

# Statistical inference with the GSS data

## Setup

### Load packages

```
#devtools::install_github("haleyjeppson/ggmosaic")
#install.packages('vcd')
#install.packages('pastecs')
#install.packages(c("ggplot2", "ggpubr", "tidyverse", "broom", "AICcmodavg"))
#install.packages('tinytex')
#tinytex::install_tinytex()

library(ggplot2)
library(dplyr)
library(statsr)
library(readr)
library(ggmosaic)
library(vcd)
library(pastecs)
library(skimr)
library(ggpubr)
library(tidyverse)
library(broom)
library(AICcmodavg)
library(car)
library(plyr)
```

### Load data

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

## Part 1: Data

### How are the observations in the sample collected?

According to the official GSS website, the GSS is a “nationally representative” survey of adults in the United States. It utilizes an “area probability design that randomly selects respondents in households across the United States.”

Due to this, it's safe to assume that the data collected is representative of the population of adults at large. However, it's likely that due to the nature of data collection, causal connections cannot be inferred, but the results can be generalized.

## Part 2: Research questions

### Investigating the data

#### Potentially interesting variables to investigate further

1. Religiosity (relig) and Education (degree). Relationship between level of education and religious affiliation?
2. Religiosity (relig) and Political Party (partyid). Relationship between religiosity and political party?  
\*\*\*
3. Religiosity (relig), Political Party (partyid) and Education level (degree).
4. Is level of education (degree) and family income (coninc) associated? \*\*\*
5. News and overall happiness.
6. What is the average # of Children adults have.
7. What is the average age of adults in the US?

```
#### Investigate the dataset
```

```
#head(gss)
#colnames(gss)
#names(gss)
#str(gss)
#glimpse(gss)
#skim_without_charts(gss)
#summary(gss)
```

## Analysis #1 - Religion and Party Affiliation (Chi-Square Test of Independence)

### Part 3: Exploratory data analysis

Does there appear to be a relationship between religiosity and political party affiliation?

The parties contained in this analysis are Democrat, Independent, Republican and Other.

For the purposes of this analysis, anyone affiliated with a religion is considered religious. Those that are nearly or strongly leaning toward a particular political party will be considered associated with that party.

#### Initial Analysis and Cleaning

```
##### Keep specific columns for analysis.
trimmed_gss_analysis_one <- gss %>%
  select(c(relig, partyid))

##### Check for N/A's.
sapply(trimmed_gss_analysis_one, function(x) sum(is.na(x)))
```

```
## relig partyid
##      233      327
```

```
##### Remove rows with N/A' (for this analysis, it was acceptable.)
trimmed_gss_analysis_one <- na.omit(trimmed_gss_analysis_one)

##### Recode data for easier analysis.
##### Combine religion and party affiliation values.

trimmed_gss_analysis_one <- trimmed_gss_analysis_one %>% mutate(
  relig = case_match(
    relig,
    "None" ~ "Not Religious",
    .default = "Religious"),
  partyid = case_match(
    partyid,
    "Strong Democrat"~"Democrat",
    "Not Str Democrat"~"Democrat",
    "Ind,Near Dem"~"Democrat",
    "Independent"~"Independent",
    "Ind,Near Rep"~"Republican",
    "Not Str Republican"~"Republican",
    "Strong Republican"~"Republican",
    "Other Party"~"Other"
  ),
)

##### Rename column headers.
colnames(trimmed_gss_analysis_one) <- c("Religiosity","Political Party")

##### Create contingency table
(partyid_relig_n <- (table(trimmed_gss_analysis_one$Religiosity,
                           trimmed_gss_analysis_one$'Political Party')))
```

```
##
##          Democrat Independent Other Republican
## Not Religious    3031      1548    195      1306
## Religious       24795     6891    663     18132
```

```
##### Export contingency table and trimmed data for external Chi-Sq Test of Independence in
Excel.
write.csv(partyid_relig_n, "party_religion.csv", row.names = TRUE)
```

## Plots

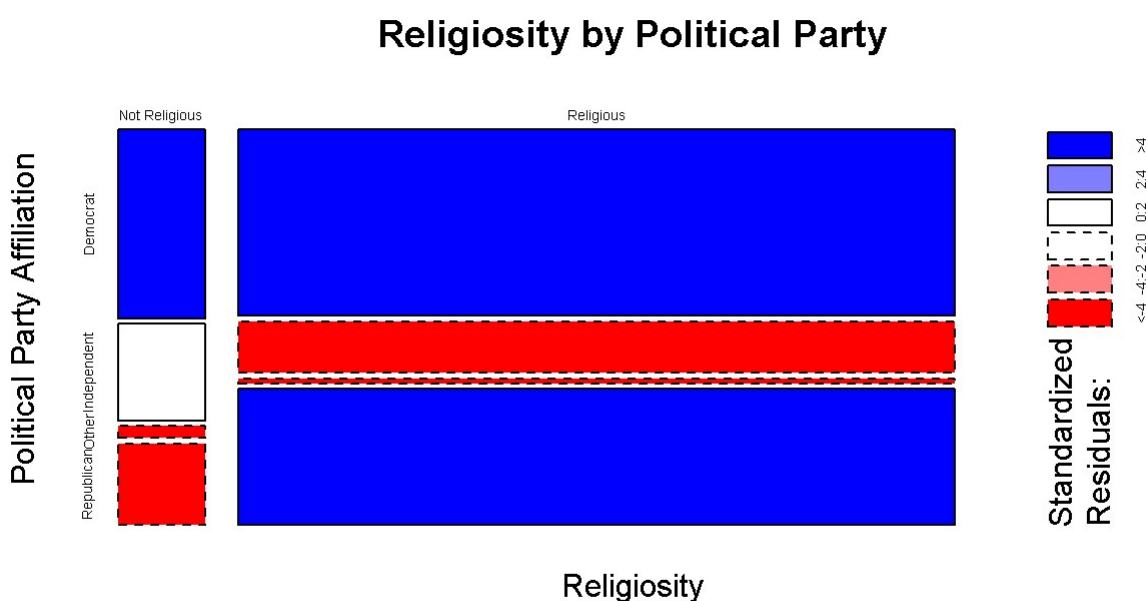
### Mosaic Plot with Residuals

The following plot shows the the association between Political Party affiliation and religiosity. The data indicate that Democrats, Independents and Other see a higher proportion of Non-Religious members, in

contrast to Republicans which see the opposite. Based on this, we can determine that political party and religiosity are dependent on one another.

The standardized residuals indicate the difference between the observed and expected values.

```
mosaicplot(partyid_relig_n, main='Religiosity by Political Party',
           col=TRUE,
           xlab = "Religiosity",
           ylab="Political Party Affiliation",
           shade = TRUE,
           off = 4,
           cex.axis = 0.45,
           margin = 1)
```



## Summary Statistics

### Contingency Tables

The summary statistics indicate a significant difference between the proportion of religious and non-religious individuals affiliated with a particular political party and those that are not

```
(partyid_relig_counts <- (addmargins(table(trimmed_gss_analysis_one$Religiosity,
                                           trimmed_gss_analysis_one$'Political Party'))))
```

```
##
##           Democrat Independent Other Republican   Sum
## Not Religious    3031      1548   195        1306  6080
## Religious       24795      6891   663        18132 50481
## Sum             27826      8439   858        19438 56561
```

```
(partyid_relig_prop <- (prop.table(addmargins(table(trimmed_gss_analysis_one$Religiosity,
                                                    trimmed_gss_analysis_one$'Political Party')))))
```

```
##
##               Democrat Independent      Other  Republican
## Not Religious 0.0133970404 0.0068421704 0.0008619013 0.0057725288
## Religious    0.1095940666 0.0304582663 0.0029304645 0.0801435618
## Sum          0.1229911069 0.0373004367 0.0037923658 0.0859160906
##
##               Sum
## Not Religious 0.0268736408
## Religious    0.2231263592
## Sum          0.2500000000
```

## Part 4: Inference

### Hypotheses

Does there appear to be a relationship between religiosity and political party affiliation?

H0 (nothing is going on) Religion and Party affiliation are independent. Religiosity does not vary by Political Party.

HA (something is going on) Religion and party affiliation are dependent. Religiosity does vary by Political Party.

### Evaluating the Hypothesis

1. Quantify how different the observed counts are from the expected counts.
2. Significant deviations from what is expected based on sampling variation (chance) alone provide strong evidence for the alternative hypothesis.

### Conditions for the Chi-Square Test of Independence

#### Independence:

1. Sampled observations are independent.
  - Random sampling was employed for the survey.
2.  $n < 10\%$  of the population if sampling without replacement.
  - Observations are less than 10% of the population.
3. Each case only contributes to one cell in the table.
  - It is safe to assume that each individual surveyed reported once, although it's not impossible someone took part more than once, and therefore is part of multiple cells. For the sake of this analysis, we can assume that each case has only contributed to one cell.

#### Sample Size:

1. Each cell must have at least 5 expected cases.

- True.

## Test the hypothesis

What is the overall religiosity in the sample?

```
50481/55703
```

```
## [1] 0.9062528
```

Test the hypothesis that religiosity and political party affiliation are associated at the 5% significance level.

- $\chi^2 = 970.29$
- $df = 3$

```
(chiSq <- pchisq(970.29, 3, lower.tail = FALSE))
```

```
## [1] 5.012349e-210
```

```
if(chiSq < 0.05) {
  print("Reject the null hypothesis")
} else {
  print("Fail to reject the null hypothesis")
}
```

```
## [1] "Reject the null hypothesis"
```

## Conclusion

The  $p_{\text{value}}$  of  $5.012349e-210$  is less than the 0.05 significance level. Therefore, we reject the null hypothesis in favour of the alternative; there is sufficient evidence indicating there is an association between Political Party affiliation and Religiosity.

# Analysis #2 - Level of respondents education and household income (ANOVA and pairwise tests (theoretical only))

## Part 3: Exploratory data analysis

Does there appear to be a relationship between the level of respondents education and their household income?

The levels of education contained in this analysis are: Limited HighSchool, HighSchool, Junior College, Bachelor, Graduate.

## Initial Analysis and Cleaning

```
##### Keep specific columns for analysis.
trimmed_gss_analysis_two <- gss %>%
  select(c(coninc, degree))

##### Check for N/A's.
sapply(trimmed_gss_analysis_two, function(x) sum(is.na(x)))
```

```
## coninc degree
## 5829 1010
```

```
##### Replace missing values with the mean of the column.
#trimmed_gss_analysis_two$coninc[is.na(trimmed_gss_analysis_two$coninc)] <- #mean(trimmed_g
ss_analysis_two$coninc, na.rm=TRUE)

#### Omit rows with NA's (acceptable due to the large amount of data collected.)
trimmed_gss_analysis_two <- na.omit(trimmed_gss_analysis_two)

##### Rename column headers.
colnames(trimmed_gss_analysis_two) <- c("Family_Income", "Degree")

# ##### Recode data for easier analysis.
# trimmed_gss_analysis_two <- trimmed_gss_analysis_two %>% mutate(
#   Degree = case_match(
#     Degree,
#     'Lt High School' ~ '1',
#     'High School' ~ '2',
#     'Junior College' ~ '3',
#     'Bachelor' ~ '4',
#     'Graduate' ~ '5',
#   )
# )

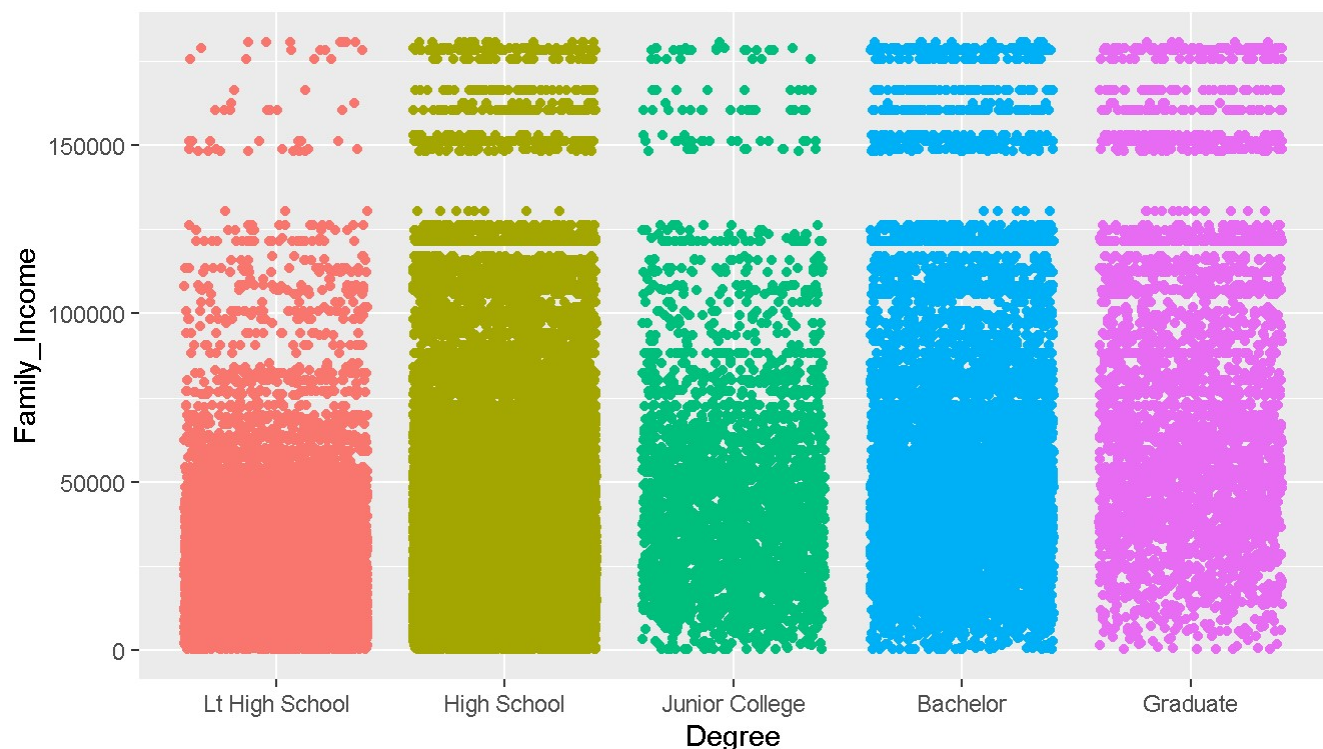
#### Export data for further analysis in Excel and additional formatting.
write.csv(trimmed_gss_analysis_two, "GSS_FamIncome_Degree_ANOVA.csv", row.names = TRUE)
```

## Plots

### Jitter Plot

The graph indicates a clear difference between various levels of respondent education and associated family earnings. The largest appears to be between Limited High School and High School; however, there are apparent differences among many of the degrees of education.

```
ggplot(trimmed_gss_analysis_two) +
  aes(x = Degree, y = Family_Income, color = Degree) +
  geom_jitter() +
  theme(legend.position = "none")
```



## Part 4: Inference

### Hypotheses

Does there appear to be a relationship between the level of respondents education and their household income?

- $H_0$ : The average family income is the same across all degrees of education. ( $\mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5$ )
- $H_A$ : The average family income differs between at least one pair of degrees of education.

### Evaluating the Hypothesis

#### Conditions for the ANOVA Test

```
res_aov <- aov(Family_Income ~ Degree,
  data = trimmed_gss_analysis_two
)
```

#### Independence

Within: sampled observations must be independent of each other 1. Random sample / assignment 2. Each  $n_j$  less than 10% of respective population 3. Always important, but sometimes difficult to check - Respondents were sampled randomly as part of the survey, therefore, it is safe to assume that the observations are independent of one another. - The observations are less than 10% of the respective populations.

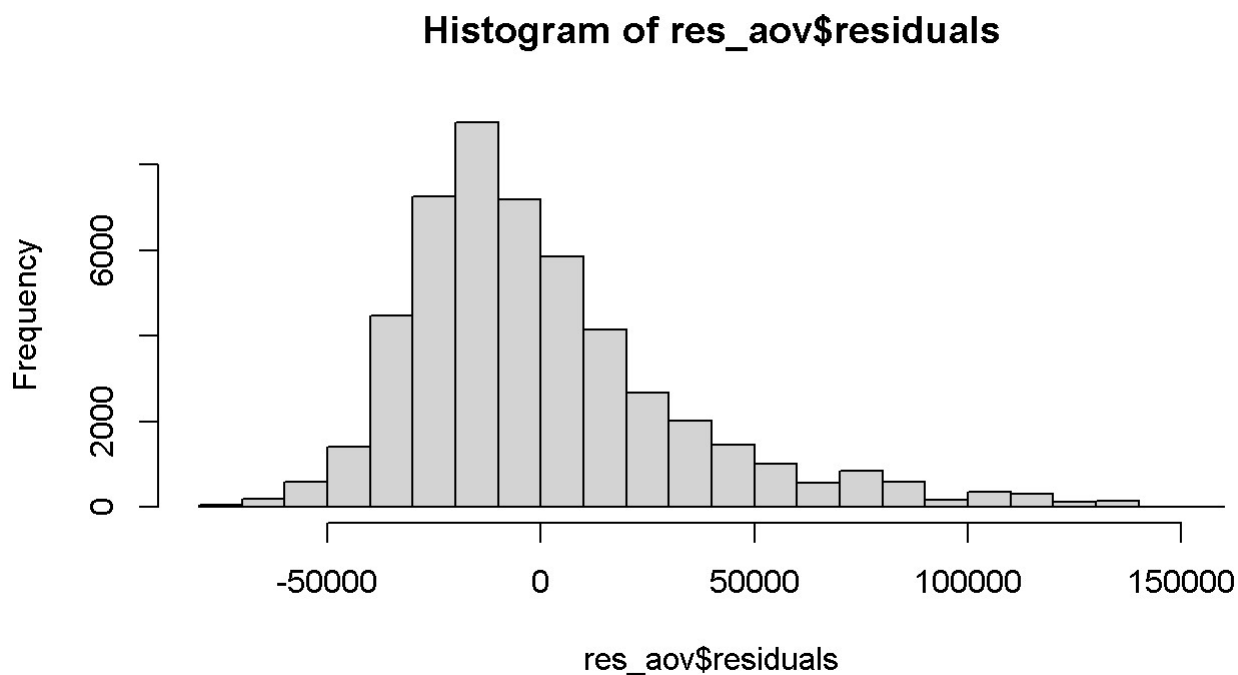
Between: groups must be independent of each other 1. Carefully consider whether the groups may be dependent - Random sampling was employed, so it's safe to assume that each observation is independent, and therefore each group is also independent of one another as well as each observation is contained within one cell.



## Approximately Normal

1. Distribution of response variable within each group should be approximately normal
  2. Especially important when sample sizes are small
- Although there is skew in the data, the sample size is very large, therefore this shouldn't be an issue.

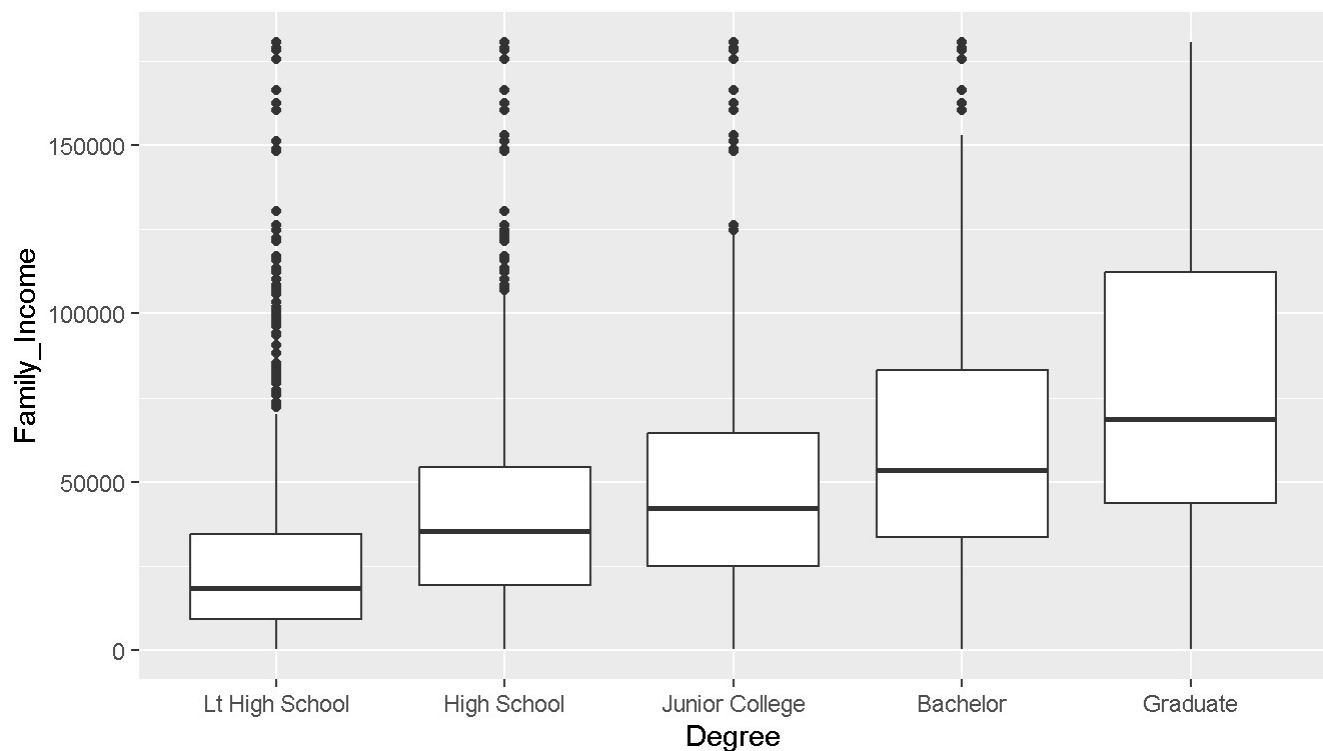
```
hist(res_aov$residuals)
```



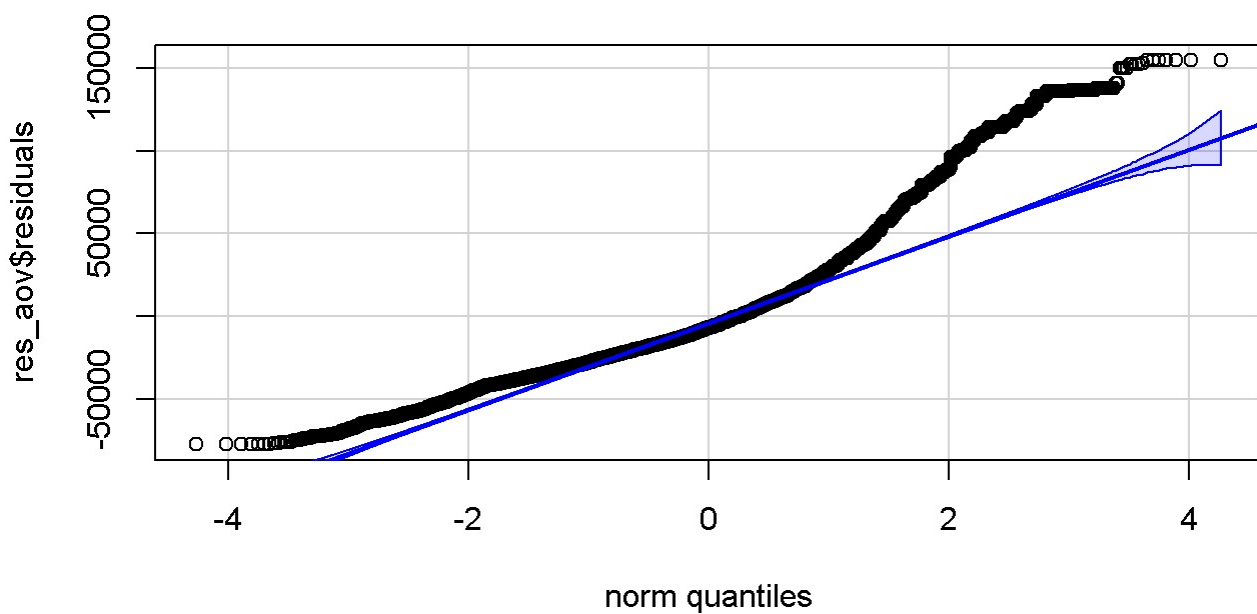
## Constant Variance

1. Variability should be consistent across groups: homoscedastic groups
  2. Especially important when sample sizes differ between groups
- Variance does not appear to be constant among all groups, therefore an alternative ANOVA test is required.

```
ggplot(trimmed_gss_analysis_two) +  
  aes(x = Degree, y = Family_Income) +  
  geom_boxplot()
```



```
qqPlot(res_aov$residuals,
  id = FALSE # id = FALSE to remove point identification
)
```



## Test the hypothesis - ANOVA

How much variability is attributed to the explanatory variable? (Education)

```
1.08127E+13/6.53791E+13
```

```
## [1] 0.1653847
```

```
print("Approximately 16% of variability is due to the explanatory variable.")
```

```
## [1] "Approximately 16% of variability is due to the explanatory variable."
```

Test the hypothesis that Family Income and respondent Education are associated at the 5% significance level.

```
oneway.test(Family_Income ~ Degree,
  data = trimmed_gss_analysis_two,
  var.equal = FALSE # assuming unequal variances
)
```

```
##
## One-way analysis of means (not assuming equal variances)
##
## data: Family_Income and Degree
## F = 2340.8, num df = 4, denom df = 10598, p-value < 2.2e-16
```

```
pf(2340.8,4,10598,lower.tail=FALSE)
```

```
## [1] 0
```

## Conclusion

- If p-value is small (less than  $\alpha$ ), reject  $H_0$ .
- The data provide convincing evidence that at least one pair of population means are different from each other (but we can't tell which one).
- If p-value is large, fail to reject  $H_0$ .
- The data do not provide convincing evidence that at least one pair of population means are different from each other, the observed differences in sample means are attributable to sampling variability (or chance).

The p\_value of 0 is less than the 0.05 significance level. Therefore, we reject the null hypothesis in favour of the alternative; there is sufficient evidence indicating that at least one pair of population means are different from one another, however, we do not know which.

## Test the hypothesis - Pairwise Testing

Modified Significance level given 5 levels: 0.02 SE for multiple pairwise comparisons: 382.773 Degrees of Freedom for multiple pairwise comparisons: 10267

1. Is there a difference between the average family income of Lt High School and High School respondents?

```
(lhs_HS <- 2 * pt(43.332, df = 10267, lower.tail = FALSE))
```

```
## [1] 0
```

```
if(lhs_HS < 0.02) {  
  print("Reject the null hypothesis")  
} else {  
  print("Fail to reject the null hypothesis")  
}
```

```
## [1] "Reject the null hypothesis"
```

2. Is there a difference between the average family income of Lt High School and Junior College respondents?

```
(lhs_JC <- 2 * pt(35.078, df = 2803, lower.tail = FALSE))
```

```
## [1] 8.33725e-224
```

```
if(lhs_JC < 0.02) {  
  print("Reject the null hypothesis")  
} else {  
  print("Fail to reject the null hypothesis")  
}
```

```
## [1] "Reject the null hypothesis"
```

3. Is there a difference between the average family income of Lt High School and Bachelor respondents?

```
(lhs_B <- 2 * pt(55.011, df = 7373, lower.tail = FALSE))
```

```
## [1] 0
```

```
if(lhs_B < 0.02) {  
  print("Reject the null hypothesis")  
} else {  
  print("Fail to reject the null hypothesis")  
}
```

```
## [1] "Reject the null hypothesis"
```

4. Is there a difference between the average family income of Lt High School and Graduate respondents?

```
(lhs_G <- 2 * pt(-39.771, df = 3565, lower.tail = FALSE))
```

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

```
if(lhs_G < 0.02) {  
  print("Reject the null hypothesis")  
} else {  
  print("Fail to reject the null hypothesis")  
}
```

```
## [1] "Fail to reject the null hypothesis"
```

5. Is there a difference between the average family income of High School and Junior College respondents?

```
(HS_JC <- 2 * pt(12.255, df = 2803, lower.tail = FALSE))
```

```
## [1] 1.131081e-33
```

```
if(HS_JC < 0.02) {  
  print("Reject the null hypothesis")  
} else {  
  print("Fail to reject the null hypothesis")  
}
```

```
## [1] "Reject the null hypothesis"
```

6. Is there a difference between the average family income of High School and Bachelor respondents?

```
(HS_B <- 2 * pt(50.723, df = 7373, lower.tail = FALSE))
```

```
## [1] 0
```

```
if(HS_B < 0.02) {  
  print("Reject the null hypothesis")  
} else {  
  print("Fail to reject the null hypothesis")  
}
```

```
## [1] "Reject the null hypothesis"
```

7. Is there a difference between the average family income of High School and Graduate respondents?

```
(HS_G <- 2 * pt(61.497, df = 3565, lower.tail = FALSE))
```

```
## [1] 0
```

```
if(HS_G < 0.02) {  
  print("Reject the null hypothesis")  
} else {  
  print("Fail to reject the null hypothesis")  
}
```

```
## [1] "Reject the null hypothesis"
```

8. Is there a difference between the average family income of Junior College and Bachelor respondents?

```
(JC_B <- 2 * pt(19.144, df = 2803, lower.tail = FALSE))
```

```
## [1] 7.112095e-77
```

```
if(JC_B < 0.02) {  
  print("Reject the null hypothesis")  
} else {  
  print("Fail to reject the null hypothesis")  
}
```

```
## [1] "Reject the null hypothesis"
```

9. Is there a difference between the average family income of Junior College and Graduate respondents?

```
(JC_G <- 2 * pt(33.828, df = 2803, lower.tail = FALSE))
```

```
## [1] 1.176363e-210
```

```
if(JC_G < 0.02) {
  print("Reject the null hypothesis")
} else {
  print("Fail to reject the null hypothesis")
}
```

```
## [1] "Reject the null hypothesis"
```

10. Is there a difference between the average family income of Bachelor and Graduate respondents?

```
(B_G <- 2 * pt(21.036, df = 3565, lower.tail = FALSE))
```

```
## [1] 1.058915e-92
```

```
if(B_G < 0.02) {
  print("Reject the null hypothesis")
} else {
  print("Fail to reject the null hypothesis")
}
```

```
## [1] "Reject the null hypothesis"
```

## Conclusion

All pairs but Lt High School vs Graduate showed a significant difference.

## Analysis #3 - Average # of Children

How many children do respondents in this sample have on average? Estimate the true, average number of children based on this sample with a 95% confidence interval.

## Part 3: Exploratory data analysis

```
##### Keep specific columns for analysis.
trimmed_gss_analysis_three <- gss %>%
  select(c(childs))

##### Check for N/A's.
sapply(trimmed_gss_analysis_three, function(x) sum(is.na(x)))
```

```
## childs
##      181
```

```
##### Replace missing values with the mean of the column.
#trimmed_gss_analysis_two$coninc[is.na(trimmed_gss_analysis_two$coninc)] <- #mean(trimmed_g
ss_analysis_two$coninc, na.rm=TRUE)

#### Omit rows with NA's (acceptable due to the large amount of data collected.)
trimmed_gss_analysis_three <- na.omit(trimmed_gss_analysis_three)

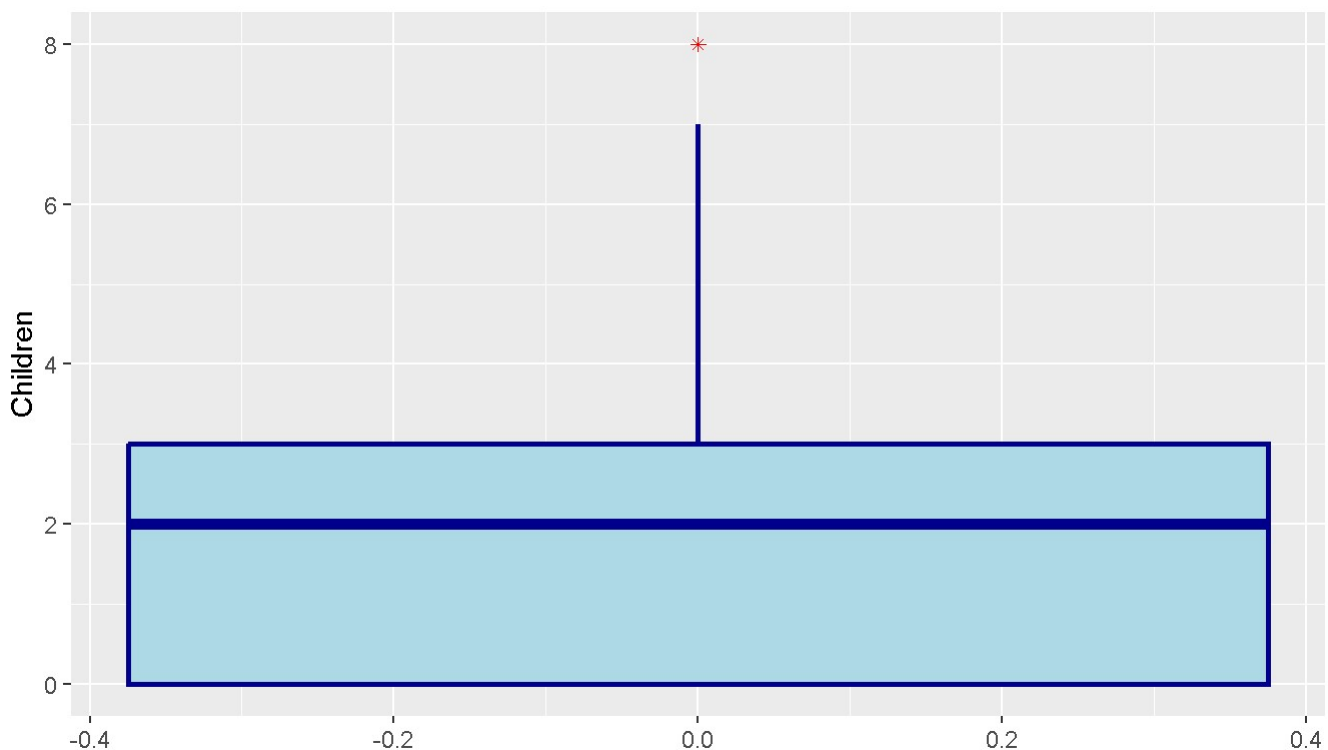
##### Rename column headers.
colnames(trimmed_gss_analysis_three) <- c("Children")

#### Export data for further analysis in Excel, if necessary.
write.csv(trimmed_gss_analysis_three, "c_data.csv", row.names = TRUE)
```

## Plots

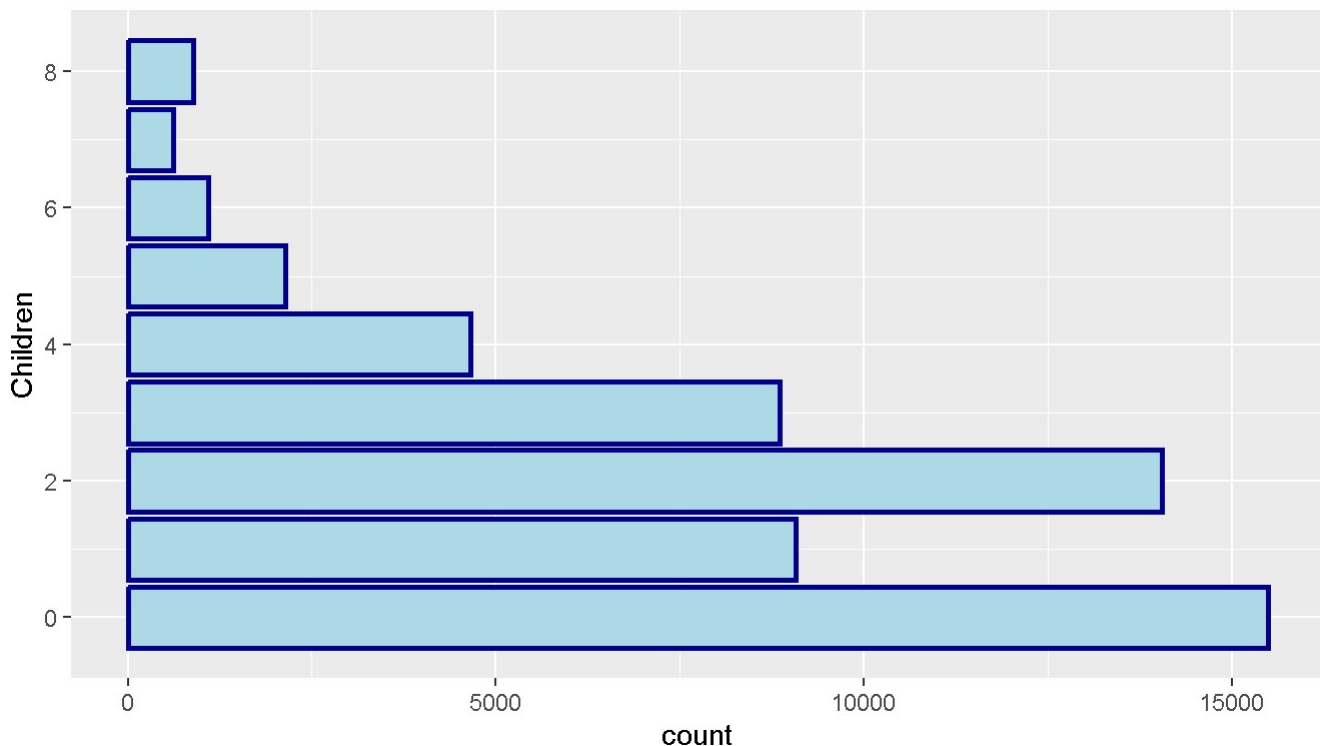
### Boxplot and Barplot

```
ggplot(trimmed_gss_analysis_three) +
  geom_boxplot(fill = "lightblue", color = "darkblue", size = 1, outlier.color = "red", out
lier.shape = 8) +
  aes(y = Children)
```





```
ggplot(trimmed_gss_analysis_three) +
  geom_bar(fill = "lightblue", color = "darkblue", linewidth = 1) +
  aes(y = Children)
```



#### #### Summary Statistics

```
summary(trimmed_gss_analysis_three)
```

```
##      Children
##  Min.   :0.000
## 1st Qu.:0.000
##  Median:2.000
##   Mean  :1.953
## 3rd Qu.:3.000
##   Max.  :8.000
```

## Part 4: Inference

### Conditions

#### Independence

1. Independence: Sampled observations must be independent.
2. random sample/assignment
3. if sampling without replacement,  $n < 10\%$  of population
  - Random sample &  $56,880 < 10\%$  of the population at large.
  - Independence cannot be guaranteed, as there is a chance multiple respondents from the same

household were surveyed, or individuals who are no longer in the same household provide data on the same children.

## Approximately Normal

1. Sample size/skew:  $n \geq 30$ , larger if the population distribution is very skewed.
  - Sample size/skew:  $56,880 \geq 30$ .
  - There is slight skew in the data, however, the sample size is large.

## Calculate the Confidence Interval

```
#### Calculate Mean, Median and Standard Deviation.
```

```
(trimmed_gss_analysis_three %>%
  summarize(Mean=mean(Children), Median=median(Children), standard_deviation=sd(Children)))
```

```
##           Mean Median standard_deviation
## 1  1.952848      2          1.791539
```

```
se = sd(trimmed_gss_analysis_three$Children/(sqrt(56880)))

CI_low = mean(trimmed_gss_analysis_three$Children - 1.96*se)
CI_high = mean(trimmed_gss_analysis_three$Children + 1.96*se)

(CI <- c(CI_low,CI_high))
```

```
## [1] 1.938125 1.967571
```

## Conclusion

We are 95% confident that adults, on average, have approximately 1.94 - 1.97 children.

## Test the hypothesis

A 95% confidence interval for the average number of children adults have was (1.94, 1.97). Based on this confidence interval, do these data support the hypothesis that adults on average have more than 1.95 children? (In this case, the sample mean.)

$H_0: \mu = 1.95$ : Adults have 1.95 children on average.  $H_A: \mu > 1.95$ : Adults have more than 1.95 children on average.

p-value:  $P(\text{observed or more extreme outcome} \mid H_0 \text{ true}) P(X > 1.95 \mid H_0: \mu = 1.95) X \sim N(\mu = 1.95, SE = 0.0075)$

- $N = 56880$
- $\bar{x} = 1.95$
- $SD = 1.79$
- $SE = 0.0075$

```
pnorm <- pnorm(((mean(trimmed_gss_analysis_three$Children - 1.95)/sd(trimmed_gss_analysis_three$Children/(sqrt(56880))))),lower.tail=FALSE)

if(pnorm < 0.05) {
  print("Reject the null hypothesis")
} else {
  print("Fail to reject the null hypothesis")
}
```

```
## [1] "Fail to reject the null hypothesis"
```

## Conclusion

Since the p\_value is high (greater than 0.05), we fail to reject the null hypothesis. There is insufficient data indicating that adults have more than 1.95 children on average.

# Analysis #4 - Average age of adults in the US

## Part 3: Exploratory data analysis

What is the average age of the US adult population? Estimate the true, average age of adults in the US, based on this sample with a 95% confidence interval.

```
##### Keep specific columns for analysis.
trimmed_gss_analysis_four <- gss %>%
  select(c(age))

##### Check for N/A's.
sapply(trimmed_gss_analysis_four, function(x) sum(is.na(x)))
```

```
## age
## 202
```

```
##### Replace missing values with the mean of the column.
#trimmed_gss_analysis_two$coninc[is.na(trimmed_gss_analysis_two$coninc)] <- #mean(trimmed_gss_analysis_two$coninc, na.rm=TRUE)

#### Omit rows with NA's (acceptable due to the large amount of data collected.)
trimmed_gss_analysis_four <- na.omit(trimmed_gss_analysis_four)

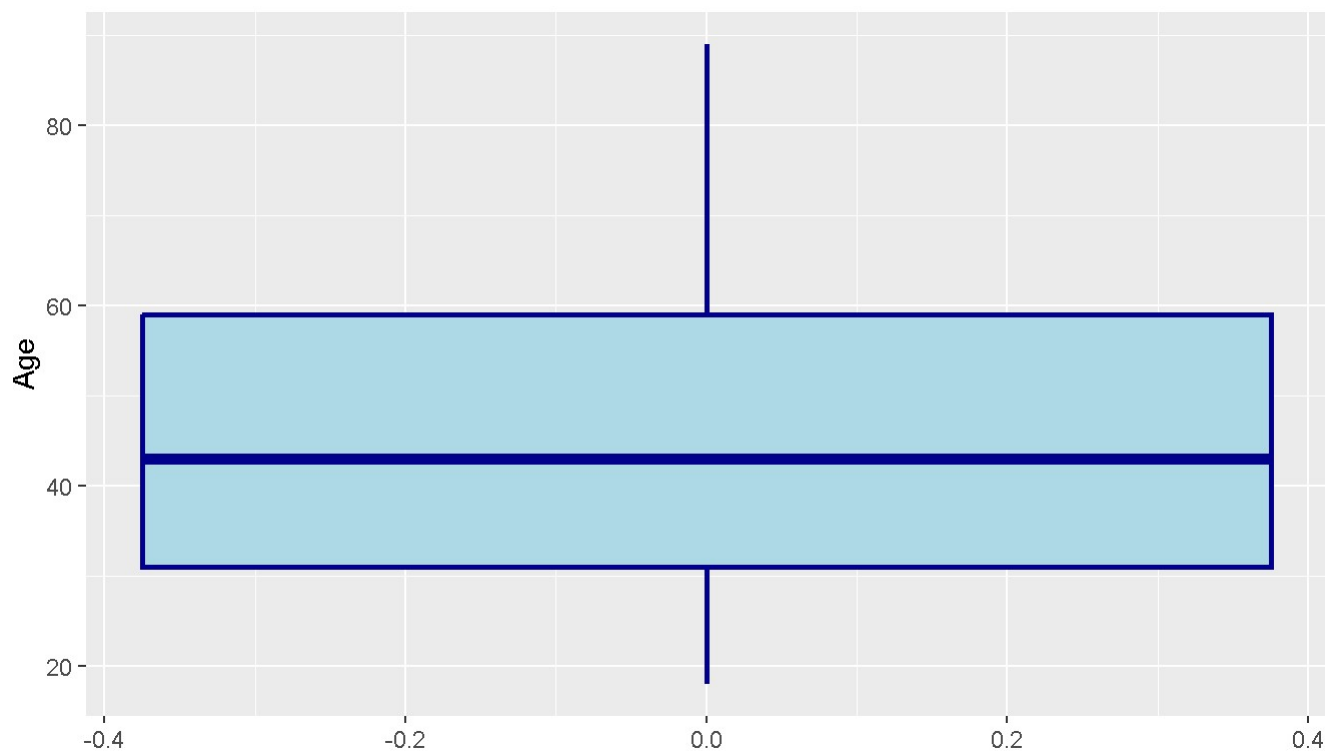
##### Rename column headers.
colnames(trimmed_gss_analysis_four) <- c("Age")

#### Export data for further analysis in Excel, if necessary.
write.csv(trimmed_gss_analysis_four, "a_data.csv", row.names = TRUE)
```

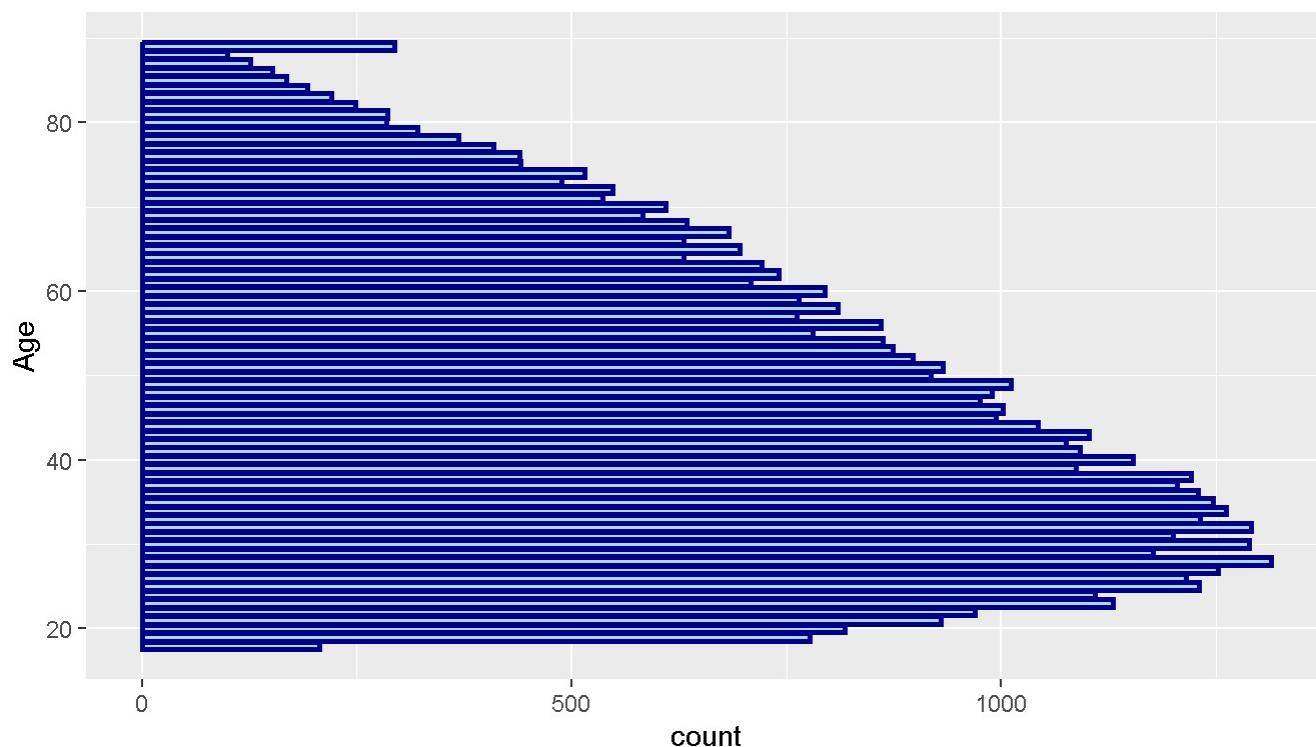
## Plots

## Boxplot and Barplot

```
ggplot(trimmed_gss_analysis_four) +  
  geom_boxplot(fill = "lightblue", color = "darkblue", size = 1, outlier.color = "red", outlier.shape = 8) +  
  aes(y = Age)
```



```
ggplot(trimmed_gss_analysis_four) +  
  geom_bar(fill = "lightblue", color = "darkblue", linewidth = 1) +  
  aes(y = Age)
```



## Summary Statistics

```
summary(trimmed_gss_analysis_four)
```

```
##      Age
##  Min.   :18.0
## 1st Qu.:31.0
##  Median :43.0
##   Mean  :45.7
## 3rd Qu.:59.0
##   Max.  :89.0
```

## Part 4: Inference

### Conditions

#### Independence

1. Independence: Sampled observations must be independent.
2. random sample/assignment
3. if sampling without replacement,  $n < 10\%$  of population
  - Random sample &  $56,859 < 10\%$  of the adult population at large.
  - One adult is surveyed per home. Being that this analysis is on the average age of adults, even if, somehow, more than one adult from the same household provided a response, it wouldn't drastically impact the integrity of the analysis.

#### Approximately Normal

1. Sample size/skew:  $n \geq 30$ , larger if the population distribution is very skewed.

- Sample size/skew:  $56,859 \geq 30$ .
- There is slight skew in the data, however, the sample size is large.

## Calculate the Confidence Interval

```
#### Calculate Mean, Median and Standard Deviation.
(trimmed_gss_analysis_four %>%
  summarize(Mean=mean(Age), Median=median(Age), standard_deviation=sd(Age)))
```

```
##      Mean Median standard_deviation
## 1 45.69795     43          17.47211
```

```
se = sd(trimmed_gss_analysis_four$Age/(sqrt(56859)))

CI_low = mean(trimmed_gss_analysis_four$Age - 1.96*se)
CI_high = mean(trimmed_gss_analysis_four$Age + 1.96*se)

(CI <- c(CI_low, CI_high))
```

```
## [1] 45.55434 45.84157
```

We are 95% confident that the average adult age in the US is between 45.55 and 45.84 years.

## Test the hypothesis

A 95% confidence interval for the average age of adults have was (45.55434 45.84157). Based on this confidence interval, do these data support the hypothesis that the average age of Adults in the US is greater than 43 (the median)?

- $H_0: \mu = 43$ : Average age of Adults in the US is 43.
- $H_A: \mu > 43$ : Average age of Adults in the US is greater than 43.

p-value:  $P(\text{observed or more extreme outcome} \mid H_0 \text{ true}) P(X > 45.7 \mid H_0: \mu = 43) X \sim N(\mu = 43, SE = 0.0075)$

- $N = 56859$
- $\bar{x} = 45.7$
- $SD = 17.47211$
- $SE = 0.07327331$

```
pnorm <- pnorm(((mean(trimmed_gss_analysis_four$Age - 43)/sd(trimmed_gss_analysis_four$Age
  /(sqrt(56859))))),lower.tail=FALSE)

if(pnorm < 0.05) {
  print("Reject the null hypothesis")
} else {
  print("Fail to reject the null hypothesis")
}
```

```
## [1] "Reject the null hypothesis"
```

## Conclusion

Since the  $p$ -value is high (greater than 0.05), we reject the null hypothesis. There is sufficient data indicating that the average age of Adults in the Us is greater than 43.