

Statistical Inference with the GSS Data Using R

Author: Arash Sadeghzadeh

Date: 30-Dec-2022

Introduction

Since 1972, the General Social Survey (GSS) has been monitoring societal change and studying the growing complexity of American society. The GSS aims to gather data on contemporary American society in order to monitor and explain trends and constants in attitudes, behaviors, and attributes. In this project, we consider a few research questions to answer based on these data. We perform inference that addresses the research questions using R.

Load packages

```
library(ggplot2)
library(dplyr)
library(statsr)
library(tidyr)
library(ggmosaic)
```

Load data

```
load("gss.Rdata")
```

Lets have a look at the summary of the data frame:

```
str(gss)

## 'data.frame':   57061 obs. of  114 variables:
##  $ caseid   : int   1 2 3 4 5 6 7 8 9 10 ...
##  $ year     : int  1972 1972 1972 1972 1972 1972 1972 1972 1972 1972 ...
##  $ age      : int   23 70 48 27 61 26 28 27 21 30 ...
##  $ sex      : Factor w/ 2 levels "Male","Female": 2 1 2 2 2 1 1 1 2 2 ...
##  $ race     : Factor w/ 3 levels "White","Black",...: 1 1 1 1 1 1 1 1 2 2 ...
##  $ hispanic : Factor w/ 28 levels "Not Hispanic",...: NA NA NA NA NA NA NA NA NA NA ...
##  $ uscitzn  : Factor w/ 4 levels "A U.S. Citizen",...: NA NA NA NA NA NA NA NA NA NA ...
##  $ educ     : int   16 10 12 17 12 14 13 16 12 12 ...
##  $ paeduc   : int   10 8 8 16 8 18 16 16 12 10 ...
##  $ maeduc   : int   NA 8 8 12 8 19 12 14 12 7 ...
##  $ speduc   : int   NA 12 11 20 12 NA NA NA NA 11 ...
##  $ degree   : Factor w/ 5 levels "Lt High School",...: 4 1 2 4 2 2 2 4 2 2 ...
##  $ vetyears : Factor w/ 5 levels "None","Less Than 2 Yrs",...: NA NA NA NA NA NA NA NA NA NA ...
##  $ sei      : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ wrkstat  : Factor w/ 8 levels "Working Fulltime",...: 1 5 2 1 7 1 1 1 2 1 ...
##  $ wrkslf   : Factor w/ 2 levels "Self-Employed",...: 2 2 2 2 2 2 2 2 2 2 ...
##  $ marital  : Factor w/ 5 levels "Married","Widowed",...: 5 1 1 1 1 5 3 5 5 1 ...
##  $ spwrksta : Factor w/ 8 levels "Working Fulltime",...: NA 7 1 1 3 NA NA NA NA 1 ...
##  $ sibs     : int    3 4 5 5 2 1 7 1 2 7 ...
```

```

## $ childs : int 0 5 4 0 2 0 2 0 2 4 ...
## $ agekdbn: int NA NA NA NA NA NA NA NA NA NA ...
## $ incom16 : Factor w/ 6 levels "Far Below Average",...: 3 4 3 3 2 3 4 3 3 1 ...
## $ born : Factor w/ 2 levels "Yes","No": NA NA NA NA NA NA NA NA NA NA ...
## $ parborn : Factor w/ 9 levels "Both In U.S",...: NA NA NA NA NA NA NA NA NA NA ...
## $ granborn: int NA NA NA NA NA NA NA NA NA NA ...
## $ income06: Factor w/ 26 levels "Under $1 000",...: NA NA NA NA NA NA NA NA NA NA ...
## $ coninc : int 25926 33333 33333 41667 69444 60185 50926 18519 3704 25926 ...
## $ region : Factor w/ 9 levels "New England",...: 3 3 3 3 3 3 3 7 7 ...
## $ partyid : Factor w/ 8 levels "Strong Democrat",...: 3 2 4 2 1 3 3 3 1 1 ...
## $ polviews: Factor w/ 7 levels "Extremely Liberal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ relig : Factor w/ 13 levels "Protestant","Catholic",...: 3 2 1 5 1 1 2 3 1 1 ...
## $ attend : Factor w/ 9 levels "Never","Lt Once A Year",...: 3 8 5 NA NA 3 8 NA 4 9 ...
## $ natspac : Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ natenvir: Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ natheal : Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ natcity : Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ natcrime: Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ natdrug : Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ nateduc : Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ natrace : Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ natarms : Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ nataid : Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ natfare : Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ natroad : Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ natsoc : Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ natmass : Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ natpark : Factor w/ 3 levels "Too Little","About Right",...: NA NA NA NA NA NA NA NA NA NA ...
## $ confinan: Factor w/ 3 levels "A Great Deal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ conbus : Factor w/ 3 levels "A Great Deal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ conclerg: Factor w/ 3 levels "A Great Deal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ coneduc : Factor w/ 3 levels "A Great Deal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ confed : Factor w/ 3 levels "A Great Deal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ conlabor: Factor w/ 3 levels "A Great Deal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ conpress: Factor w/ 3 levels "A Great Deal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ conmedic: Factor w/ 3 levels "A Great Deal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ contv : Factor w/ 3 levels "A Great Deal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ conjudge: Factor w/ 3 levels "A Great Deal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ consci : Factor w/ 3 levels "A Great Deal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ conlegis: Factor w/ 3 levels "A Great Deal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ conarmy : Factor w/ 3 levels "A Great Deal",...: NA NA NA NA NA NA NA NA NA NA ...
## $ joblose : Factor w/ 5 levels "Very Likely",...: NA NA NA NA NA NA NA NA NA NA ...
## $ jobfind : Factor w/ 3 levels "Very Easy","Somewhat Easy",...: NA NA NA NA NA NA NA NA NA NA ...
## $ satjob : Factor w/ 4 levels "Very Satisfied",...: 3 NA 2 1 NA 2 1 3 2 2 ...
## $ richwork: Factor w/ 2 levels "Continue Working",...: NA NA NA NA NA NA NA NA NA NA ...
## $ jobinc : Factor w/ 5 levels "Most Impt","Second",...: NA NA NA NA NA NA NA NA NA NA ...
## $ jobsec : Factor w/ 5 levels "Most Impt","Second",...: NA NA NA NA NA NA NA NA NA NA ...
## $ jobhour : Factor w/ 5 levels "Most Impt","Second",...: NA NA NA NA NA NA NA NA NA NA ...
## $ jobpromo: Factor w/ 5 levels "Most Impt","Second",...: NA NA NA NA NA NA NA NA NA NA ...
## $ jobmeans: Factor w/ 5 levels "Most Impt","Second",...: NA NA NA NA NA NA NA NA NA NA ...
## $ class : Factor w/ 5 levels "Lower Class",...: 3 3 2 3 2 3 3 2 2 2 ...
## $ rank : int NA NA NA NA NA NA NA NA NA NA ...
## $ satfin : Factor w/ 3 levels "Satisfied","More Or Less",...: 3 2 1 3 1 2 2 3 2 3 ...
## $ finalter: Factor w/ 3 levels "Better","Worse",...: 1 3 1 3 1 1 1 1 2 3 ...

```

```
## $ finrela : Factor w/ 5 levels "Far Below Average",...: 3 4 3 3 4 4 4 3 3 2 ...
## $ unemp   : Factor w/ 2 levels "Yes","No": NA NA NA NA NA NA NA NA NA NA ...
## $ govaid  : Factor w/ 2 levels "Yes","No": NA NA NA NA NA NA NA NA NA NA ...
## $ getaid  : Factor w/ 2 levels "Yes","No": NA NA NA NA NA NA NA NA NA NA ...
## $ union   : Factor w/ 4 levels "R Belongs","Spouse Belongs",...: NA NA NA NA NA NA NA NA NA NA ...
## $ getahead: Factor w/ 4 levels "Hard Work","Both Equally",...: NA NA NA NA NA NA NA NA NA NA ...
## $ parsol  : Factor w/ 5 levels "Much Better",...: NA NA NA NA NA NA NA NA NA NA ...
## $ kidssol : Factor w/ 6 levels "Much Better",...: NA NA NA NA NA NA NA NA NA NA ...
## $ abdefect: Factor w/ 2 levels "Yes","No": 1 1 1 2 1 1 1 1 1 1 ...
## $ abnomore: Factor w/ 2 levels "Yes","No": 1 2 1 2 1 1 2 1 2 2 ...
## $ abhlth  : Factor w/ 2 levels "Yes","No": 1 1 1 1 1 1 1 1 1 1 ...
## $ abpoor  : Factor w/ 2 levels "Yes","No": 1 2 1 1 1 1 2 1 2 1 ...
## $ abrape  : Factor w/ 2 levels "Yes","No": 1 1 1 1 1 1 1 1 NA 1 ...
## $ absingle: Factor w/ 2 levels "Yes","No": 1 1 1 1 1 1 1 1 2 2 ...
## $ abany   : Factor w/ 2 levels "Yes","No": NA NA NA NA NA NA NA NA NA NA ...
## $ pillok  : Factor w/ 4 levels "Strongly Agree",...: NA NA NA NA NA NA NA NA NA NA ...
## $ sexeduc : Factor w/ 3 levels "Favor","Oppose",...: NA NA NA NA NA NA NA NA NA NA ...
## $ divlaw  : Factor w/ 3 levels "Easier","More Difficult",...: NA NA NA NA NA NA NA NA NA NA ...
## $ premarsx: Factor w/ 5 levels "Always Wrong",...: 4 1 1 1 3 3 4 3 4 1 ...
## $ teensex : Factor w/ 5 levels "Always Wrong",...: NA NA NA NA NA NA NA NA NA NA ...
## $ xmarsex : Factor w/ 5 levels "Always Wrong",...: NA NA NA NA NA NA NA NA NA NA ...
## $ homosex : Factor w/ 5 levels "Always Wrong",...: NA NA NA NA NA NA NA NA NA NA ...
## $ suicide1: Factor w/ 2 levels "Yes","No": NA NA NA NA NA NA NA NA NA NA ...
## $ suicide2: Factor w/ 2 levels "Yes","No": NA NA NA NA NA NA NA NA NA NA ...
## $ suicide3: Factor w/ 2 levels "Yes","No": NA NA NA NA NA NA NA NA NA NA ...
## $ suicide4: Factor w/ 2 levels "Yes","No": NA NA NA NA NA NA NA NA NA NA ...
## [list output truncated]
```

The numbers of the data in the data frame are as follows:

```
nrow(gss)
```

```
## [1] 57061
```

Research questions

Is there any difference between the average income of women and men?

Note that the total family income and the sex have been reported in columns “coninc” and “sex”, respectively. Let us first compute the average income for men and women:

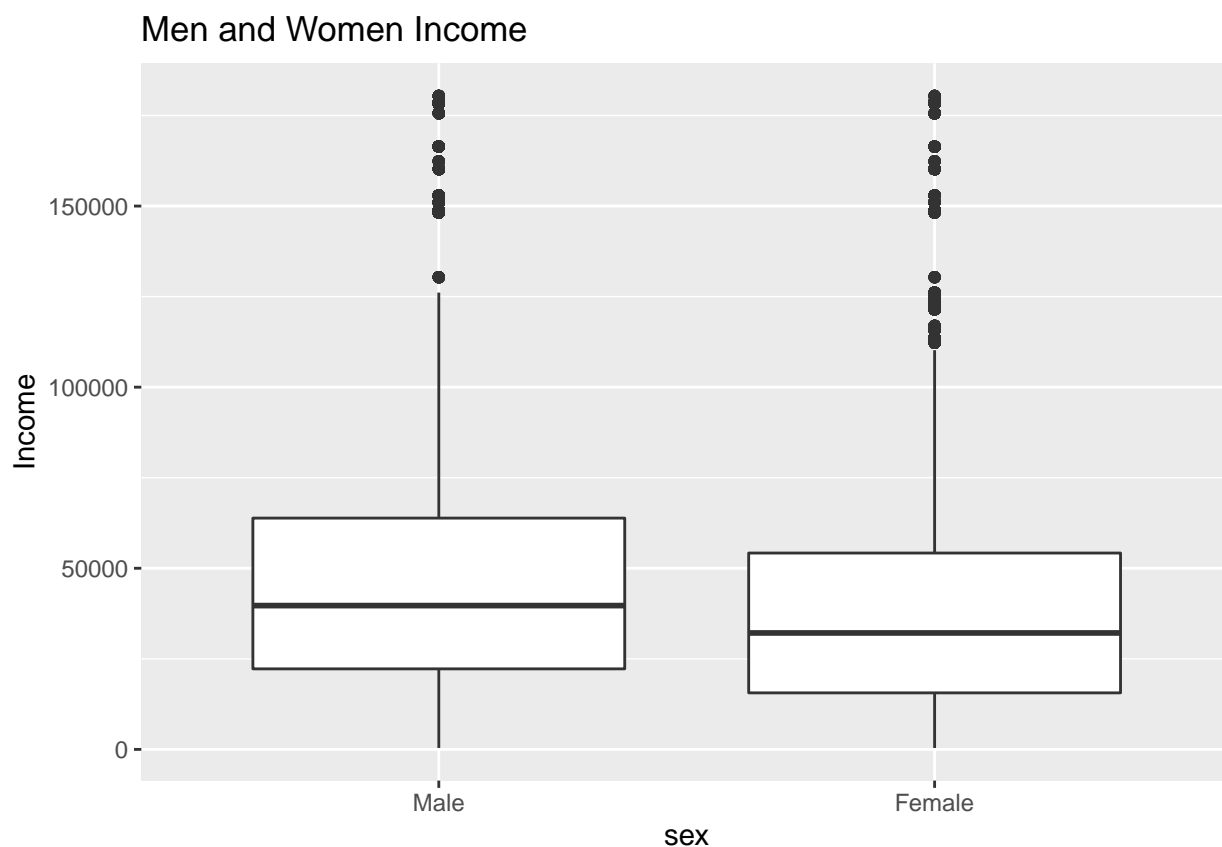
```
gss %>% group_by(sex) %>% summarize(Average.Income = round(mean(coninc, na.rm=TRUE)))
```

```
## # A tibble: 2 x 2
##   sex      Average.Income
##   <fct>          <dbl>
## 1 Male           48764
## 2 Female         41020
```

Lets have a look at the related boxplot of men and women income:

```
ggplot(data=gss, aes(x=sex, y=coninc)) + geom_boxplot() + ylab("Income") +
ggtitle("Men and Women Income")
```

```
## Warning: Removed 5829 rows containing non-finite values (stat_boxplot).
```



Now, let's compute the confidence interval for these averages. We first compute the 95% confidence interval for the average of the men's income:

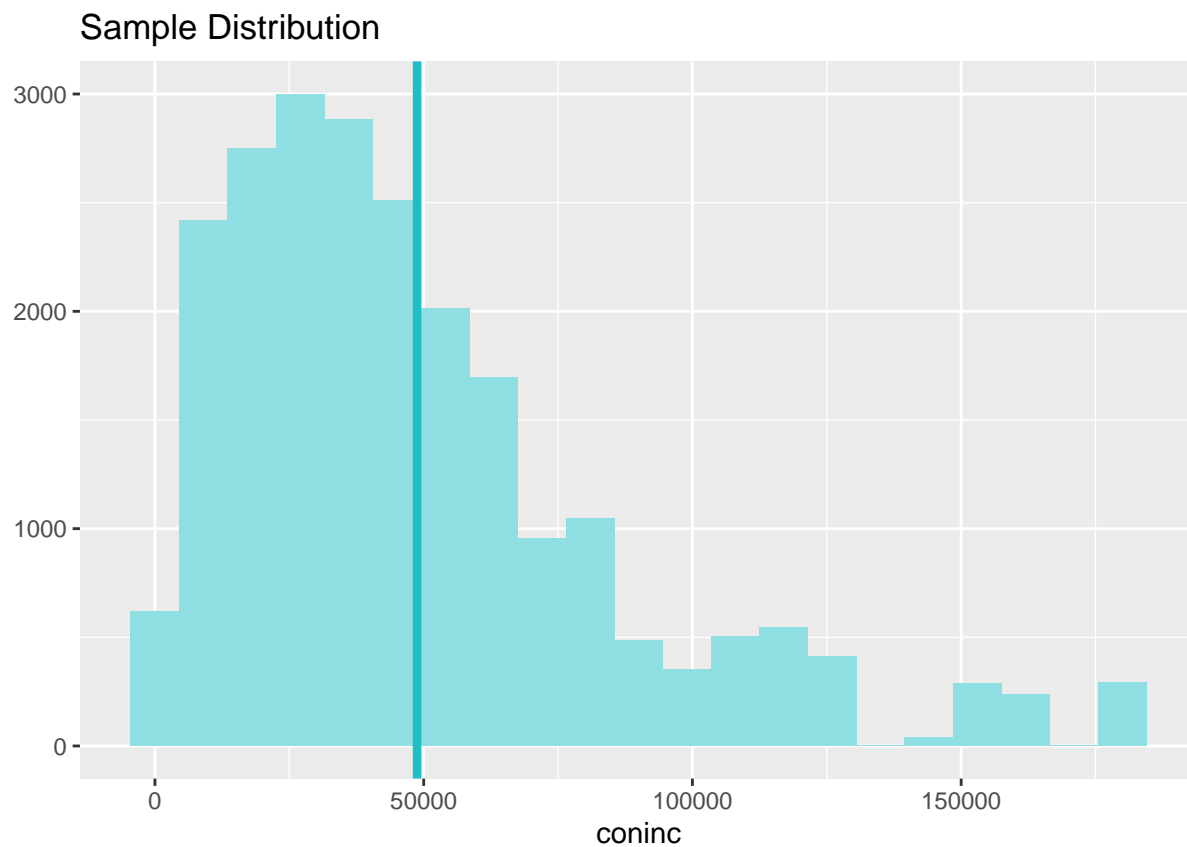
```
gss_men <- gss %>% filter(sex=="Male")
```

```
gss_men %>% summarize(Number.of.Samples.Men = n())
```

```
##   Number.of.Samples.Men
## 1                25146
```

```
inference(y=coninc, data=gss_men, type="ci", statistic="mean", conf_level=0.95,
          method="theoretical")
```

```
## Single numerical variable
## n = 23043, y-bar = 48763.6453, s = 36916.3394
## 95% CI: (48286.9729 , 49240.3177)
```



Now, let us compute the related 95% confidence interval for the average income for women as follows:

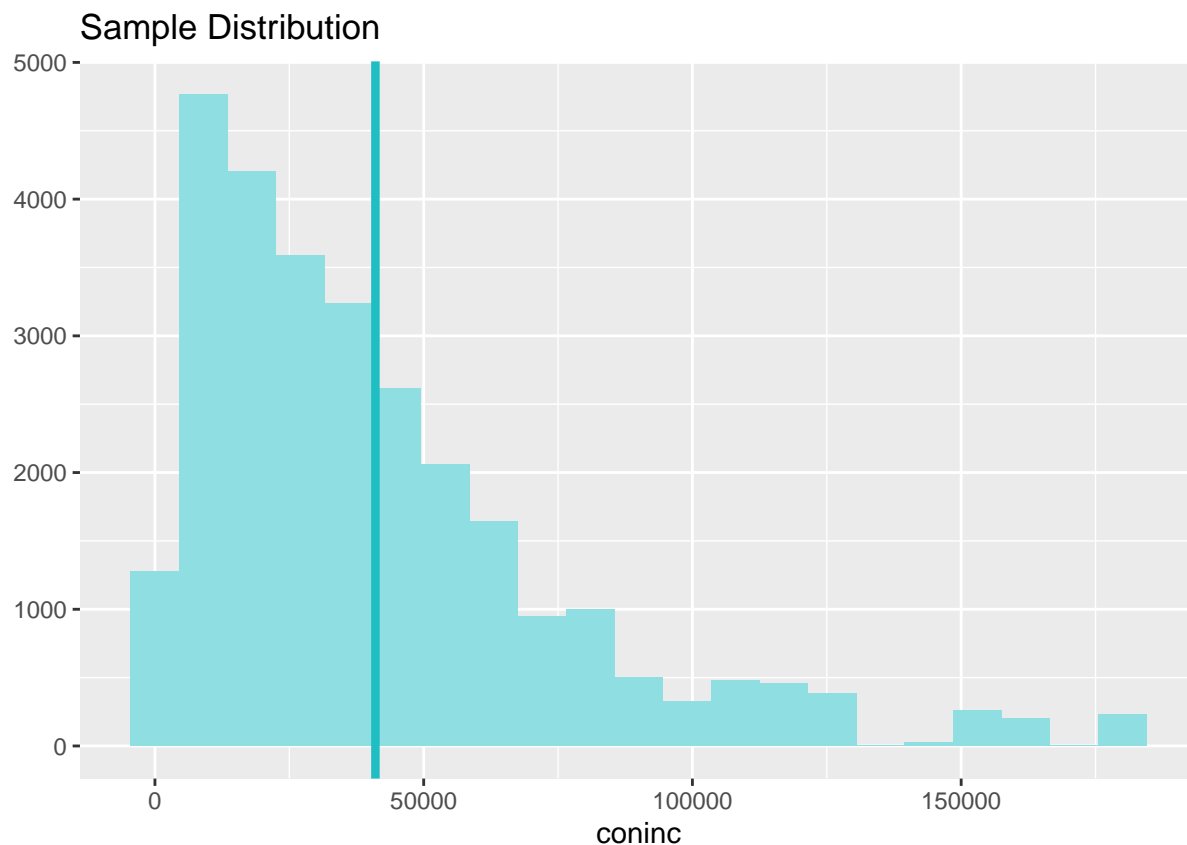
```
gss_women <- gss %>% filter(sex=="Female")
```

```
gss_women %>% summarize(Number.of.Samples.Women = n())
```

```
##   Number.of.Samples.Women
## 1                      31915
```

```
inference(y=coninc, data=gss_women, type="ci", statistic="mean", conf_level=0.95,
          method="theoretical")
```

```
## Single numerical variable
## n = 28189, y-bar = 41020.2199, s = 34728.8358
## 95% CI: (40614.7888 , 41425.6511)
```



Now, let us evaluate whether these data support the hypothesis that men on average receive more salary than women?

So, we have the following hypotheses:

H_0 (Null hypothesis): Men and women on average receive the same salary.

H_A (Alternative hypothesis): Men receive more salary than women on average.

```
df_withoutNA <- gss %>% drop_na(coninc)
inference(data=df_withoutNA, y=coninc, x=sex, type="ht", statistic="mean",
          method="theoretical", alternative="greater", sig_level=0.95)
```

```
## Warning: Missing null value, set to 0
```

```
## Response variable: numerical
```

```
## Explanatory variable: categorical (2 levels)
```

```
## n_Male = 23043, y_bar_Male = 48763.6453, s_Male = 36916.3394
```

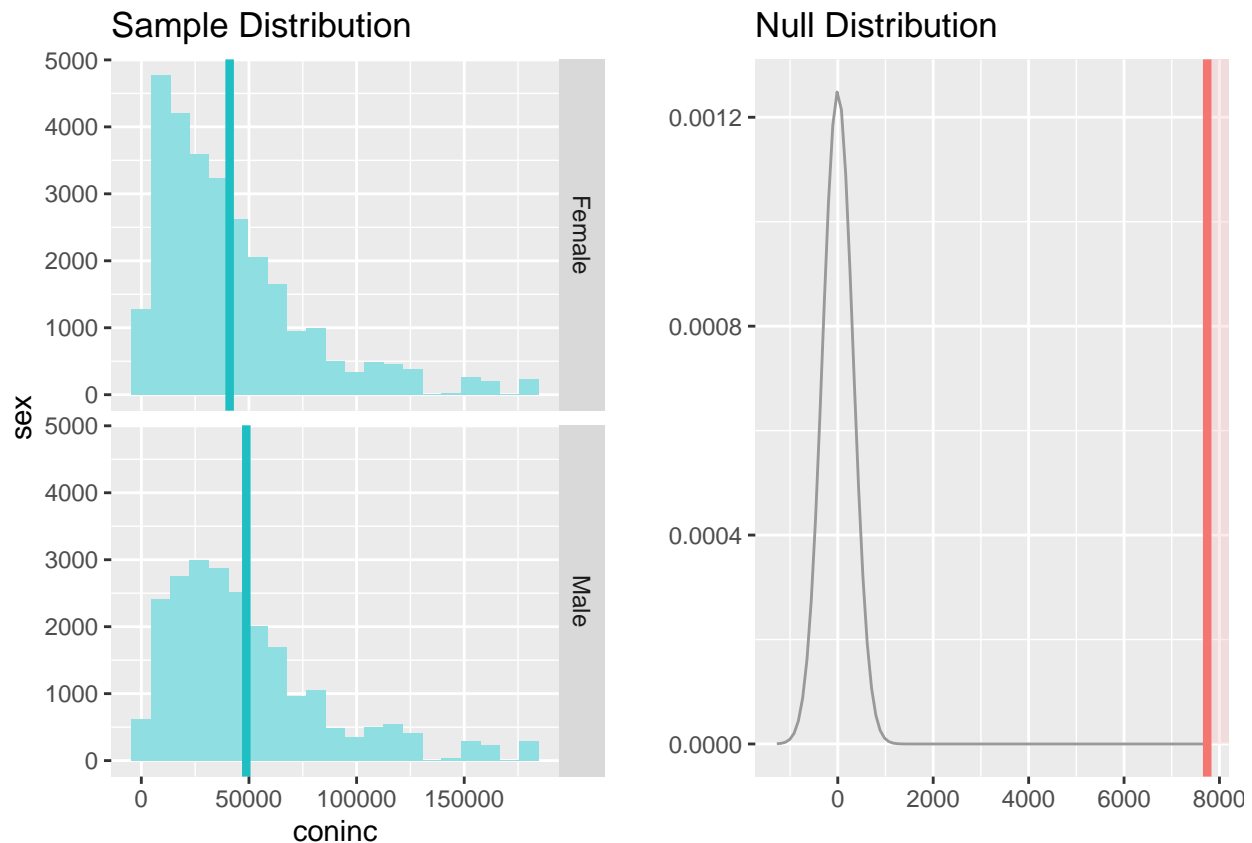
```
## n_Female = 28189, y_bar_Female = 41020.2199, s_Female = 34728.8358
```

```
## H0: mu_Male = mu_Female
```

```
## HA: mu_Male > mu_Female
```

```
## t = 24.2541, df = 23042
```

```
## p_value = < 0.0001
```



Since p-value is small (less than 0.05), the data provide convincing evidence that men on average receive more salary than women.

Now, just for the further evaluation, we compute the difference between the average salaries of 1000 random samples of men and 1000 random samples of women. To obtain the distribution of this difference, we consider the salary difference computation for 5000 different set of samples of men and women as follows:

```
set.seed(1979)
sample_salary_men <- gss %>% filter(sex=="Male") %>%
  rep_sample_n(size=1000, reps=5000, replace=TRUE) %>%
  summarize(mean.salary = mean(coninc, na.rm=TRUE))

sample_salary_women <- gss %>% filter(sex=="Female") %>%
  rep_sample_n(size=1000, reps=5000, replace=TRUE) %>%
  summarize(mean.salary = mean(coninc, na.rm=TRUE))

diff_salary <- sample_salary_men - sample_salary_women
ggplot(data=diff_salary, aes(x=mean.salary)) + geom_histogram(binwidth=300)
```

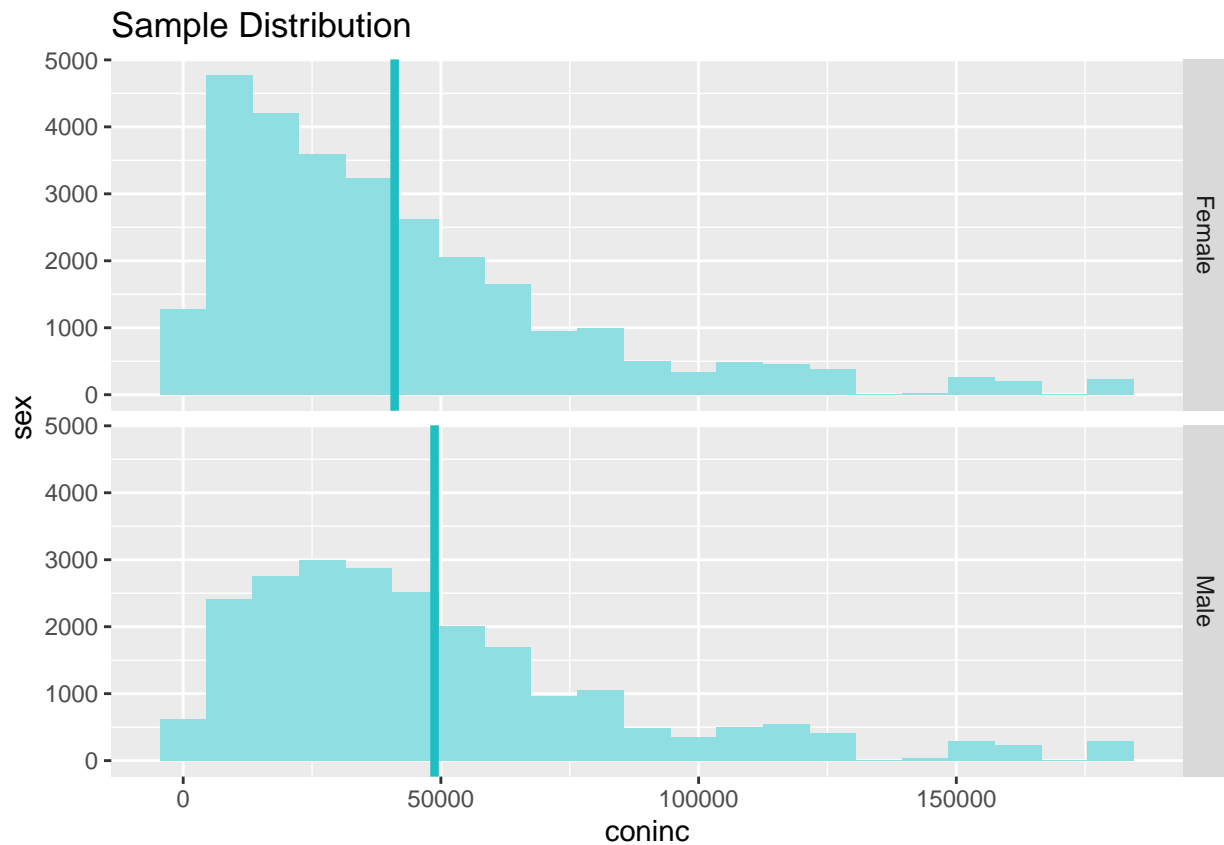


Not surprisingly, the salary difference has a nearly normal distribution.

Now, we compute the 95% confidence interval for the difference between the average salary for men and women as follows:

```
df_withoutNA <- gss %>% drop_na(coninc)
inference(data=df_withoutNA, y=coninc, x=sex, type="ci", statistic="mean",
          method="theoretical", sig_level=0.95)
```

```
## Response variable: numerical, Explanatory variable: categorical (2 levels)
## n_Male = 23043, y_bar_Male = 48763.6453, s_Male = 36916.3394
## n_Female = 28189, y_bar_Female = 41020.2199, s_Female = 34728.8358
## 95% CI (Male - Female): (7117.6505 , 8369.2002)
```

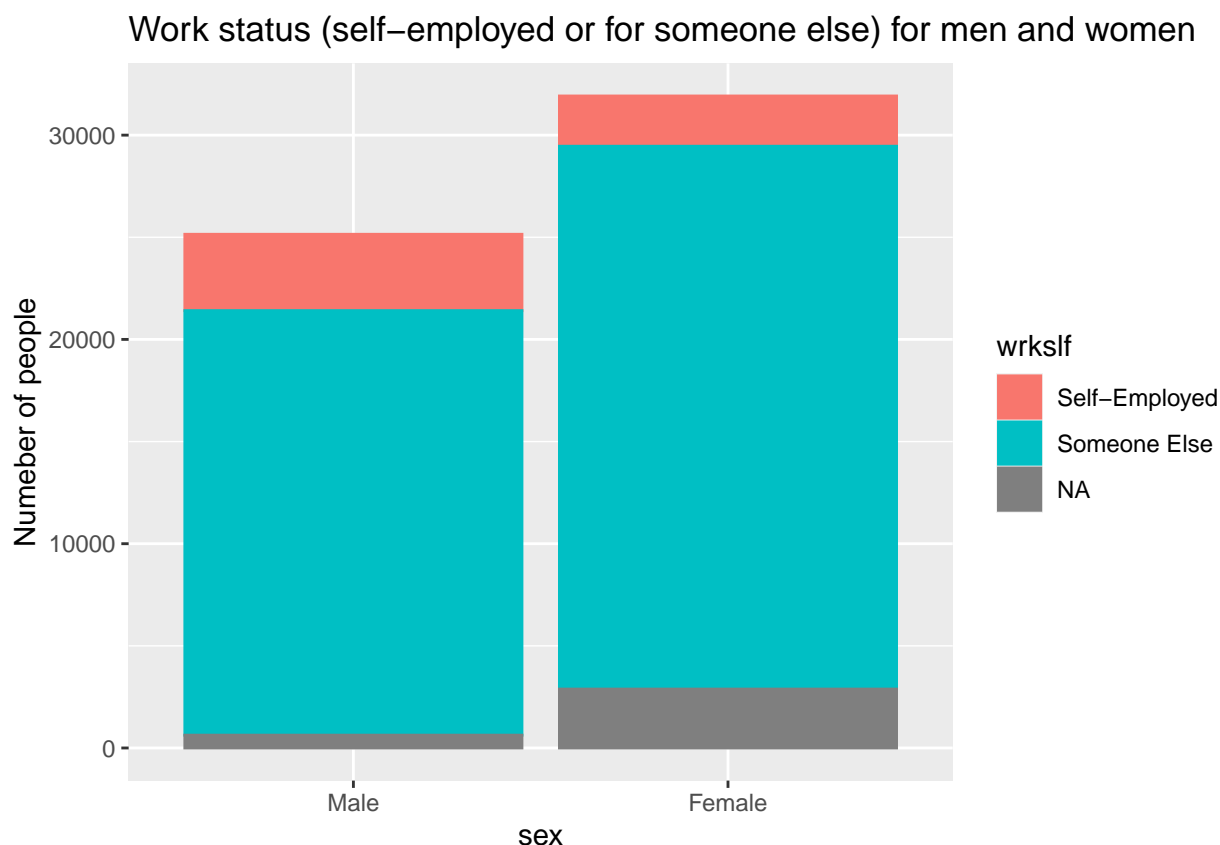



Therefore, we are 95% confident that the average salary for men is about 7117 up to 8369 more than the average salary for women.

Estimate how men and women at large compare with respect to being self-employed?

Note that “wrkslf” column contains the data regarding being self-employed or employed by someone else. In what follows, we illustrate the related barplots:

```
c4 = c("A", "B", "C")
ggplot(data=gss, aes(x=sex, color=wrkslf, fill=wrkslf)) + geom_bar() +
  ylab("Number of people") +
  ggtitle("Work status (self-employed or for someone else) for men and women")
```



Lets compute the proportionals of being self-employed for men and women as follows:

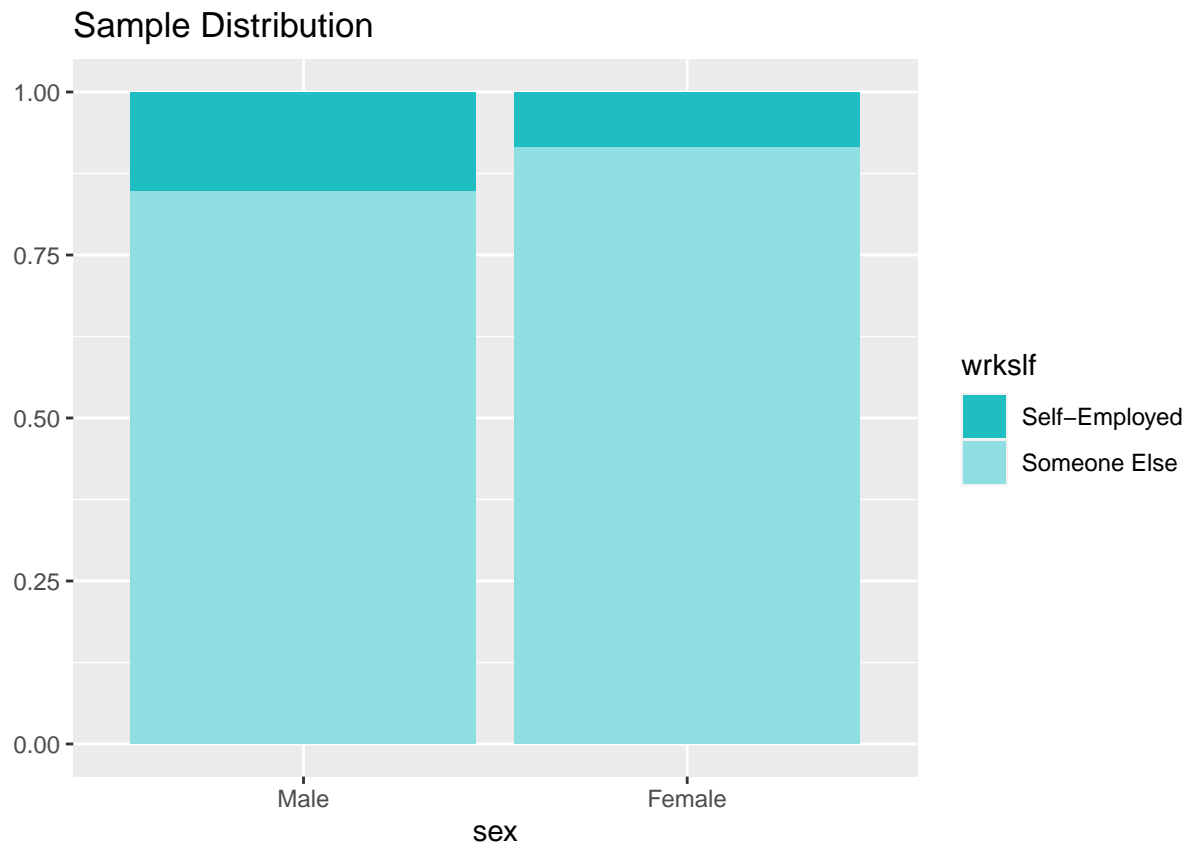
```
gss %>% drop_na(wrkslf) %>%
  mutate(self.employed=ifelse(wrkslf=="Self-Employed", 1, 0)) %>%
  group_by(sex) %>%
  summarize(prop.self.employed = sum(self.employed)/n(),
            num.self.employed = sum(self.employed),
            total.num = n())
```

```
## # A tibble: 2 x 4
##   sex    prop.self.employed num.self.employed total.num
##   <fct>          <dbl>          <dbl>      <int>
## 1 Male             0.152             3738      24518
## 2 Female           0.0847            2459      29031
```

Now, we evaluate the related 95% confidence interval for the difference of being self-employed between men and women:

```
inference(data=gss, y=wrkslf, x=sex, statistic="proportion",
          success="Self-Employed", method="theoretical", type="ci",
          sig_level=0.95)
```

```
## Response variable: categorical (2 levels, success: Self-Employed)
## Explanatory variable: categorical (2 levels)
## n_Male = 24518, p_hat_Male = 0.1525
## n_Female = 29031, p_hat_Female = 0.0847
## 95% CI (Male - Female): (0.0622 , 0.0733)
```



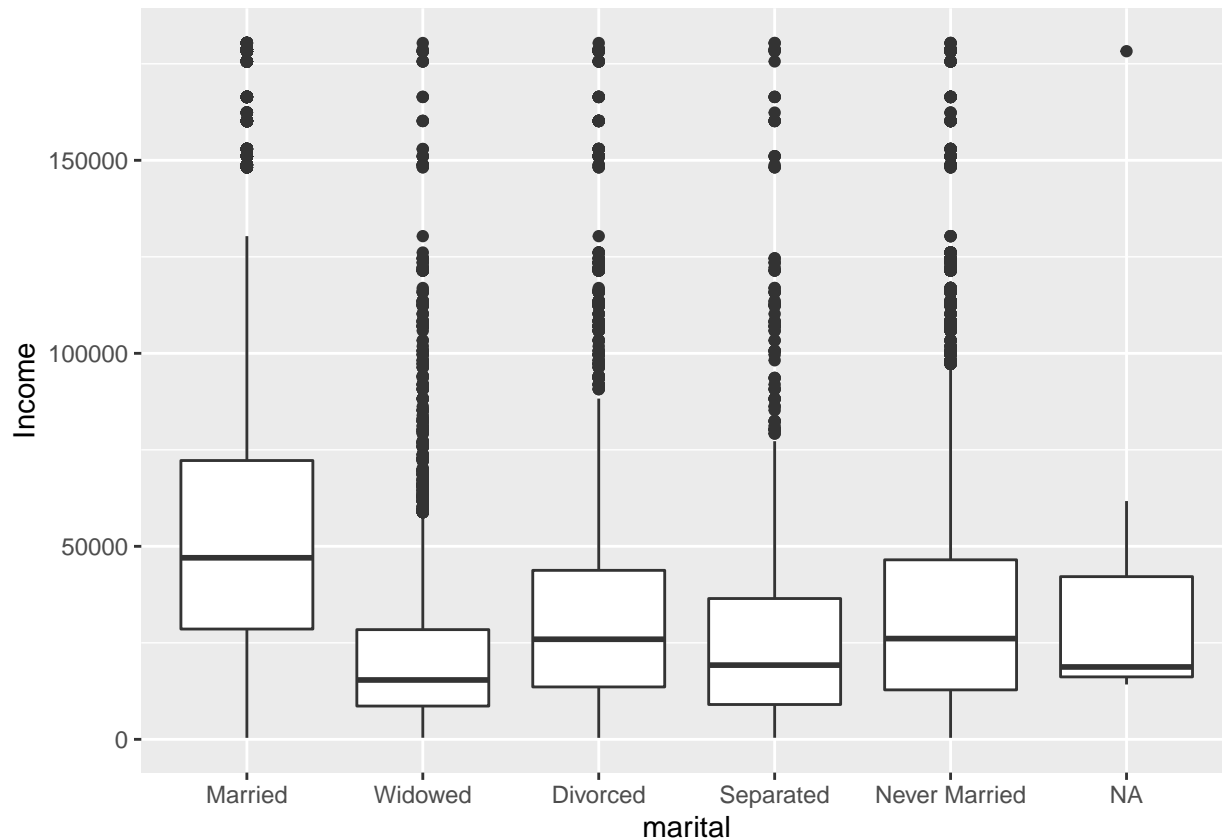
This implies that we are 95% confident that between 6 to 7 percent men are more self-employed than women.

Are there differences between the average salaries of people having different marital status?

Note that the marital status of the respondents are reported in “marital” column of the data frame. First, we illustrate the boxplots of the income with respect to the marital status.

```
ggplot(data=gss, aes(x=marital, y=coninc)) + geom_boxplot() + ylab("Income")
```

```
## Warning: Removed 5829 rows containing non-finite values (stat_boxplot).
```



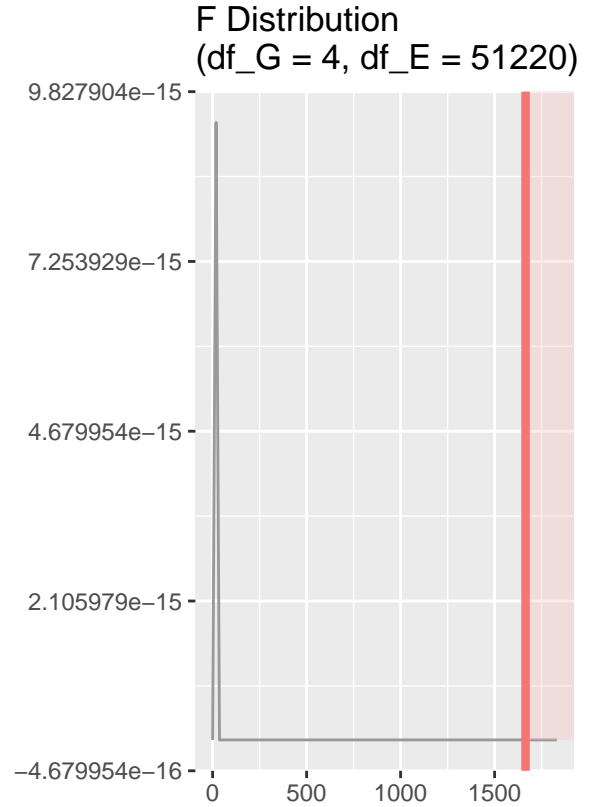
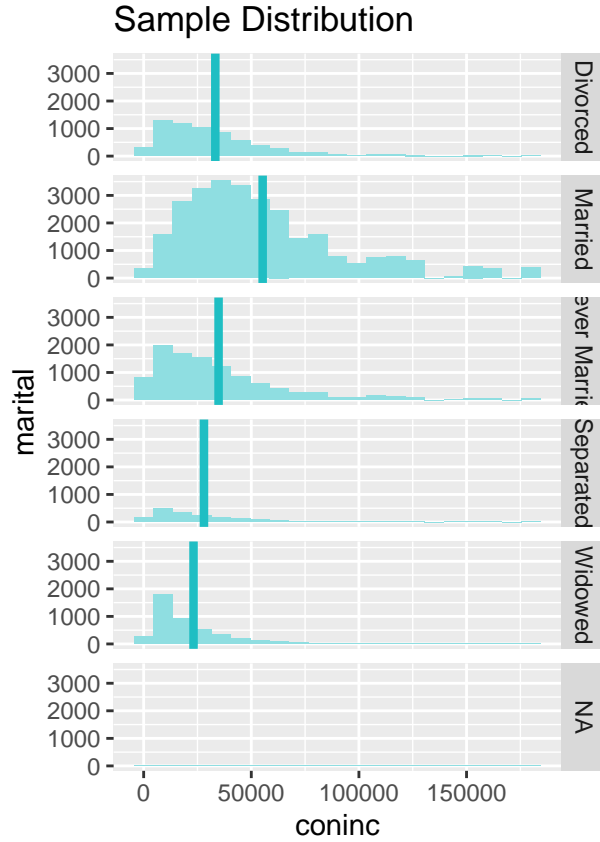
Therefore, we have the problem of comparing independent means, which can be addressed using ANOVA. For this example we have considered the significance-level $\alpha = 0.01$.

```
gss <- gss %>% drop_na(coninc)

inference(data=gss, y=coninc, x=marital, statistic="mean", method="theoretical",
          type="ht", alternative="greater", conf_level=0.99 )

## Response variable: numerical
## Explanatory variable: categorical (5 levels)
## n_Married = 27931, y_bar_Married = 55280.2263, s_Married = 37173.3465
## n_Widowed = 4569, y_bar_Widowed = 23165.4612, s_Widowed = 24106.9344
## n_Divorced = 6516, y_bar_Divorced = 33293.19, s_Divorced = 28768.4118
## n_Separated = 1780, y_bar_Separated = 28010.3556, s_Separated = 28907.1412
## n_Never Married = 10429, y_bar_Never Married = 34804.8646, s_Never Married = 31581.7963
##
## ANOVA:
##           df      Sum_Sq    Mean_Sq      F    p_value
## marital      4 7608229263438.02 1902057315859.5 1664.5183 < 0.0001
## Residuals 51220 58529469597965.5 1142707333.0333
## Total      51224 66137698861403.5
##
## Pairwise tests - t tests with pooled SD:
## # A tibble: 10 x 3
##   group1      group2    p.value
##   <chr>      <chr>      <dbl>
## 1 Widowed    Married      0
## 2 Divorced   Married      0
```

```
## 3 Divorced      Widowed  3.05e- 54
## 4 Separated     Married  2.49e-236
## 5 Separated     Widowed  2.91e-  7
## 6 Separated     Divorced  5.15e-  9
## 7 Never Married Married    0
## 8 Never Married Widowed  1.32e- 83
## 9 Never Married Divorced  4.63e-  3
## 10 Never Married Separated 4.68e- 15
```



The obtained value for p-value is small (less than 0.01);, therefore, one can conclude that at least one pair of means are different.

The pairwise test can reveal the means of which pair are different or not different. Testing many pairs of group is called multiple comparisons. The “Benferroni correction” suggests that a more stringent significance level is more appropriate for these tests. Actually, α should be adjusted by the number of comparisons being considered:

$$(\text{Benerroni correction}) \alpha^* = \frac{\alpha}{K}, \quad K : \text{number of comparisons}, \quad k : \text{number of means}, \quad K = \frac{k(k-1)}{2} \quad (1)$$

This implies that for this example, we have to adjust the significance level to the new significance level as $\alpha^* = \frac{0.01}{10} = 0.001$, $K = \frac{5 \times 4}{2} = 10$.

Now, the pairwise comparison shows that for all pairs the means are different but the p-value between “never married” and “divorced” is greater than the p-value=0.001; thus one can conclude that for all pairs the means are different unless for the pair “never married” and “divorced” based on the considered significance-level.

Are attitudes toward sex education and subjective social class independent?

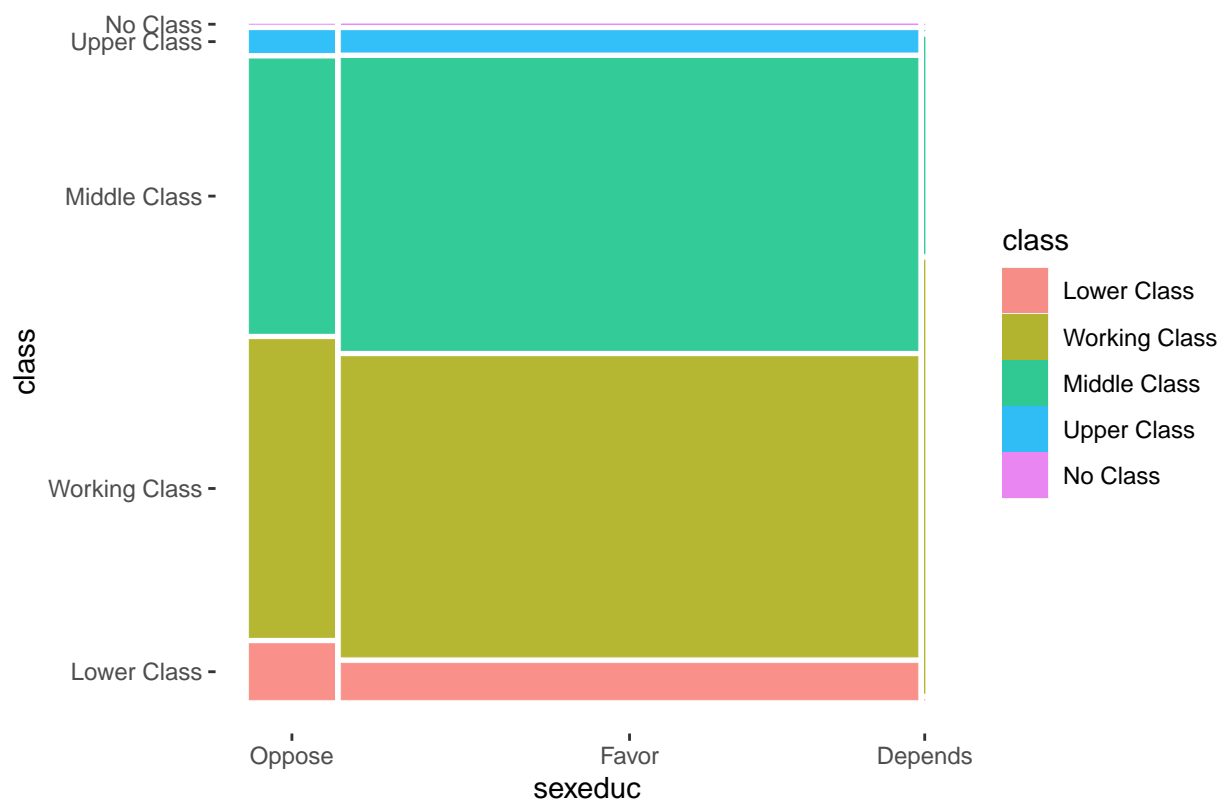
Note that attitude toward sex education is reported in “sexeduc” column, and the subjective social class is given in “class” column.

```
table(gss$sexeduc, gss$class)
```

```
##
##           Lower Class Working Class Middle Class Upper Class No Class
##   Favor           1306           10780           10482           765           1
##   Oppose           297            1579            1455            116           0
##   Depends           0              6              3              0           0
```

We illustrated the related data in a mosaic plot as follows:

```
# Change the order of sexeduc factors for a better representation
gss_modified <- gss %>% mutate(sexeduc=factor(sexeduc, levels=c("Oppose", "Favor", "Depends")))
ggplot(data=gss_modified) +
  geom_mosaic(aes(x=product(sexeduc), fill=class), na.rm=TRUE) +
  theme_mosaic()
```



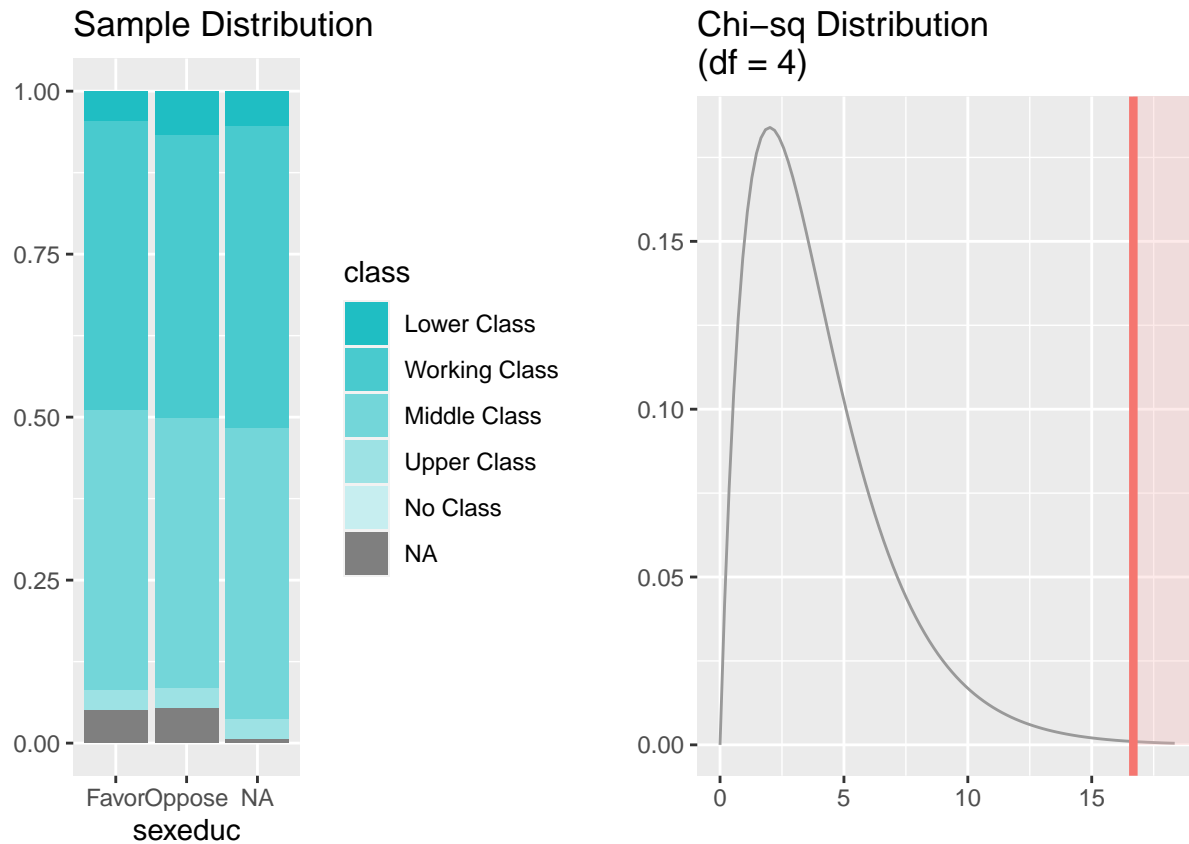
This problem can be addressed by using chi-square independence test. We have two categorical variables one of which has more than two levels. One can resort to “inference” to tackle this problem as follows:

```
gss <- gss %>% drop_na(abany)

inference(data=gss, y=class, x=sexeduc, statistic="proportion", method="theoretical",
          type="ht", alternative="greater", success="Yes")
```

```
## Warning in chisq.test(x, y, correct = FALSE): Chi-squared approximation may be
## incorrect
```

```
## Response variable: categorical (5 levels)
## Explanatory variable: categorical (3 levels)
## Observed:
##      y
## x      Lower Class Working Class Middle Class Upper Class No Class
## Favor      566      5501      5326      380      1
## Oppose      119      765      729      55      0
##
## Expected:
##      y
## x      Lower Class Working Class Middle Class Upper Class No Class
## Favor    599.99926    5488.4603    5303.6431    381.02143 0.8759113
## Oppose     85.00074     777.5397     751.3569     53.97857 0.1240887
##
## H0: sexeduc and class are independent
## HA: sexeduc and class are dependent
## chi_sq = 16.68, df = 4, p_value = 0.0022
```



Since the obtained value for p-value is smaller than 0.05, therefore one can conclude that the attitude toward sex education and the subjective social class are not independent.