictors in the model are being held constant.

We can also interpret the coefficient of our indicator function in a similar way. For example, the predicted years of education of a student increases by .0256 when their high school is located in an urban location as opposed to a non-urban location when their achievement test score, average state tuition for a four-year college and distance from a four-year college all remain the same.

Finally, thinking about the intercept in context, we can see that the intercept wouldn't really make sense to interpret because even a student that scored a zero on their achievement test, had an average state tuition for a 4-year college of zero dollars, lived zero miles from a 4-year college and went to a non-urban high school would still have at least 12 years of education since they had to have graduated in order to be part of this data set and our intercept is less than 12. We cannot use this model to determine anything about students who did not graduate.

## Confidence Intervals for Regression Coefficients

We recognize that these regression coefficient values given by our summary are simply point estimates so we'll need to create confidence intervals for them as well. We can create confidence intervals for them using the command `confint(model)`. This will generate 95% confidence intervals for all of the coefficients in our model.

```
##              2.5 %      97.5 %
## (Intercept)  8.84909261  9.432937106
## score        0.09034461  0.100848185
## tuition      -0.27695296 -0.008301033
## distance     -0.06938312 -0.028062221
## urbanyes     -0.08630305  0.137541361
```

As an example, we can interpret the confidence intervals for achievement test score and the indicator variable for whether a student's high school is located in an urban area. We interpret the confidence interval for achievement test score, a quantitative variable, like this: we are 95% confident that an increase by 1 point in achievement test score will increase the number of years of education a student achieves 6 years after high school graduation by between .09 and .1 years on average, all other variables held constant. We interpret the confidence interval for our indicator variable, a categorical variable, like this: we are 95% confident that a student going to high school in an urban area will achieve between -.086 and .138 years of education more than a student going to a high school in a non-urban area on average, all other variables held constant. Again, we cannot interpret our intercept because it does not make sense within the limits of our data set.

## T-Tests

From our summary, we can also determine the significance of each of our predictor variables through the results of a t-test. Like I discussed in my simple linear regression R guide, a t-test tests the null hypothesis that  $\beta = 0$  against the alternative hypothesis that  $\beta \neq 0$ , where  $\beta$  is the regression coefficient. The only difference for multiple linear regression is that we have multiple  $\beta$ s, so a different t-test is being run for each regression coefficient. In our summary table, we see the t-statistic for each of our predictor variables as well as the intercept and a corresponding p-value for each of the t-statistics. The results of the t-tests tell us that all of the variables in our model are significant predictors of a student's number of years of education except whether their high school is in an urban or non-urban area since that is the only p-value greater than .05. This means we can remove our urban variable from the model because it has no significant relationship with number of years of education that a student achieves 6 years after high school graduation. This is a similar conclusion to the one that was drawn from our automated model building procedure earlier.

## Making Predictions

Next, we can make some predictions using our model. We can use our regression equation to predict the years of education a student would have attained if they had a test score of 52, the average college tuition for their state was \$900, they lived five miles from a four-year college and they attended an urban high school.

```
## [1] 13.98511
```

From our equation, we can see that this student would have likely pursued just shy of 2 years of education after graduating high school. Now, let's compare that value with a similar student who only differed by where they went to high school. This student did not go to high school in an urban area.

```
## [1] 13.95951
```

From our equation, it seems that this student would have also pursued just shy of 2 years of education after graduating high school. However, take note of the difference between the two predicted years of education.

```
## [1] 0.0256
```

The difference between the two educational attainment levels is the value of the regression coefficient estimated for the urban indicator variable. Since this coefficient is small, the difference between the two students is also small. We know that providing a point estimate is usually not adequate when providing a prediction based on a regression model because it gives an answer that is very specific and likely not exactly correct. Therefore, we might want to make a prediction interval. We can do this using the command `predict(model, new.dataframe, interval)` `predict(model, new.dataframe, interval)`, similar to how we made prediction intervals for simple linear regression models. Our interval would be specified as “predict”

The only difference for the multiple linear regression models is that the new data frame must contain specified values for all of the variables in your regression model. So, our new data frame, corresponding to the example student that attended a high school in an urban area, would look like this:

```
## score tuition distance urban
## 1    52      0.9      0.5   yes
```

Now, we can use this data frame to create a 95% confidence interval for the student with a high school in an urban area that we made a point prediction for previously.

```
##          fit          lwr          upr
## 1 13.98492 10.88636 17.08349
```

We can see that the point estimate given by this prediction interval is slightly different than the value we calculated, however, this can be attributed to a rounding error. We see that our 95% prediction interval is [10.886, 17.083], thus, I am 95% confident that an individual student with an achievement score of 52, an average state tuition for four-year colleges of \$900, living five miles from a four-year college and going to a high school in an urban area will have achieved a total of between 10.886 and 17.083 years of education 6 years after graduating from high school. In this case, we have a slightly odd interval because we know that any student graduating from high school will attain at least 12 years of education 6 years after graduating from high school.

## Confidence Intervals

In addition to prediction intervals, we can also create confidence intervals for the mean years of education students with specific characteristics might attain 6 years after high school graduation. We can use a similar `predict(model, new.dataframe, interval)` command to do this, changing the interval argument from “predict” to “confidence”. We will use the same data frame that we created previously for the prediction interval.

```
##          fit          lwr          upr
## 1 13.98492 13.89019 14.07965
```

We get the same point prediction as we did for the prediction interval, which is to be expected. Also as we would expect, we see a smaller interval than we did for the prediction interval. From this interval we can say we are 95% confident that the mean number of years of education a student would attain 6 years after high school graduation, given that they had an achievement test score of 52, an average state tuition for four-year colleges of \$900, live five miles from the nearest four-year college and went to a high school in an urban area, is between 13.89 and 14.08 years.

## Adding Interaction Terms

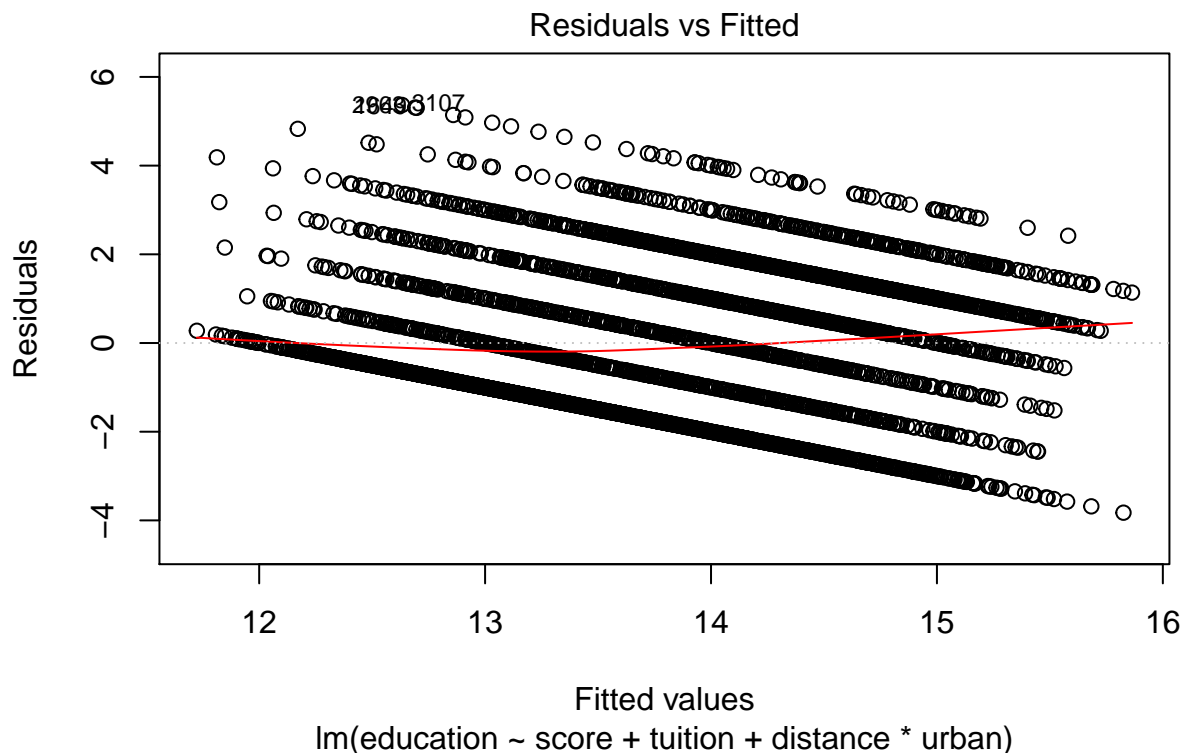
### Creating a New Model

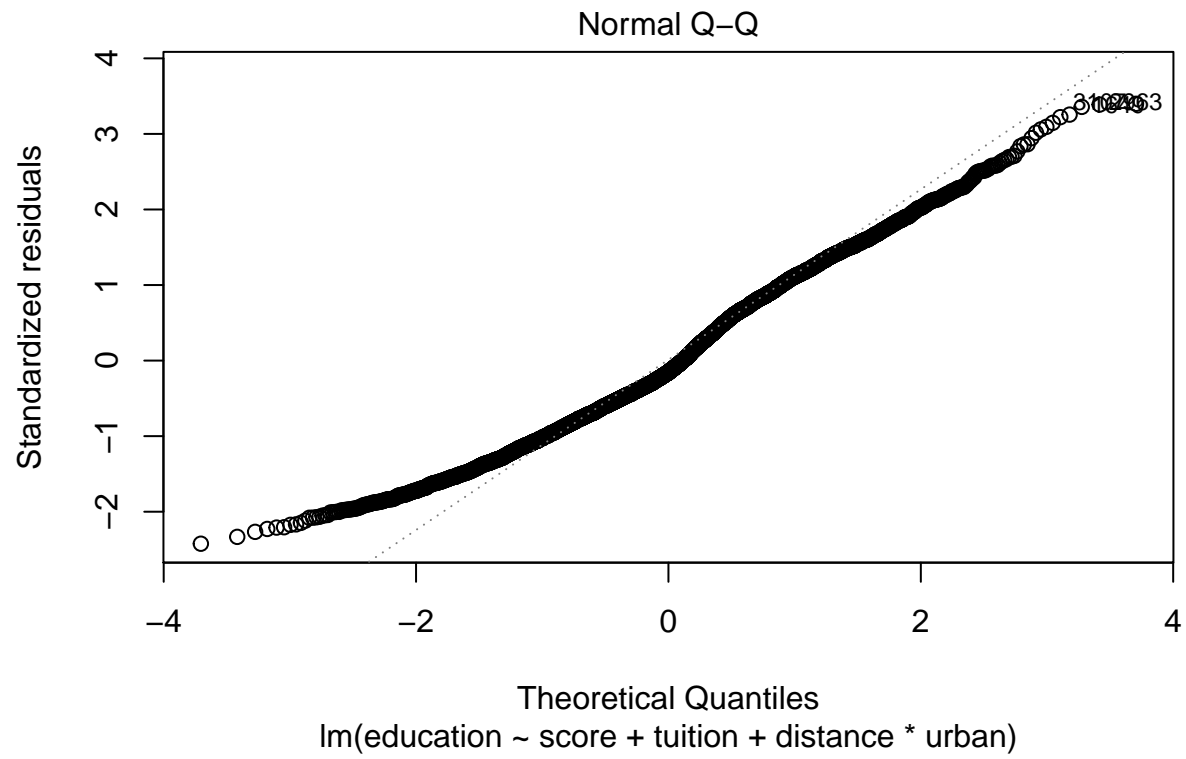
Another attribute we can add to our models is interaction terms. For this model we will use an interaction between our distance distance and urban variables. By that I mean that the relationship between the education level that a student achieves and the distance they are from a four-year college depends on whether the student's high school is located in an urban area or not. We can add the interaction term by using a \* between the two predictor variables included in our interaction term in the `lm()` command instead of +, like this:

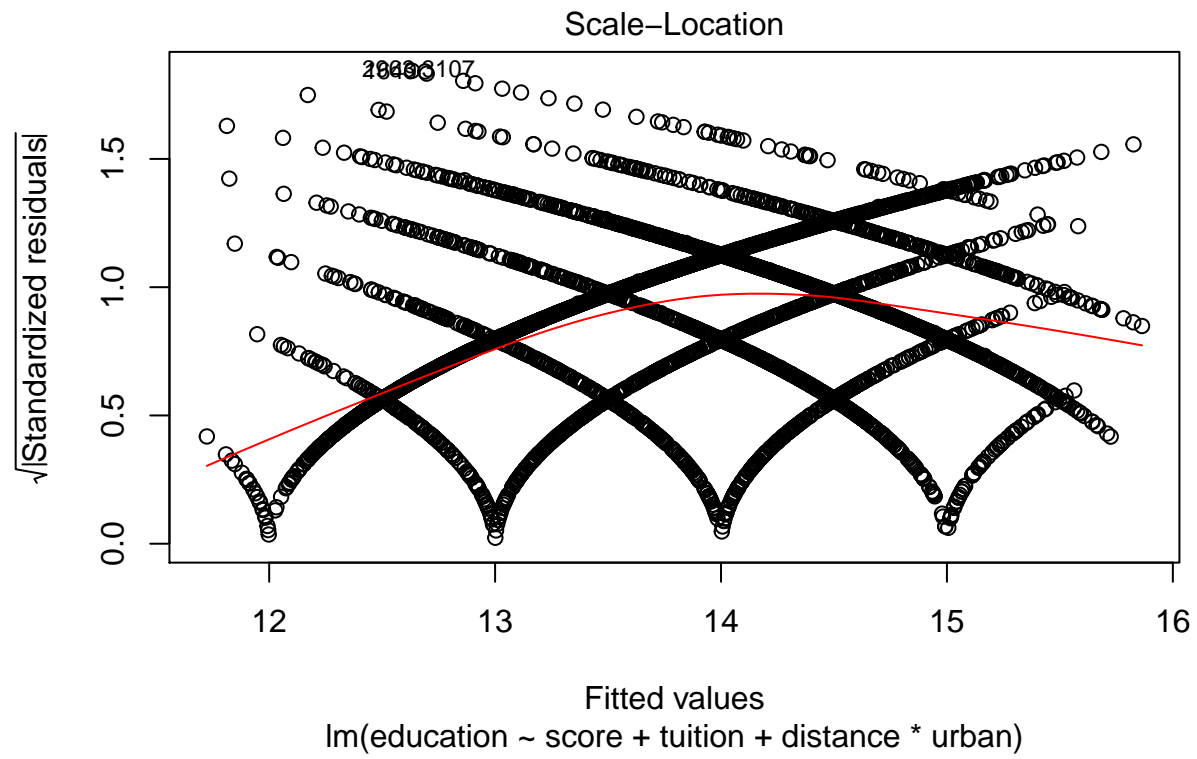
The \* denotes an interaction between two variables, while a + indicates that there is no interaction between the two variables it sits between.

### Checking Model Assumptions

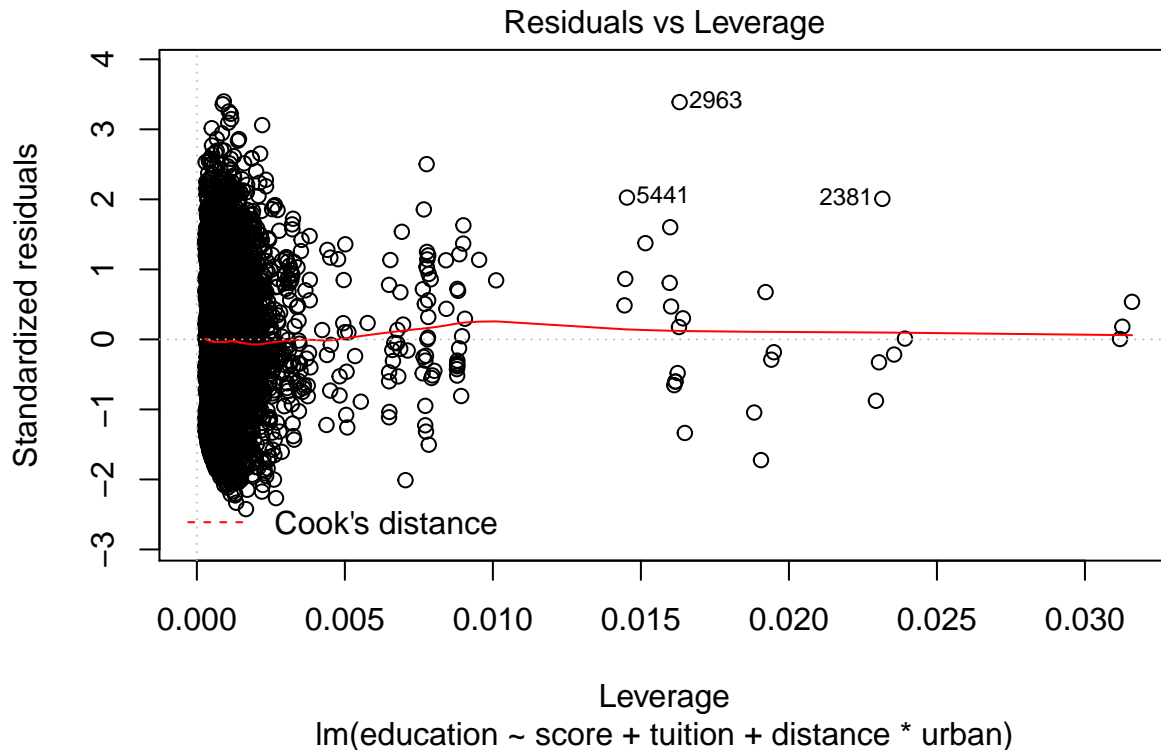
Since we have created a second model, we need to check our model assumptions again. We interpreted all of the diagnostic plots for the first model so we will just focus on the most important plot, the residual plot, for this one. We always need to check the residual plot to make sure the mean function is appropriate. We will use the same `plot()` command as before and so all four plots will still be generated











We see a similar residual plot to our first residual plot, with the same over-/underprediction problem. Again, there appears to be no strong evidence for heteroskedasticity. Finally, we see no outliers on this plot. Additionally, the other 3 diagnostic plots look very similar to the diagnostic plots generated for our first model.

### The Regression Equation

```
##
## Call:
## lm(formula = education ~ score + tuition + distance * urban,
##     data = CollegeDistance)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.8258 -1.1814 -0.2471  1.2193  5.3674
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    9.141319   0.149012  61.346 < 2e-16 ***
## score          0.095592   0.002680  35.665 < 2e-16 ***
## tuition       -0.142507   0.068555  -2.079  0.0377 *
## distance      -0.048802   0.010627  -4.592 4.5e-06 ***
## urbanyes       0.022658   0.076409   0.297  0.7668
## distance:urbanyes 0.004735   0.081204   0.058  0.9535
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.58 on 4733 degrees of freedom
## Multiple R-squared:  0.221, Adjusted R-squared:  0.2201
## F-statistic: 268.5 on 5 and 4733 DF, p-value: < 2.2e-16
```

We see from our summary output that our new regression equation looks like this:

$$\widehat{\text{EducationalAttainment}} = 9.141 + .096 \cdot \text{AchievementScore} - .143 \cdot \text{Tuition} - .049 \cdot \text{Distance} + .026 \cdot I(\text{urban} = \text{yes}) + .005 \cdot \text{Distance} \cdot I(\text{urban} = \text{yes})$$

Again, we can split our regression equation into two equations, one for a student going to a high school in an urban area and one for a student going to a high school in a non-urban area. For student's going to school in urban areas:

$$\widehat{\text{EducationalAttainment}} = (9.141 + .026) + .096 \cdot \text{AchievementScore} - .143 \cdot \text{Tuition} - (.049 - .005) \cdot \text{Distance}$$

For student's not going to school in urban areas:

$$\widehat{\text{EducationalAttainment}} = 9.141 + .096 \cdot \text{AchievementScore} - .143 \cdot \text{Tuition} - .049 \cdot \text{Distance}$$

This time we see that the difference between the two equations is not just the intercept, but also the coefficient for the distance a student lives from a four-year college since our urban variable effects both of those locations in the equation.

## The T-test Interpretation

Our summary also tells us the results of the t-test performed for our interaction term (distance:urbanyes). We can see in the table that the t-statistic is .058 and the corresponding p-value is .9535. This means the relationship between a student's number of years of education and distance from a four-year college does not significantly depend on whether that student's high school is located in an urban area. It's important to note that, although in this example we used an interaction between a categorical and a quantitative variable, interaction terms can also be made using two categorical or two quantitative variables.

## F-tests

Another test that our summary output displays is the results of an F-test for the model. The F-test tests the null hypothesis that  $\beta_1 = \beta_2 = \dots = \beta_n = 0$  against the alternative hypothesis that at least one of the regression coefficients,  $\beta$ s, is not equal to 0. We can see from our summary output that the F-statistic for this model is 268.5 and this is compared to an F-table with 5 and 4733 degrees of freedom to obtain a p-value of less than  $2.2 \cdot 10^{-16}$ . This means that at least one of the variables in our model has a significant relationship with our predictor.

## Partial F-tests

A final test we can look at to understand the significance of the variables in our model is the partial F-test. The partial F-test tests the null hypothesis that  $\beta_m = \dots = \beta_p = 0$ , a subset of your regression coefficients is equal to 0, against the alternative hypothesis that at least one  $\beta$  from the subset of  $\beta$ s is not equal to 0. This test allows us to test the significance of a subset of our model variables. It differs from a t-test in that a t-test only allows us to test the significance of one of our variables at a time.

We can use a partial F-test to understand the significance of our interaction term since our interaction term is a subset of our predictor variables. This is essentially the same as a t-test since our subset is only one term in our regression equation. To run the test, we will use the command `anova(complete.model, reduced.model)`, where `complete.model` is our original model and `reduced.model` is our original model without the variables we are interested in testing the significance of.

```
## Analysis of Variance Table
##
## Model 1: education ~ score + tuition + distance + urban
## Model 2: education ~ score + tuition + distance * urban
##   Res.Df  RSS Df Sum of Sq    F Pr(>F)
## 1    4734 11815
## 2    4733 11815   1  0.008489 0.0034 0.9535
```

Our output tells us that our partial F-statistic is .0034 and, from that, we get a p-value of .9535. This p-value is not significant at even a high significance level of .1 so we can conclude that the relationship between the education level that a student achieves and the distance they are from a four-year college does not depend on whether the student's high school is located in an urban area or not.

## Discussion and Conclusions:

Conclude your findings, limitations, and suggest areas for future work.

## APPENDIX

### References

Chetty, R., & Hendren, N. (2016). The Impacts of Neighborhoods on Intergenerational Mobility I: Childhood Exposure Effects. doi: 10.3386/w23001 Chetty, R., Hendren, N., Lin, F., Majerovitz, J., & Scuderi, B. (2016). Childhood Environment and Gender Gaps in Adulthood. doi: 10.3386/w21936 Møllegaard, S., & Jæger, M. M. (2015). The effect of grandparents' economic, cultural, and social capital on grandchildren's educational success. *Research in Social Stratification and Mobility*, 42, 11–19. doi: 10.1016/j.rssm.2015.06.004  
References (to be converted to APA format for project): Relative national intergenerational social mobility: [http://www.ecineq.org/ecineq\\_nyc17/FILESx2017/CR2/p256.pdf](http://www.ecineq.org/ecineq_nyc17/FILESx2017/CR2/p256.pdf) Chetty paper on geography of intergenerational mobility: [http://www.equality-of-opportunity.org/assets/documents/mobility\\_geo.pdf](http://www.equality-of-opportunity.org/assets/documents/mobility_geo.pdf) Obama quote: <https://talkpoverty.org/2015/12/17/american-dream-zip-codes-affordable-housing/> CollegeDistance data info: <http://rdocumentation.org/packages/AER/versions/1.2-7/topics/CollegeDistance> [http://wps.pearsoned.co.uk/wps/media/objects/12401/12699039/empirical/empex\\_tb/CollegeDistance\\_DataDescription.pdf](http://wps.pearsoned.co.uk/wps/media/objects/12401/12699039/empirical/empex_tb/CollegeDistance_DataDescription.pdf)

### Code

```
knitr::opts_chunk$set(echo = TRUE)
library(AER)
library(reshape2)
library(dplyr)
library(ggplot2)
library(psych)
library(skimr)
data(CollegeDistance)
head(CollegeDistance)
names(CollegeDistance)

psych::describe(CollegeDistance)
skimr::skim(CollegeDistance)

plot(CollegeDistance[c(3,10, 11)])

#Continuous Variables

cont_vars <- CollegeDistance %>% dplyr::select(-c(gender, ethnicity, fcollege, mcollege, home, urban, i

melted <- melt(cont_vars)

ggplot(melted, aes(value)) + geom_bar(aes(fill = variable, col = variable), alpha = 0.5, show.legend = F

cor(CollegeDistance[c(3,10, 11)])
score <- lm(score ~ distance + tuition, data = CollegeDistance)
1/(1-summary(score)$r.squared)
plot(CollegeDistance$score, CollegeDistance$education, xlab = "Years of Education", ylab = "Achievement
newscore <- jitter(CollegeDistance$score, factor = 2)
```

```

newed <- jitter(CollegeDistance$education, factor = 2)
plot(newscore, newed, xlab = "Years of Education", ylab = "Achievement Test Score", pch = as.numeric(Co
library(MASS)
starting.model <- lm(education ~ score + urban + distance + tuition, data = CollegeDistance)
simple.model <- lm(education ~ 1, data = CollegeDistance)
stepAIC(starting.model, scope = list(upper = starting.model, lower = simple.model), direction = "backwa

schoolmod <- lm(education ~ score + tuition + distance + urban, data = CollegeDistance)
plot(schoolmod)
summary(schoolmod)
confint(schoolmod)
u <- (9.141+.0256) + .0956*52 - .1426*.900 - .0487*.5
u
nu <- 9.141 + .0956*52 - .1426*.900 - .0487*.5
nu
abs(u-nu)

newdata <- data.frame(score = 52, tuition = .9, distance = .5, urban = 'yes')
newdata
predict(schoolmod, newdata, interval = "predict")
predict(schoolmod, newdata, interval = "confidence")

schoolmod2 <- lm(education ~ score + tuition + distance*urban, data = CollegeDistance)

plot(schoolmod2)
summary(schoolmod2)
anova(schoolmod, schoolmod2)

```