# Level of education and age when the first child born

## Introduction:

In this research I would like to investigate a question: is there a relationship between the level of education of the person and the age when he decide he is ready to have a first child.

It's often been told that modern people tend to have families and children later. But from limited sample of my circle of people I couldn't jump to such a conclusion. Moreover it look's for me that what people do and try to achieve in lives are very important factor in their decision to start a family. Therefore I think that researching this particular question on such a great data would be interesting and valuable.

## Data:

In this research I'm using the modified dataset (all missing values have been recoded as "NA") from General Social Survey (GSS). GSS has been monitoring societal change and studying the growing complexity of American society since 1972. The GSS aims to gather data on contemporary American society in order to monitor and explain trends and constants in attitudes, behaviors, and attributes; to examine the structure and functioning of society in general as well as the role played by relevant subgroups; to compare the United States to other societies in order to place American society in comparative perspective and develop cross-national models of human society; and to make high-quality data easily accessible to scholars, students, policy makers, and others, with minimal cost and waiting.

Units of observation (cases) in this study is a set of answers of respondent (set from each inquired has a numerical identification number) to interview's questions and a year when this respondent was questioned ( numerical 1972-2012).

I'm going to use 2 variables: degree: highest degree achieved (categorical variable with set of values {Lt High School, High School, Junior College, Bachelor, Graduate}) agekdbrn: age when 1st child born ( numerical with values from the interval [9,65])

My study is observational since it is not an experement. In an observational study, researchers collect data in a way that does not directly interfere with how the data arise, which is our case (merely observe). In an experiment, on the other hand, Researchers randomly assign subjects to various treatments, and can therefore establish causal connections between explanatory and response variables (not an our case).
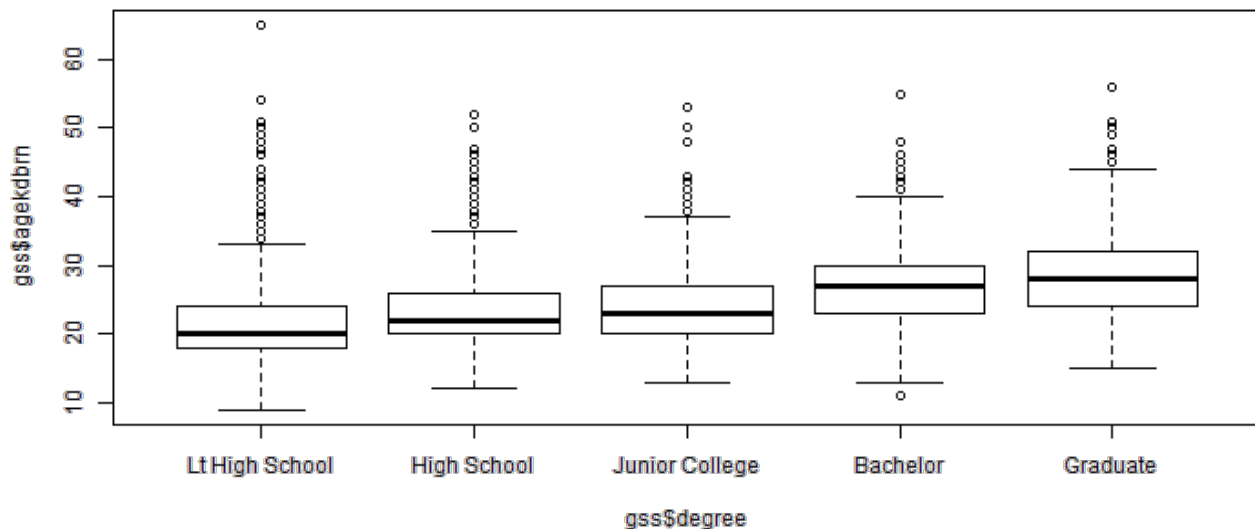
Population of interest is residents of the United States. There are some possible sources of bias: Convenience sample: individuals who are easily accessible are more likely to be included in the survey Non-response: if only a non-rendom fraction of the randomly sampled people respond to a survey such that the sample is no longer representative of the population Voluntary response: occures when the sample consist of people who volunteer to respond because they have strong opinions on the issue. As long as all this 3 possible sourses are avoided data looks reliable and it could be generalized to US population.

Based on an observational study (which we have), we can only establish an association, in

other words correlation, between the explanatory and response variables (NOT a causal connections, since we not have an experement). Because there may be other variables that we didn't control for in this study that contribute to the observed difference.

## Exploratory data analysis:

```
plot(gss$agekdbrn ~ gss$degree)
```



In this plot we can see that it is looks like we have a relationship. People with higher education level prefer to have kids later.

```
by(gss$agekdbrn, gss$degree, summary)
```

```
## gss$degree: Lt High School
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
##      9      18      20      21      24      65    8926
## ------------------------------------------------------------
## gss$degree: High School
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
##     12      20      22      23      26      52   19669
## ------------------------------------------------------------
## gss$degree: Junior College
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
##     13      20      23      24      27      53    1709
## ------------------------------------------------------------
## gss$degree: Bachelor
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
##     11      23      27      27      30      55    5367
## ------------------------------------------------------------
## gss$degree: Graduate
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
##    15.0    24.0    28.0    28.2    32.0    56.0    2430
```

Comparing age when first child born vs education level we can see increase (or rare invariance) of Median, Mean and both quartiles while education level growth. It support my conclusion based on the boxplot that there is a tendency of more educated people prefer have kids later.

# Inference:

So I want to perform a test to prove my hypothesis that there are difference between education level of a person and time when he\she decided to have a baby. I divide our sample into 5 groups according to theier education level. So my data is presented as one numerical and one categorical variable (with 5 levels). That means that I only can perform a hypothesis test. ANOVA and pairwise tests are suitable for a such situation.

At first I have to make a null hypothesis which always state that there are nothing going on and that on average age when first child born are the same for all education levels: H0: Mu1=..=Mu5. Alternative hypothesis: Ha: at least one pair of means(as a measure of average value) are different from each other.

To perform such test we have to check conditions.

1.Independence: a) within group independence: This means that the sample observations must be independent of each other in every group.

```
data = na.omit(subset(gss, select = c(agekdbrn, degree)))
a = subset(data, degree == "Lt High School")
b = subset(data, degree == "High School")
c = subset(data, degree == "Junior College")
d = subset(data, degree == "Bachelor")
e = subset(data, degree == "Graduate")
dim(a)
```

```
## [1] 2896    2
```

```
dim(b)
```

```
## [1] 9618    2
```

```
dim(c)
```

```
## [1] 1361    2
```

```
dim(d)
```

```
## [1] 2635    2
```
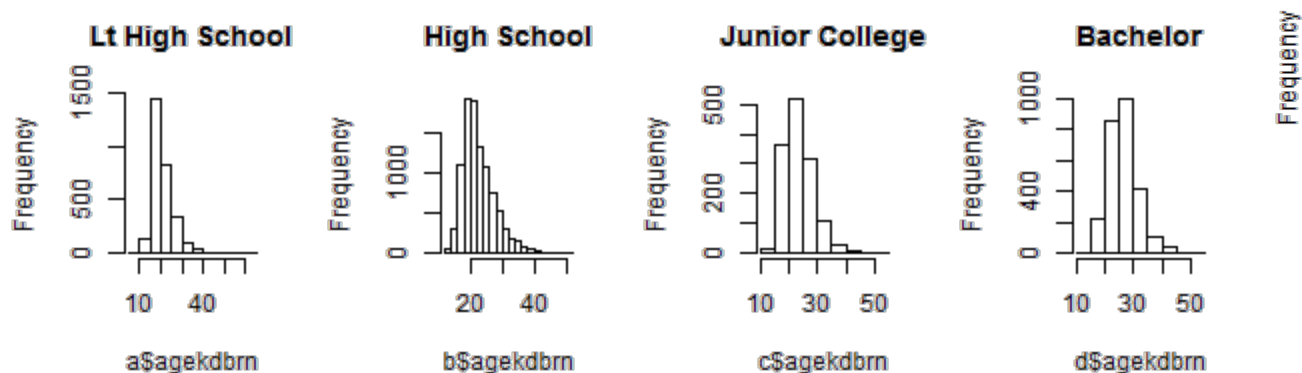
```
dim(e)
```

```
## [1] 1440    2
```

Since we have a random sample of people and numbers of people in each groups (as shown above) are less than 10% of their respective population we can conclude that this condition is met.
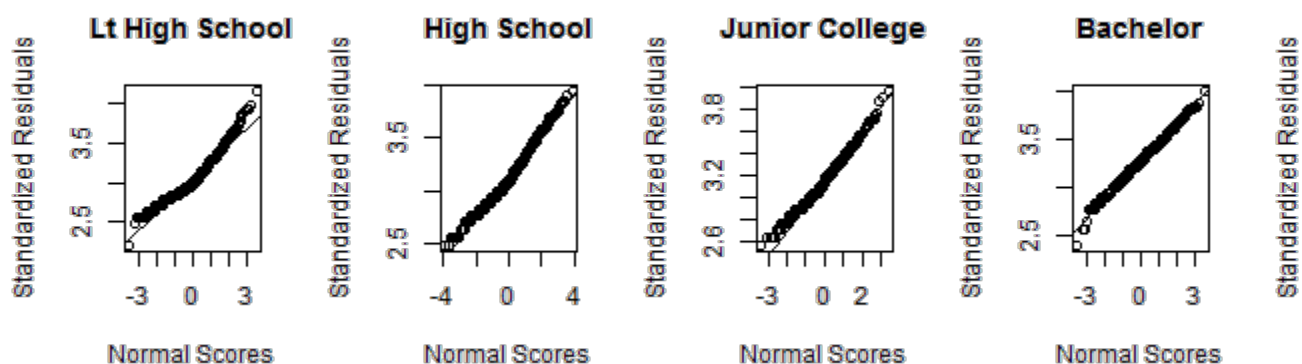
b) between group independence: This condition is also satisfied since every person could be included in only one group, so there is no pairing between the groups.

2.Approximately normal condition. The distribution of the response variable within each group should be approximately normal.

**Graduate**

From the graphics above we can see that it's not the case and our distribution are skewed. Hence we have to take logarithm to "normalize data."



Now our data set is sutible and condition of "approximate normality" is met.

3.The last condition is the condition of constant variance, which says that variability should be consistent across groups. We can check this condition using side-by-side box plots, where it's going to be pretty easy visually to compare the variability across the groups, as well as by looking at the standard deviations for each group (log): 0.218 for "Lt High School";
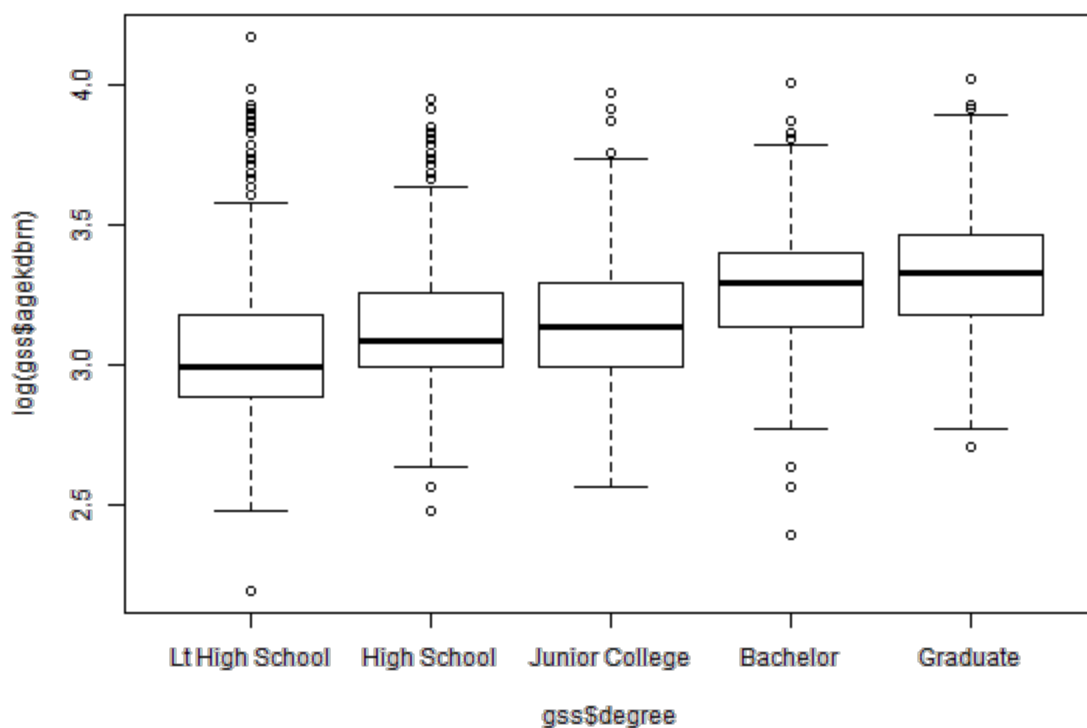
0.1994 for "High School";

0.2052 for "Junior College";

0.1862 for "Bachelor";

0.2007 for "Graduate";

From that we can see that we have roughly equal variability across the groups and that our third condition is satisfied. But I have to mension that my previous hypothesis changed: H0: ln(Mu1)=..=ln(Mu5). Alternative hypothesis: Ha: at least one pair of means(as a measure of average value) are different from each other.
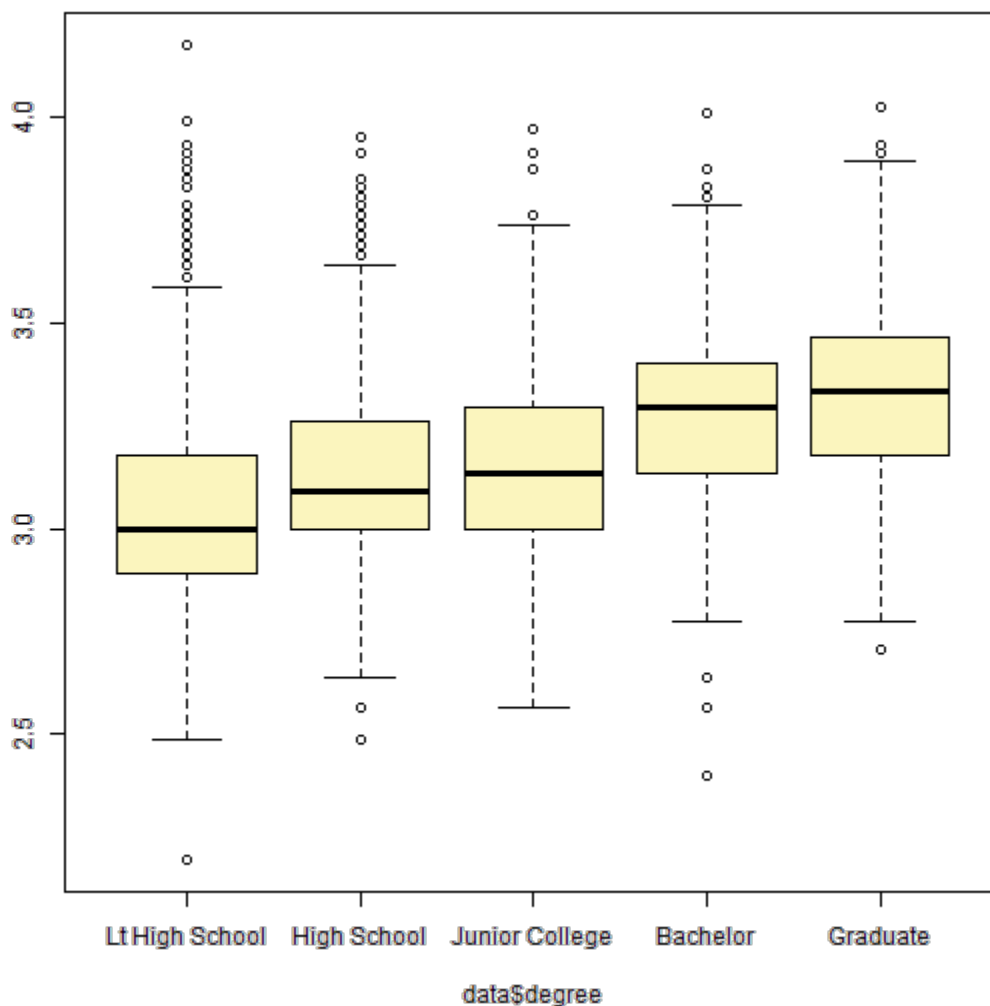
```
plot(log(gss$agekdbrn) ~ gss$degree)
```

```
load(url("http://bit.ly/dasi_gss_ws_cl"))
source("http://bit.ly/dasi_inference")
inference(y = log(data$agekdbrn), x = data$degree, est =
"mean", type = "ht",
    alternative = "greater", method = "theoretical", inf_plot
= FALSE, sum_stats = TRUE)
```

```
## Warning: package 'BHH2' was built under R version 3.0.3
```

```
## Response variable: numerical, Explanatory variable:
categorical
## ANOVA
## Summary statistics:
## n_Lt High School = 2896, mean_Lt High School = 3.034, sd_Lt
High School = 0.218
## n_High School = 9618, mean_High School = 3.116, sd_High
School = 0.1994
## n_Junior College = 1361, mean_Junior College = 3.155,
sd_Junior College = 0.2052
## n_Bachelor = 2635, mean_Bachelor = 3.275, sd_Bachelor =
0.1862
## n_Graduate = 1440, mean_Graduate = 3.319, sd_Graduate =
0.2007
```

```
## H_0: All means are equal.
## H_A: At least one mean is different.
## Analysis of Variance Table
##
## Response: y
##              Df Sum Sq Mean Sq F value Pr(>F)
## x             4    131    32.8     810 <2e-16
## Residuals 17945    726     0.0
##
## Pairwise tests: t tests with pooled SD
##                Lt High School High School Junior College
Bachelor
## High School              0          NA
NA        NA
## Junior College           0           0
NA        NA
## Bachelor                 0           0
0         NA
## Graduate                 0           0
0         0
```

Performing this test we get p-value < 2e-16, which coresponds to probability of at least as large differences between the groups variabilities if, in fact, the means of all groups are equal. Since p-value < significance level (0.05) we can reject Ho and said that at least one log(mean) is different. Pairwise tests shows that all p-values for different paires are insignificant(->0) and less then modified significance level (0.05/10=0.005). Which states that neither pair of log(mean) are equal.

# Conclusion:

Despite the fact that we was working with logarithm, we can generalize our conclusions to means itself. So ANOVA test states that my initial guess was right and that difference in anerage age when people get their first child correlate with their level of education. Furthermore this association is strong between every education level groups. Unfortunately I doesn't see the way to perform an experement in this field to find out is there a cousational relationship. But still I get an interesting result which could be use to perform a better social programs for young families etc. Maybe it could be helpful to consider some so needed reforms (For example in USA parental leave is extremely short, approximately 3 months. In my country it is 3 years!) Anyway, it was a fun research with great result which could be beginning poin for a bunch of new studies) Thank you for your time!

# References:

1. Smith, Tom W., Michael Hout, and Peter V. Marsden. General Social Survey, 1972-2012 [Cumulative File]. ICPSR34802-v1. Storrs, CT: Roper Center for Public Opinion Research, University of Connecticut /Ann Arbor, MI: Inter-university Consortium for Political and Social Research [distributors], 2013-09-11. doi:10.3886/ICPSR34802.v1 (https://d396qusza40orc.cloudfront.net

/statistics%2Fproject%2Fgss1.html)

2. https://d396qusza40orc.cloudfront.net/statistics/lab_resources/Rcommands.pdf

3. http://www.rstudio.com/ide/docs/authoring/using_markdown?version=0.98.501&mode=desktop

4. https://class.coursera.org/statistics-001

## Appendix:

```
head(data, 55)
```

```
##       agekdbrn         degree
## 29389      21     High School
## 29391      25 Lt High School
## 29392      23 Lt High School
## 29395      17 Lt High School
## 29396      23     High School
## 29397      17     High School
## 29398      17     High School
## 29401      29        Bachelor
## 29402      32        Bachelor
## 29406      34        Bachelor
## 29410      35        Graduate
## 29415      32     High School
## 29416      14 Lt High School
## 29418      23        Bachelor
## 29419      20     High School
## 29421      20     High School
## 29423      27        Bachelor
## 29429      38        Bachelor
## 29431      20        Graduate
## 29432      29     High School
## 29435      18     High School
## 29436      25        Bachelor
## 29437      19 Lt High School
## 29438      19     High School
## 29439      18     High School
## 29444      26        Bachelor
## 29446      27     High School
## 29448      23 Lt High School
## 29449      24     High School
## 29451      20     High School
## 29452      24     High School
## 29458      21 Lt High School
## 29459      19 Lt High School
## 29460      21     High School
## 29461      32 Lt High School
## 29462      24     High School
## 29466      18        Bachelor
## 29467      32        Graduate
## 29469      24        Graduate
## 29474      26     High School
## 29475      28     High School
## 29476      29     High School
## 29478      30 Lt High School
## 29479      31     High School
## 29480      21 Lt High School
## 29482      32        Graduate
## 29484      29        Graduate
## 29487      18     High School
```

```
## 29490          18 Junior College
## 29492          23    High School
## 29494          32 Lt High School
## 29495          25    High School
## 29500          27       Graduate
## 29504          26 Lt High School
## 29505          28 Lt High School
```