# Quantitative Data Analysis and Interpretation
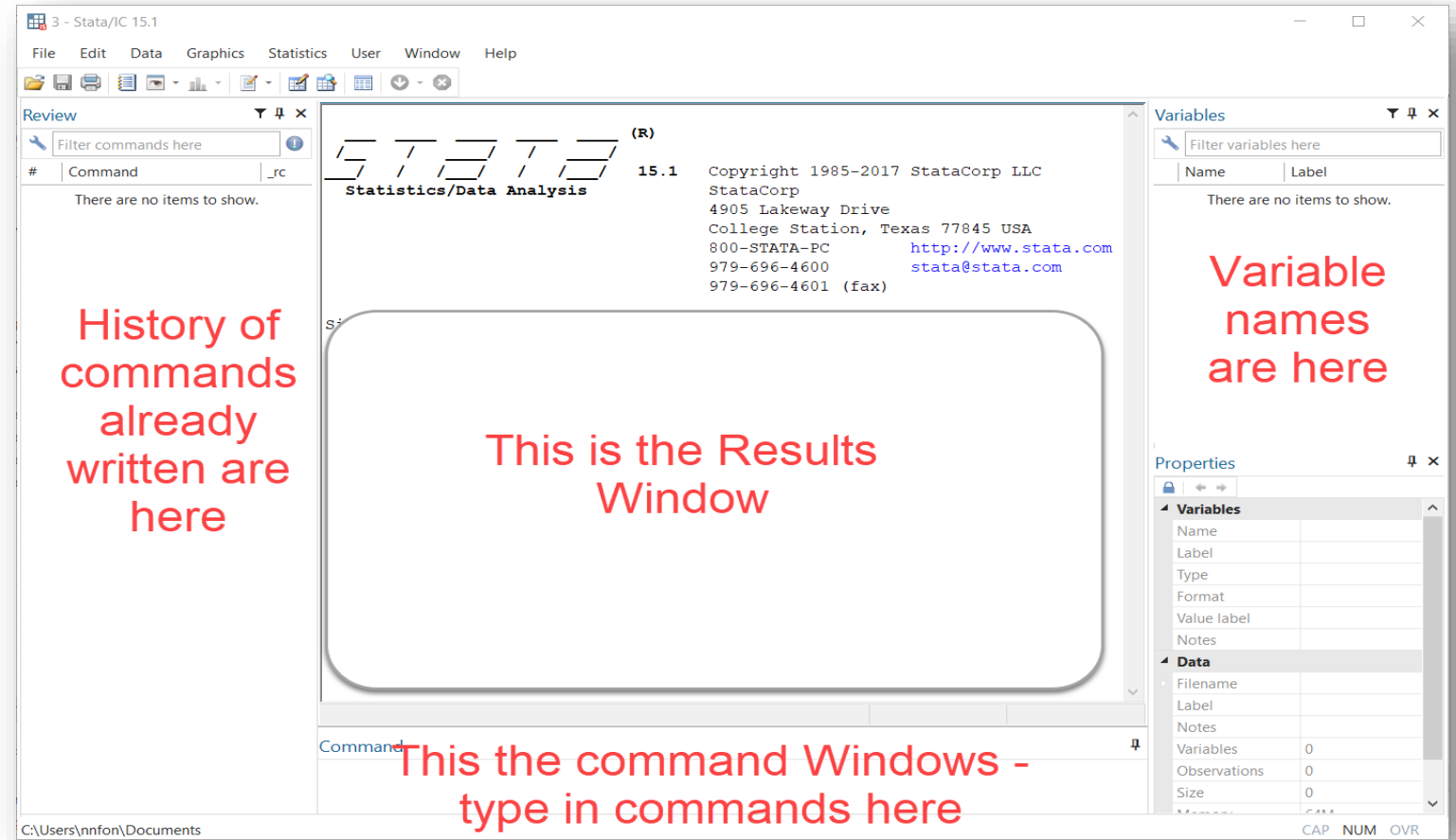
Neba Nfonsang

University of Denver

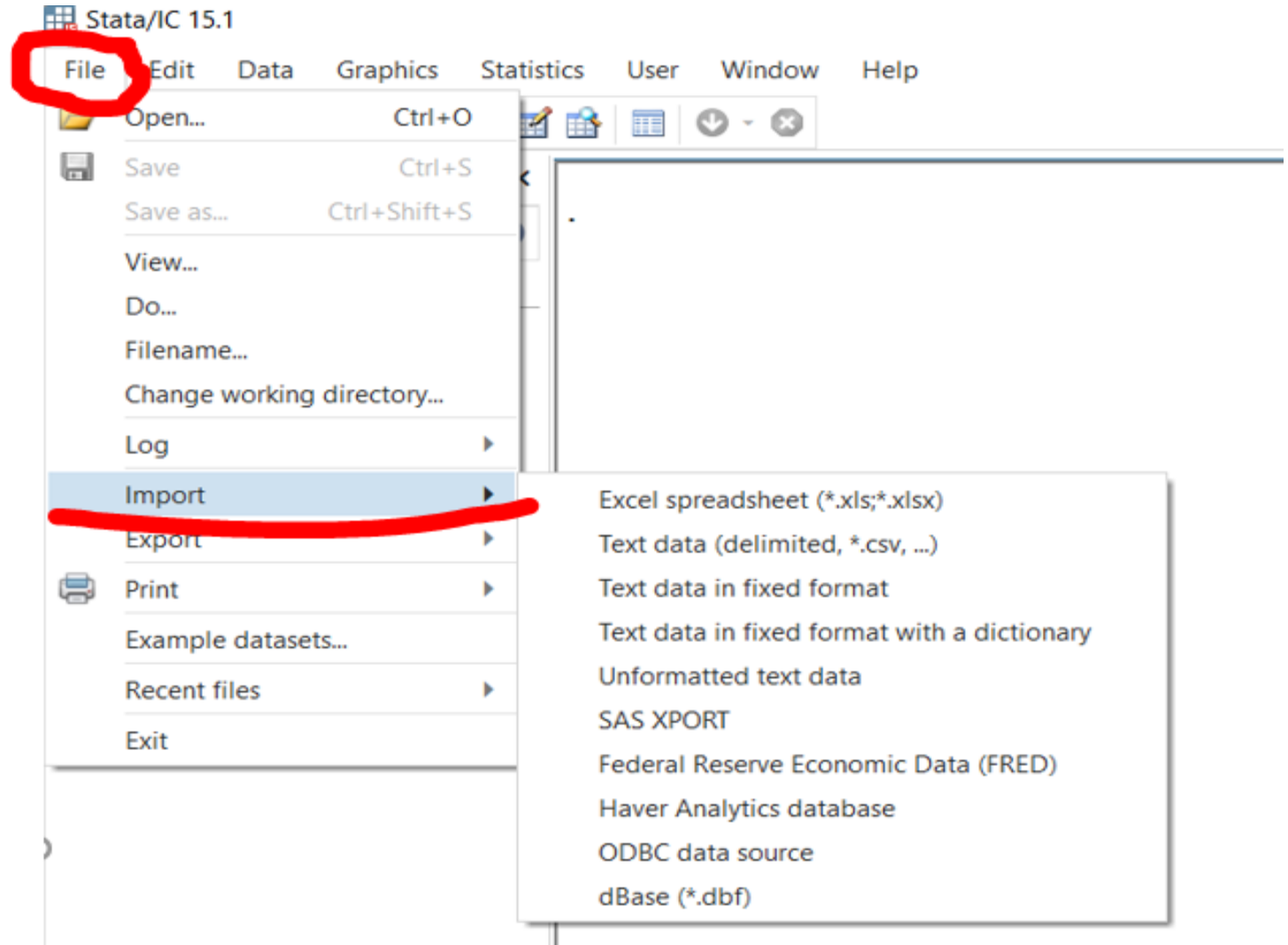# Statistical Data Analysis in Stata

# Stata Interface

- When you open Stata, you would see a Stata interface as shown.

- Basically, Stata can be used for:
  - ☐ Graphing
  - ☐ Data analysis
  - ☐ Data management



You can point and click or write commands to do a task in Stata.

# Import a Data File into Stata

- To import a data file:
    - ☐ Go to the **File** menu
    - ☐ Click on **Import**
    - ☐ Click on the file type you want to import. For example, **Text data (delimited, *.csv,…)**

# Import a Data File into Stata

- After clicking on the file type such as csv, this dialogue box will open:
  - ☐ Click the **Browse** button.
  - ☐ Navigate to the file on your computer and double click on the file.
  - ☐ Click **Submit** or **Ok** to import the file into Stata.
  - ☐ To use a command line, type: **import delimited** pathname

```
. import delimited C:\Users\nnfon\Desktop\hsb2.csv
```

# Data Editor



To check the data in Stata:
  Click on the **Data Editor** icon on the
  Stata Interface.
  You can edit the data from the data editor
  such as changing column names, etc.

# Run Descriptive Statistics

# Data Used

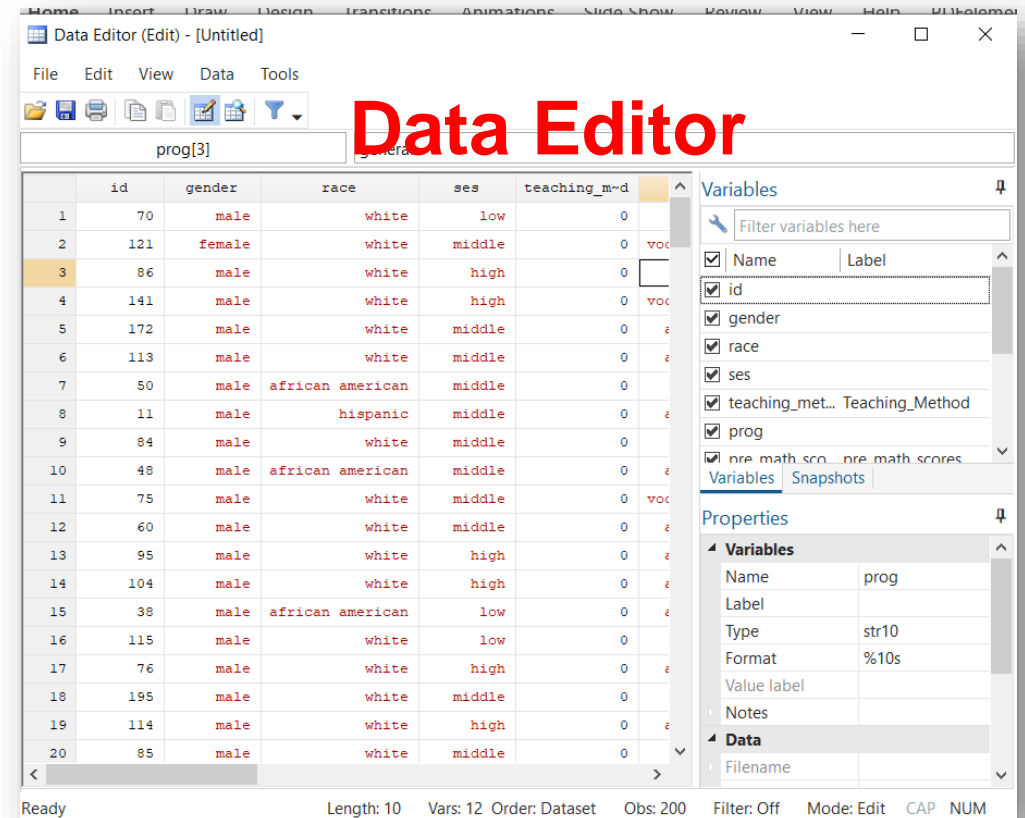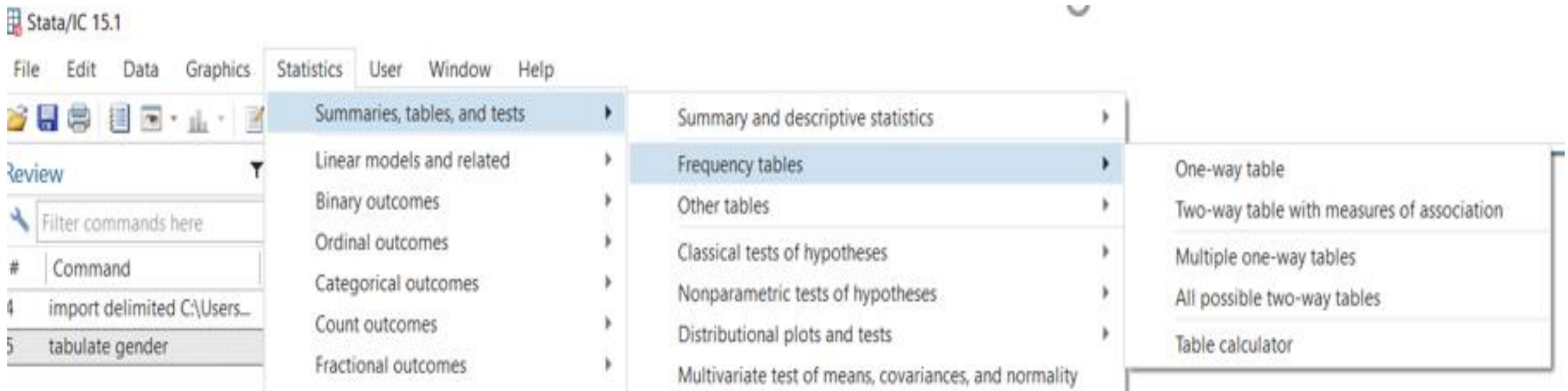| | id | gender | race | ses | teaching_m~d | prog | pre_math_s~e | post_math_~e | diff_pre_p~h | pre_reading | post_reading | diff_pre_p~g |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 70 | male | white | low | 0 | general | 57 | 52 | −5 | 41 | 47 | 6 |
| 2 | 121 | female | white | middle | 0 | vocational | 68 | 59 | −9 | 53 | 63 | 10 |
| 3 | 86 | male | white | high | 0 | general | 44 | 33 | −11 | 54 | 58 | 4 |
| 4 | 141 | male | white | high | 0 | vocational | 63 | 44 | −19 | 47 | 53 | 6 |
| 5 | 172 | male | white | middle | 0 | academic | 47 | 52 | 5 | 57 | 53 | −4 |
| 6 | 113 | male | white | middle | 0 | academic | 44 | 52 | 8 | 51 | 63 | 12 |
| 7 | 50 | male | african american | middle | 0 | general | 50 | 59 | 9 | 42 | 53 | 11 |
| 8 | 11 | male | hispanic | middle | 0 | academic | 34 | 46 | 12 | 45 | 39 | −6 |
| 9 | 84 | male | white | middle | 0 | general | 63 | 57 | −6 | 54 | 58 | 4 |
| 10 | 48 | male | african american | middle | 0 | academic | 57 | 55 | −2 | 52 | 50 | −2 |
| 11 | 75 | male | white | middle | 0 | vocational | 60 | 46 | −14 | 51 | 53 | 2 |
| 12 | 60 | male | white | middle | 0 | academic | 57 | 65 | 8 | 51 | 63 | 12 |
| 13 | 95 | male | white | high | 0 | academic | 73 | 60 | −13 | 71 | 61 | −10 |
| 14 | 104 | male | white | high | 0 | academic | 54 | 63 | 9 | 57 | 55 | −2 |
| 15 | 38 | male | african american | low | 0 | academic | 45 | 57 | 12 | 50 | 31 | −19 |
| 16 | 115 | male | white | low | 0 | general | 42 | 49 | 7 | 43 | 50 | 7 |
| 17 | 76 | male | white | high | 0 | academic | 47 | 52 | 5 | 51 | 50 | −1 |
| 18 | 195 | male | white | middle | 0 | general | 57 | 57 | 0 | 60 | 58 | −2 |
| 19 | 114 | male | white | high | 0 | academic | 68 | 65 | −3 | 62 | 55 | −7 |
| 20 | 85 | male | white | middle | 0 | general | 55 | 39 | −16 | 57 | 53 | −4 |

# Summary or Frequency Table

- For categorical data, frequencies can be obtained as follows:
  - ☐ Go to **Statistics → Summaries, tables, and test → Frequency tables → One-way table**

# Summary or Frequency Table



Select the categorical variable to be summarized

Results usually include the command that can be used to generate the results as well.

. **tabulate gender**

| gender | Freq. | Percent | Cum. |
|--------|-------|---------|------|
| female | 109 | 54.50 | 54.50 |
| male | 91 | 45.50 | 100.00 |
| Total | 200 | 100.00 | |

# Summary or Frequency Table

■ Syntax for generating a frequency table can also be run on the command windows as:

- ☐ **tabulate** followed by the variable name
- ☐ For example, **tabulate race**

```
. tabulate race

           race |      Freq.     Percent        Cum.
----------------+-----------------------------------
african american |         20       10.00       10.00
           asian |         11        5.50       15.50
        hispanic |         24       12.00       27.50
           white |        145       72.50      100.00
----------------+-----------------------------------
           Total |        200      100.00

. tabulate prog

           prog |      Freq.     Percent        Cum.
----------------+-----------------------------------
        academic |        105       52.50       52.50
         general |         45       22.50       75.00
      vocational |         50       25.00      100.00
----------------+-----------------------------------
           Total |        200      100.00
```

# Run Descriptive Statistics: **summarize**

- The **summarize** command can be used to generate:
  - ☐ Number of observations
  - ☐ Mean
  - ☐ Standard deviation
  - ☐ Minimum value
  - ☐ Maximum value for one or more variables.

```
. summarize pre_reading
```

| Variable | Obs | Mean | Std. Dev. | Min | Max |
|---|---|---|---|---|---|
| pre_reading | 200 | 52.645 | 9.368448 | 33 | 75 |

```
. summarize pre_math_score
```

| Variable | Obs | Mean | Std. Dev. | Min | Max |
|---|---|---|---|---|---|
| pre_math_s~e | 200 | 52.23 | 10.25294 | 28 | 76 |

# Run Descriptive Statistics: **tabstat**

- The **tabstat** command can be used to run descriptive statistics of your choice. Run descriptive statistics for a single variable as follows:

  - ☐ **tabstat variablename, statistics(put desired statistics here)**

```
. tabstat pre_math_score, statistics(count, mean, median, var, sd, min, max, skew, kurt)
```

| variable | N | mean | p50 | variance | sd | min | max | skewness | kurtosis |
|---|---|---|---|---|---|---|---|---|---|
| pre_math_s~e | 200 | 52.23 | 50 | 105.1227 | 10.25294 | 28 | 76 | .1948373 | 2.363052 |

# Run Descriptive Statistics

- Run descriptive statistics for multiple variables as follows:
  - **tabstat pre_math_score post_math_score diff_pre_post_math pre_reading post_reading diff_pre_post_reading , statistics(count, mean, median, var, sd, min, max, skew, kurt)**

| stats | pre_ma~e | post_m~e | diff_p~h | pre_re~g | post_r~g | diff_p~g |
|---|---|---|---|---|---|---|
| N | 200 | 200 | 200 | 200 | 200 | 200 |
| mean | 52.23 | 52.775 | .545 | 52.645 | 51.85 | -.795 |
| p50 | 50 | 54 | 0 | 52 | 53 | -1 |
| variance | 105.1227 | 89.84359 | 78.97284 | 87.76781 | 98.02764 | 68.78691 |
| sd | 10.25294 | 9.478586 | 8.886666 | 9.368448 | 9.900891 | 8.293787 |
| min | 28 | 31 | -24 | 33 | 26 | -24 |
| max | 76 | 67 | 21 | 75 | 74 | 31 |
| skewness | .1948373 | -.4784158 | -.090865 | .2844115 | -.1872277 | .3067901 |
| kurtosis | 2.363052 | 2.238527 | 2.520387 | 2.337319 | 2.428308 | 3.563576 |

When variable names are long, Stata automatically reduces the length of the variable name in the results by using a tilde (~)

# Run Descriptive Statistics

■ Descriptive statistics can be split by a categorical variable. For example, you can split the descriptive statistics by gender:

```
. tabstat diff_pre_post_math , by(gender) statistics(count, mean, median, var, sd, min, max, skew, kurt)

Summary for variables: diff_pre_post_math
    by categories of: gender
```

| gender | N | mean | p50 | variance | sd | min | max | skewness | kurtosis |
|---|---|---|---|---|---|---|---|---|---|
| female | 109 | 3.256881 | 2 | 65.7297 | 8.107385 | -16 | 21 | -.0148109 | 2.324837 |
| male | 91 | -2.703297 | -3 | 76.16654 | 8.727345 | -24 | 15 | -.0305409 | 2.439951 |
| Total | 200 | .545 | 0 | 78.97284 | 8.886666 | -24 | 21 | -.090865 | 2.520387 |

# Run Descriptive Statistics

```
. tabstat diff_pre_post_math , by(ses) statistics(count, mean, median, var, sd, min, max, skew, kurt)

Summary for variables: diff_pre_post_math
     by categories of: ses
```
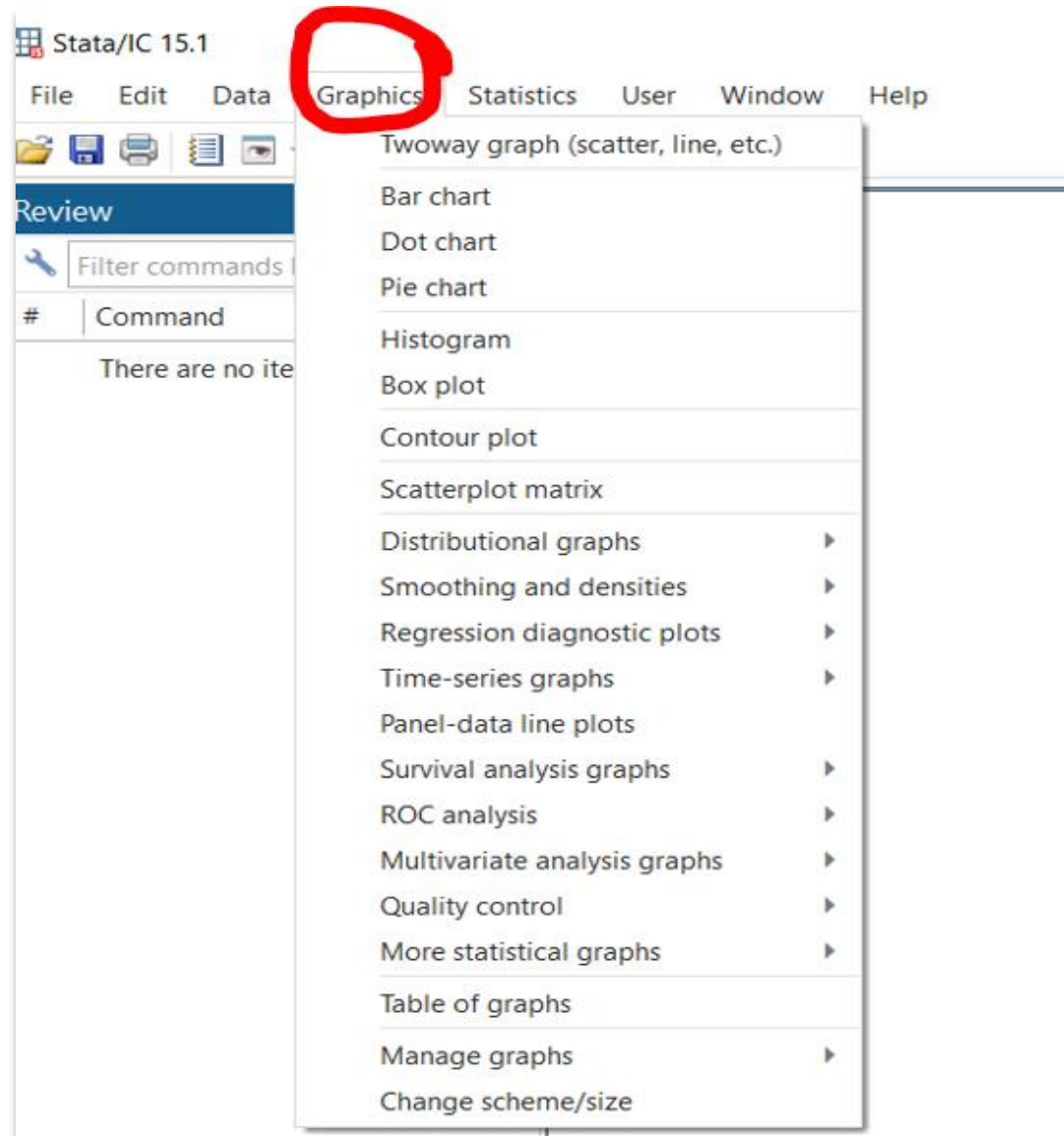
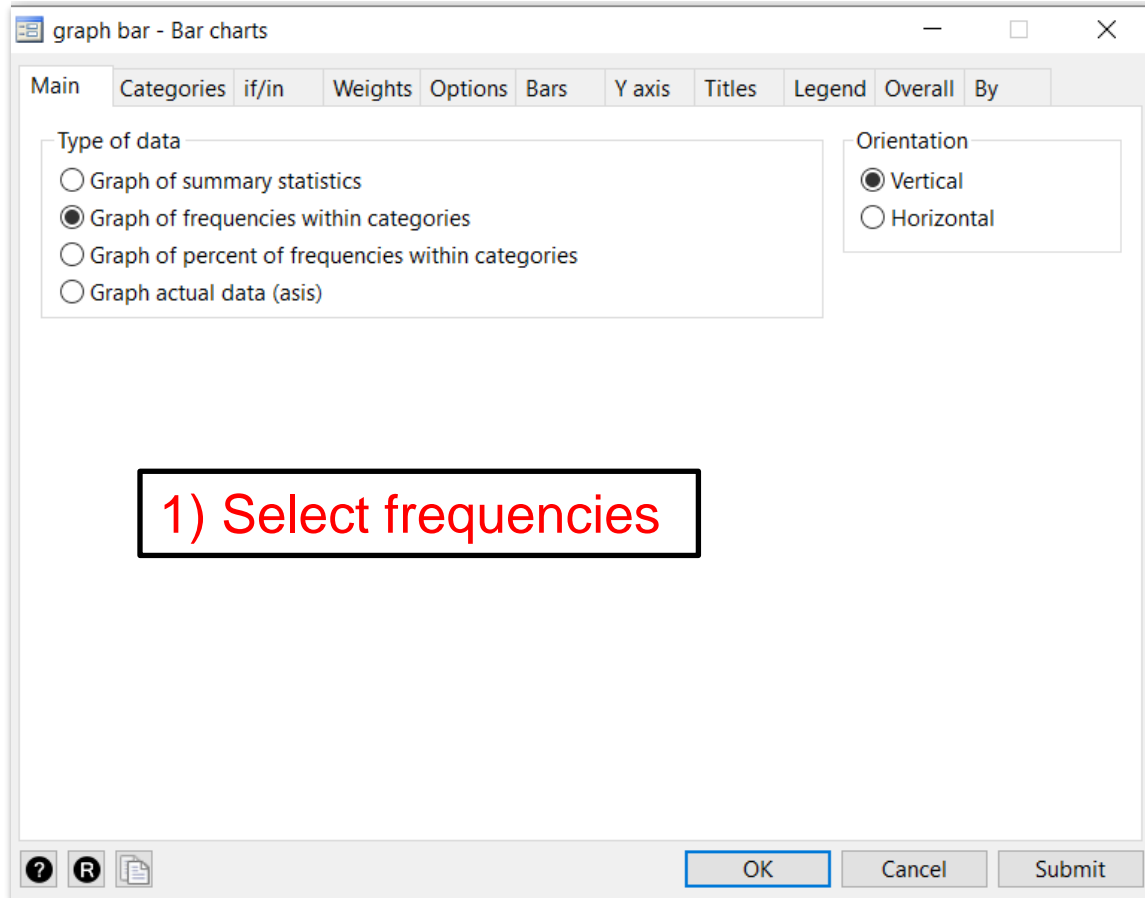| ses | N | mean | p50 | variance | sd | min | max | skewness | kurtosis |
|---|---|---|---|---|---|---|---|---|---|
| high | 58 | -.5862069 | -1 | 93.01875 | 9.644623 | -24 | 21 | .2244462 | 2.683172 |
| low | 47 | 2.340426 | 2 | 77.53377 | 8.805326 | -20 | 18 | -.3346219 | 2.690051 |
| middle | 95 | .3473684 | 1 | 70.39933 | 8.390431 | -19 | 18 | -.2013571 | 2.363308 |
| Total | 200 | .545 | 0 | 78.97284 | 8.886666 | -24 | 21 | -.090865 | 2.520387 |

# Data Visualization

# Graphics Menu

- Different types of plots can be plotted from the graphics menu.

- Click on the **Graphics** Menu, then click the desired graph, plot or chart.

# Bar Chart

**graph bar - Bar charts**

| Main | Categories | if/in | Weights | Options | Bars | Y axis | Titles | Legend | Overall | By |

☑ Group 1

Grouping variable:

[ gender ▾ ]   [ Properties ]

## 2) Enter categorical variable

☐ Group 2

Grouping variable:

[ ▾ ]   [ Properties ]

---

**graph bar - Bar charts**

| Main | Categories | if/in | Weights | Options | Bars | Y axis | Titles | Legend | Overall | By |

### Type of data

○ Graph of summary statistics
◉ Graph of frequencies within categories
○ Graph of percent of frequencies within categories
○ Graph actual data (asis)

**Orientation**
◉ Vertical
○ Horizontal

## 1) Select frequencies

[ OK ]  [ Cancel ]  [ Submit ]

---

**graph bar - Bar charts**

| Main | Categories | if/in | Weights | Options | Bars | Y axis | Titles | Legend | Overall | By |

Title:

[ A Bar Chart Showing the Frequencies of Gender Data ]   [ Properties ]

Subtitle:

[ ]   [ Properties ]

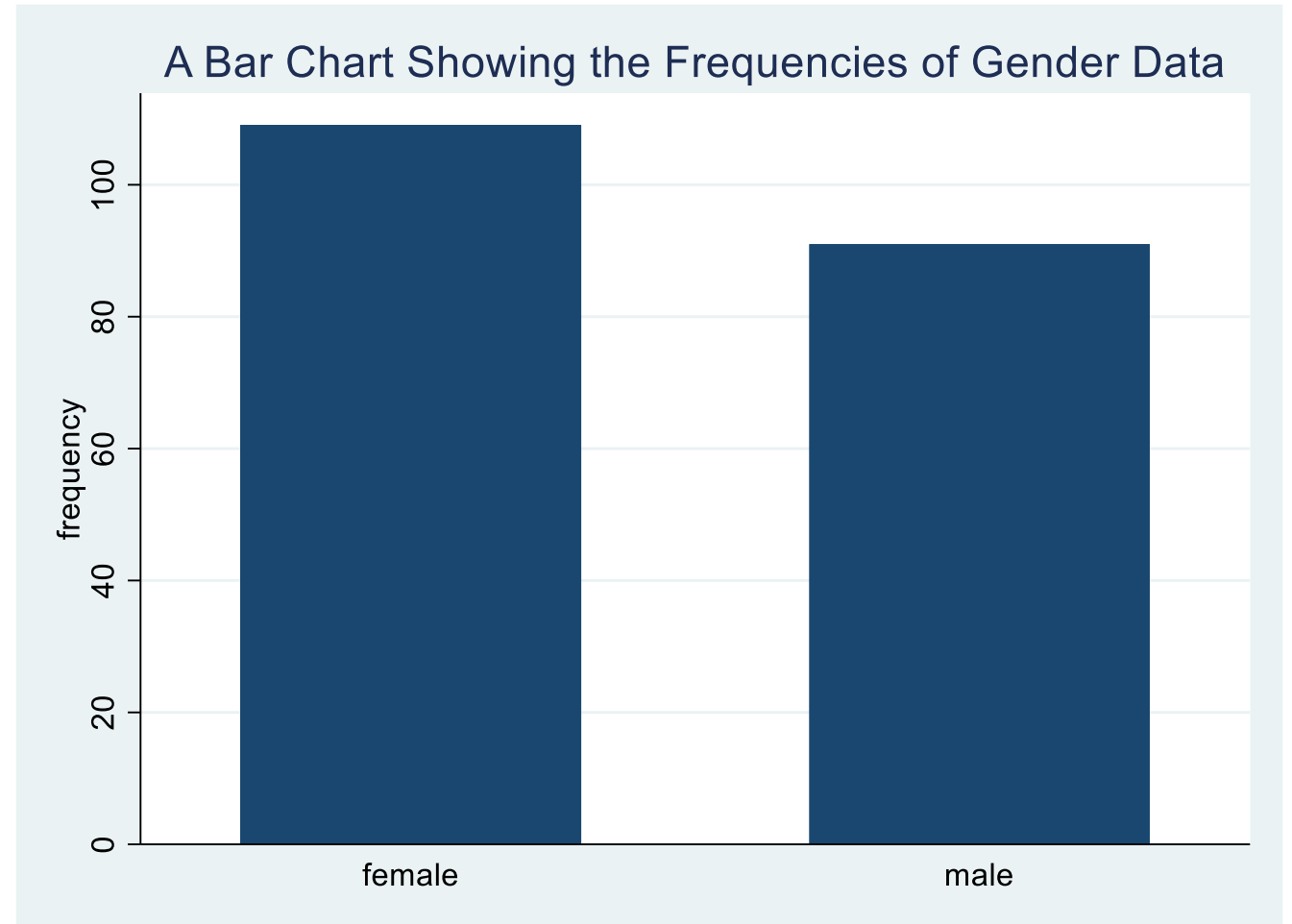Caption:

[ ]   [ Properties ]

Note:

[ ]   [ Properties ]

## 3) Include graph title, etc.

[ OK ]  [ Cancel ]  [ Submit ]

# Bar Chart

■ A bar chart is a visual display of frequencies of categorial data.



A Bar Chart Showing the Frequencies of Gender Data

syntax

```
graph bar (count), over(gender) title(A Bar Chart Showing the Frequencies of Gender Data)
```

# Histogram



histogram - Histograms for continuous and categorical variables

Main | if/in | Weights | Density plots | Add plots | Y axis | X axis | Titles | Legend | Overall | By

**Data**

Variable:
pre_math_score

◉ Data are continuous
○ Data are discrete

*Select a numerical variable*
*Click **Data are continuous***

**Bins**

☐ 10 ⏶⏷ Number of bins
☐ Width of bins
☐ Lower limit of first bin

**Y axis**

◉ Density
○ Fraction
○ Frequency
○ Percent

Bar properties

☐ Add height labels to bars
Bar label properties

☐ Recalculate bin sizes when by() is specified

OK | Cancel | Submit

---

histogram - Histograms for continuous and categorical variables

Main | if/in | Weights | Density plots | Add plots | Y axis | X axis | Titles | Legend | Overall | By

☑ Add normal-density plot      Line properties
☐ Add kernel density plot      Line properties

**Smoothing options**

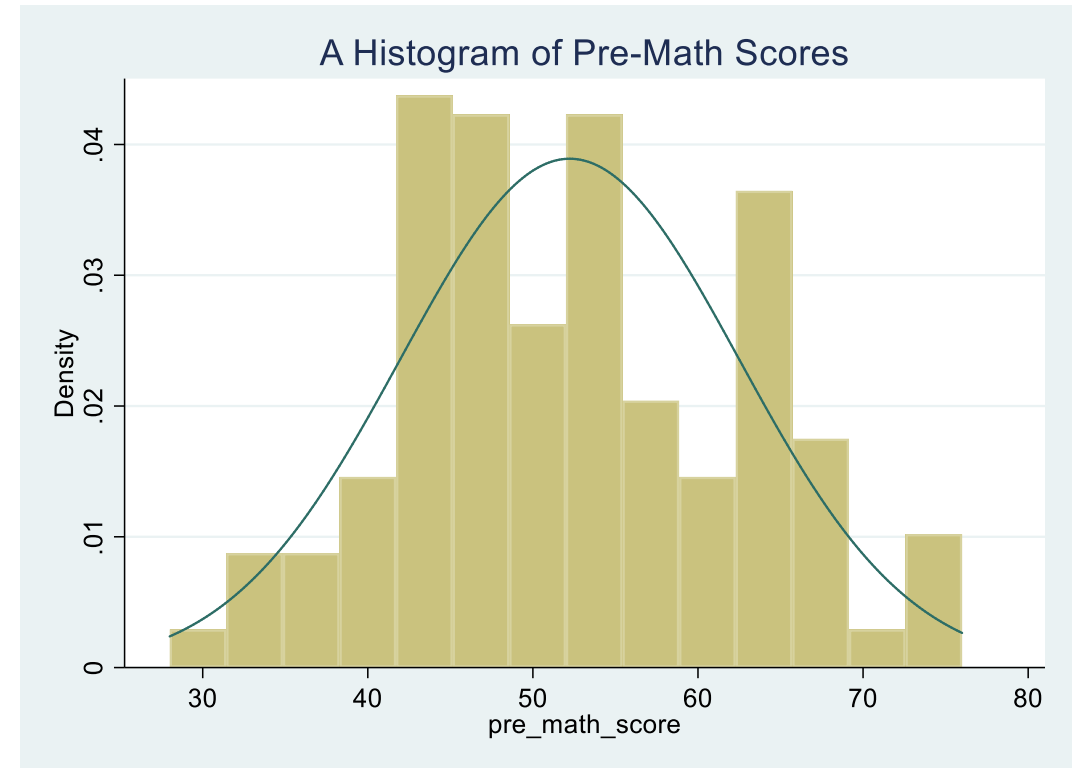☐ Override optimal width
☐ Half-width of kernel

Kernel function:
epanechnikov

*Check **Add normal-density plot to the histogram***

*Always make sure the x-axis, y-axis and title are included in to the plot*
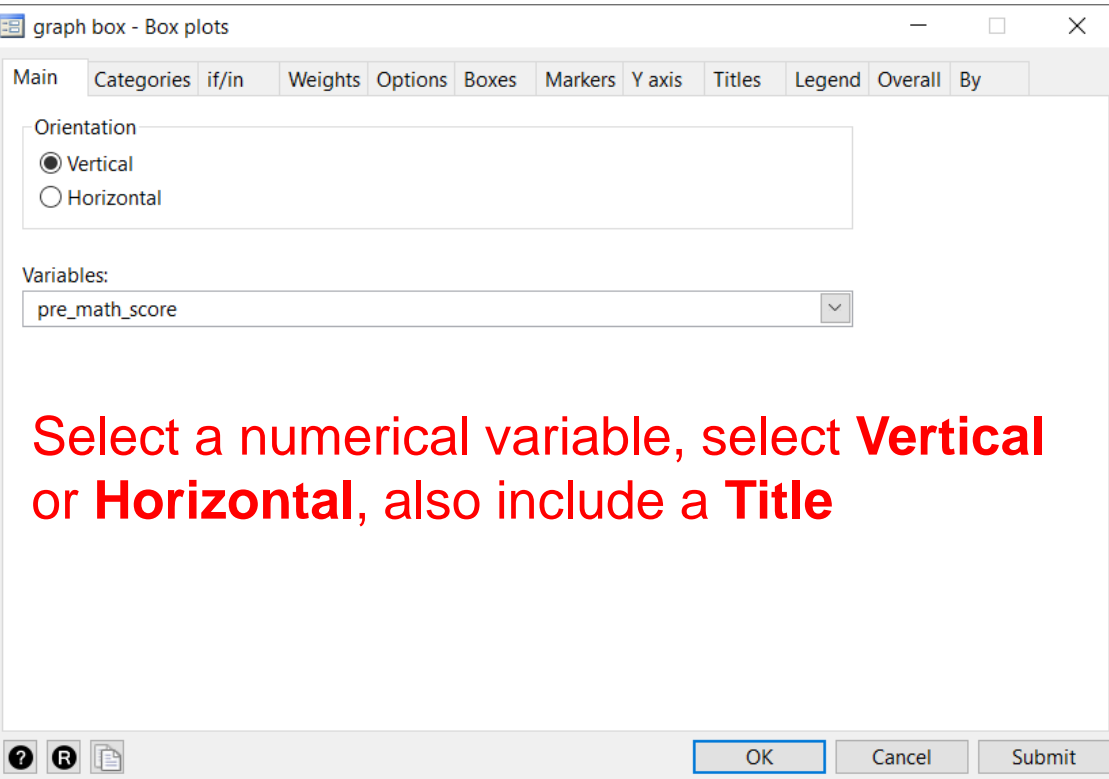
OK | Cancel | Submit

# Histogram

- The normal density line on the histogram plot helps us to see whether the plotted score is approximately normally distributed.

- Based on a visual inspection of the histogram, the post math scores appears to be nearly normally distributed.



A Histogram of Pre-Math Scores

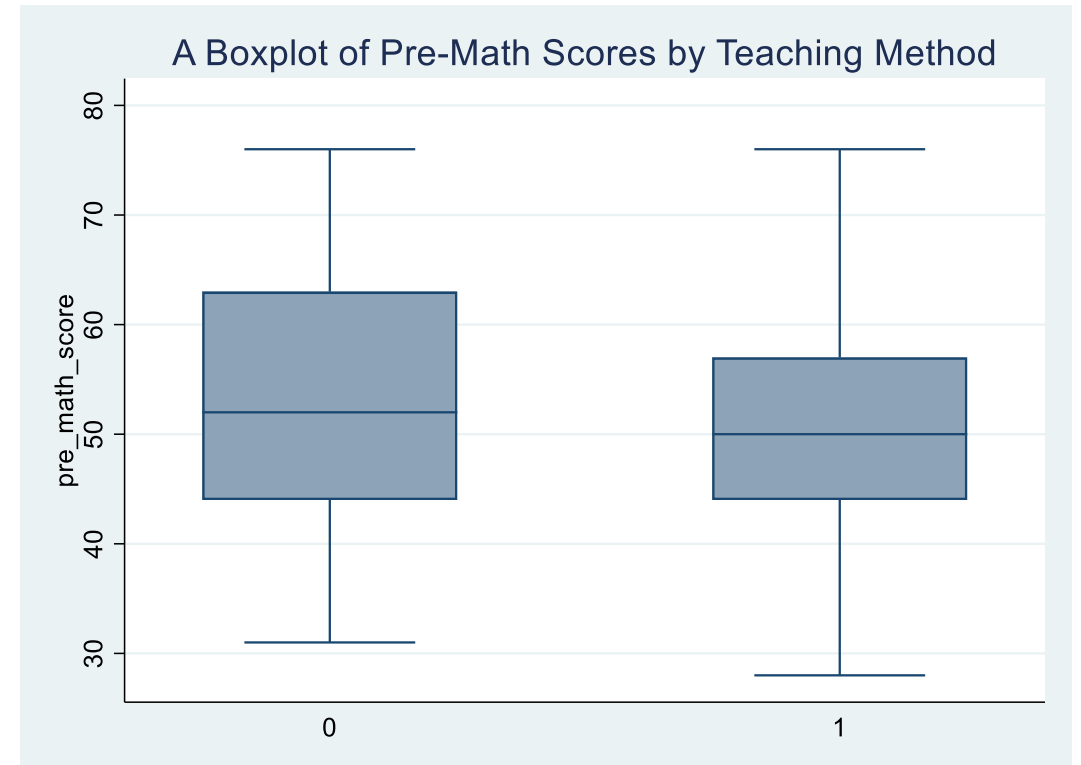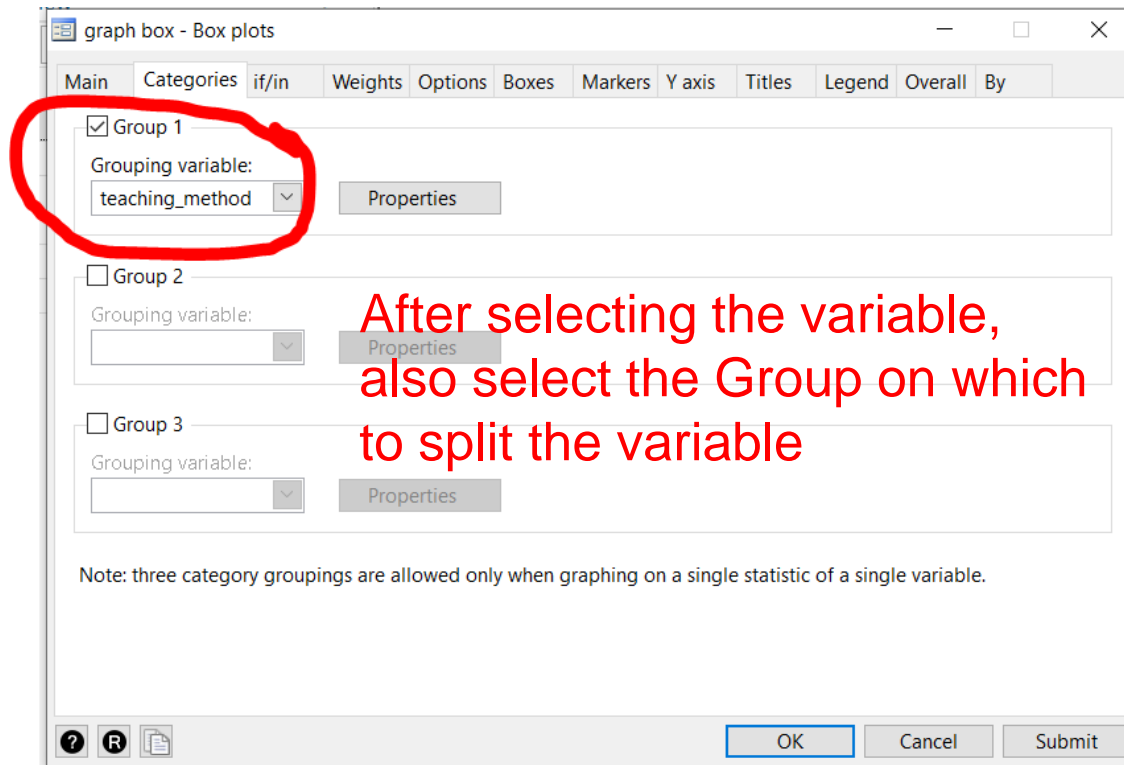syntax `histogram pre_math_score, normal title(A Histogram of Pre-Math Scores)`

# Boxplot



Select a numerical variable, select **Vertical** or **Horizontal**, also include a **Title**



syntax    `graph box pre_math_score, title(A Boxplot of Pre-Math Scores)`

# Boxplot By Category



After selecting the variable, also select the Group on which to split the variable

A Boxplot of Pre-Math Scores by Teaching Method

syntax   `graph box pre_math_score, over(teaching_method) title(A Boxplot of Pre-Math Scores by Teaching Method)`

# T-Test

- ## Research Question:
  - ☐ Is there a statistically significant difference in math score of a hybrid college class and a traditional college class?
- ## Null Hypothesis:
  - ☐ There is no statistically significant difference between the means of the math scores of a hybrid college class and a traditional college class.
- ## Alternative Hypothesis:
  - ☐ The is a statistically significant difference between the means of a hybrid college class and a traditional college class.

# T-Test

An alternative way of asking the research question and writing the hypotheses.

- Research Question:
  - Does the hybrid teaching method significantly increase the math score of college students? This is a directional
  - Does the hybrid teaching method have a significant effect on the math score of college students? This is a nondirectional (usually preferred)

- Null Hypothesis:
  - The hybrid teaching method does not have a significant effect on the math score of college students.

- Alternative Hypothesis:
  - The hybrid teaching method has a significant effect on the math score of college students.

# T-test

- Three assumptions need to be met for the results of the t-test to be valid for inference:
  - Independence
  - Normality
  - Homogeneity or equality of variances

- The independence assumption is usually assumed met.
- The normality assumption can be tested using:
  - Histogram, boxplot, skewness, and/or kurtosis
- Homogeneity of Variance test:
  - Levene's test, Bartlett's test .

# T-Test

- Note that the unit of analysis for a t-test with pre and post scores is the change in pre and post scores.

- So, we will test whether average change in pre and post scores for the hybrid class is significantly different from the average change in pre and post scores for the traditional classroom.

- Let's start by first checking the normality and homogeneity of variance assumption.

# T-Test

- Normality assumption test:
  - ☐ Check whether the pre math scores, post math scores are normality distributed for each group.
  - ☐ You could as well just check if the difference between the pre and post scores are normally distributed.
  - ☐ Let's use the histogram, boxplot, and descriptive statistics to check the normality assumption.

# T-Test

Descriptive statistics

| stats | pre_ma~e | post_m~e | diff_p~h |
|---|---|---|---|
| N | 200 | 200 | 200 |
| mean | 52.23 | 52.775 | .545 |
| p50 | 50 | 54 | 0 |
| variance | 105.1227 | 89.84359 | 78.97284 |
| sd | 10.25294 | 9.478586 | 8.886666 |
| min | 28 | 31 | -24 |
| max | 76 | 67 | 21 |
| skewness | .1948373 | -.4784158 | -.090865 |
| kurtosis | 2.363052 | 2.238527 | 2.520387 |

The descriptive statistics here is for the entire data on the outcome variable (difference in pre and post math score) and variables (pre and post math scores) used in computing the outcome.

syntax

```
tabstat pre_math_score post_math_score diff_pre_post_math , statistics(count, mean, median, var, sd, min, max, skew, kurt)
```

# T-Test

Descriptive Statistics on diff_pre_post_math by Teaching Method

| teaching_method | N | mean | p50 | variance | sd | min | max | skewness | kurtosis |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 100 | -2.38 | -3 | 74.54101 | 8.633714 | -24 | 15 | -.0547344 | 2.457884 |
| 1 | 100 | 3.47 | 2.5 | 66.91828 | 8.18036 | -16 | 21 | -.0475962 | 2.346356 |
| Total | 200 | .545 | 0 | 78.97284 | 8.886666 | -24 | 21 | -.090865 | 2.520387 |

We want to run descriptive statistics on the outcome variable for each group in the research question.
The outcome variable is the diff_pre_post_math variable which measures the improvement in math scores. The group variable is the categorical variable, teaching method

syntax
```
tabstat diff_pre_post_math , by(teaching_method) statistics(count, mean, median, var, sd, min, max, skew, kurt)
```

# T-Test

Descriptive Statistics on diff_pre_post_math by Teaching Method

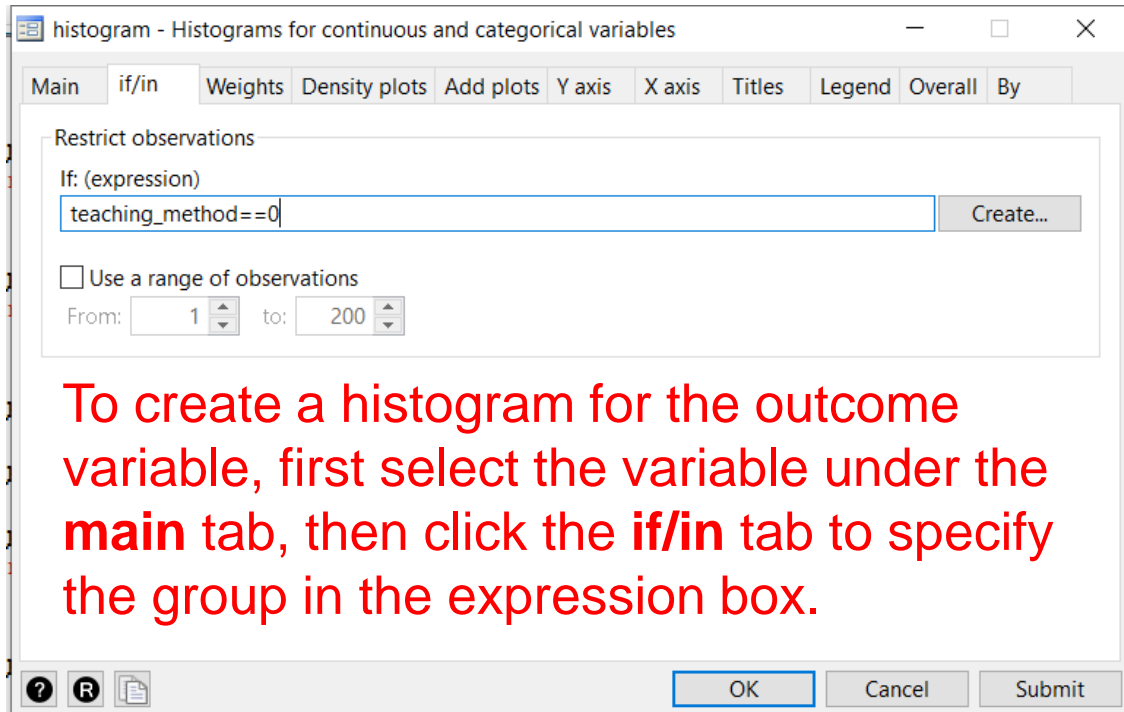| teaching_method | N | mean | p50 | variance | sd | min | max | skewness | kurtosis |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 100 | -2.38 | -3 | 74.54101 | 8.633714 | -24 | 15 | -.0547344 | 2.457884 |
| 1 | 100 | 3.47 | 2.5 | 66.91828 | 8.18036 | -16 | 21 | -.0475962 | 2.346356 |
| Total | 200 | .545 | 0 | 78.97284 | 8.886666 | -24 | 21 | -.090865 | 2.520387 |

**Interpretation of descriptive statistics:**
The skewness for the diff_pre_post_math is between -1 and 1 indicating that this outcome variable is normally distributed.
The mean and median are also approximately the same indicating normality as well.
(Note that in Stata, the median is indicated as p50 which stands for the 50[th] percentile.
You need to edit "p50" to "median" when you present your results)
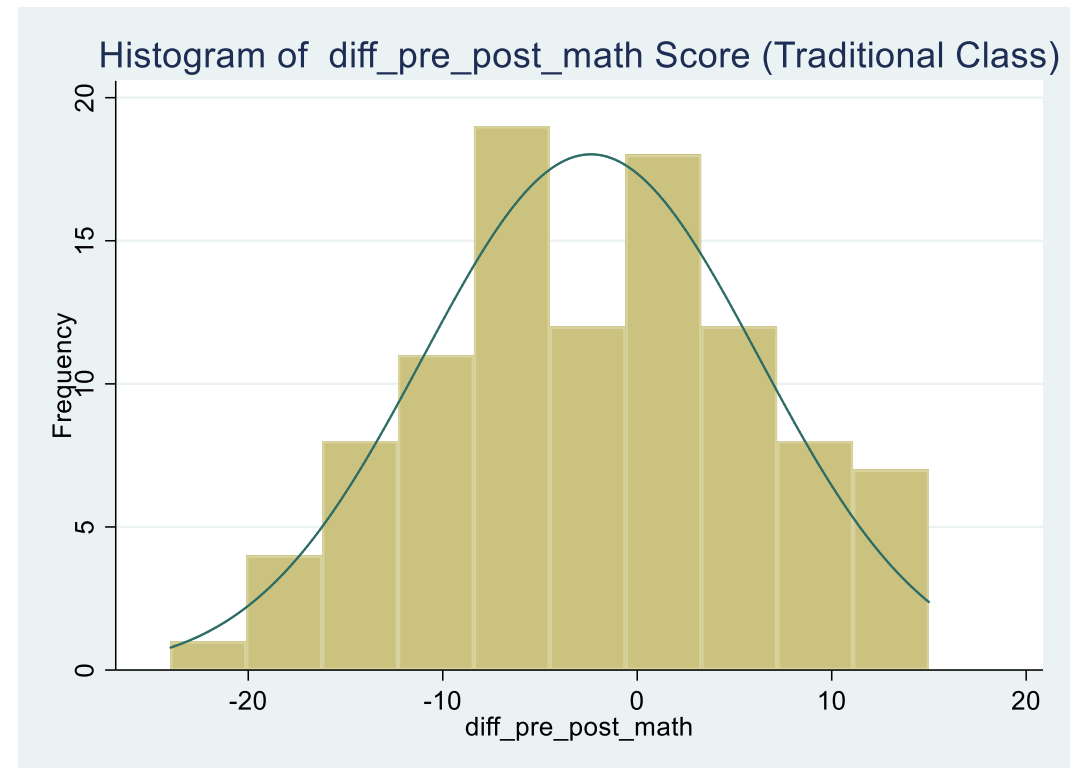
# T-Test

We also need to generate histograms for the groups (traditional and hybrid) in our research question on the outcome variable (diff_pre_post_math).



To create a histogram for the outcome variable, first select the variable under the **main** tab, then click the **if/in** tab to specify the group in the expression box.

This histogram appears to be normally distributed

syntax

```
histogram diff_pre_post_math if teaching_method==0, frequency normal title(Histogram of  diff_pre_post_math Score (Traditional Class))
```

# T-Test



Create a histogram on the outcome variable for the hybrid class ( that is, for cases with teaching_method values of 1)



This histogram appears to be normally distributed

# T-Test



For the teaching methods,
0=traditional,
1=hybrid

A Boxplot of Change in Math Scores by Teaching Method

The boxplots indicate the that the outcome is normally distributed for the hybrid and traditional class (whisker lengths are nearly equal for each boxplot).

syntax

```
graph box diff_pre_post_math, over(teaching_method) title(A Boxplot of Change in Math Scores by Teaching Method)
```

# T-Test

## Homogeneity or equality of variance tests in Stata (using the statistics tab in the menu)

**sdtest**  <span style="color:red">Use the sdtest if your data is normally distributed</span>

Statistics > Summaries, tables, and tests > Classical tests of hypotheses > Variance-comparison test

**sdtesti**

Statistics > Summaries, tables, and tests > Classical tests of hypotheses > Variance-comparison test calculator

**robvar**

Statistics > Summaries, tables, and tests > Classical tests of hypotheses > Robust equal-variance test

<span style="color:red">Use the robvar test even when your data is not normality distributed. It is robust to normality.</span>

See Stata Documentation

# T-Test

## Equality of variance tests in Stata (using syntax)

Equality of standard deviations (variances) test for v1 comparing the two groups defined by `catvar1`

`sdtest v1, by(catvar1)` Use this if your data is normal

Robust equality of variances test for v1 comparing the groups defined by `catvar1`

`robvar v1, by(catvar1)` Use this if your data is not normal

Compare the variances of v2 and v3

`sdtest v2 == v3`

Use this style of syntax only if the groups in the data were separated as two categorical variables under two different columns. We will not use this since our traditional and hybrid groups were not separated into different columns.

We would rather use the first or second style that specifies **by(categorical variable)** since the groups in our research question are both captured in a single categorical variable (teaching method).

# T-Test

## Equality of Variance Test

The p-value is greater 0.05 indicating that the null hypothesis is supported hence the variances of the groups are equal. Therefore, the equality of variance assumption is met.

syntax

```
. sdtest diff_pre_post_math, by(teaching_method)

Variance ratio test
```

| Group | Obs | Mean | Std. Err. | Std. Dev | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| 0 | 100 | -2.38 | .8633714 | 8.633714 | -4.093116 | -.6668839 |
| 1 | 100 | 3.47 | .818036 | 8.18036 | 1.846839 | 5.093161 |
| combined | 200 | .545 | .6283822 | 8.886666 | -.6941424 | 1.784142 |

```
    ratio = sd(0) / sd(1)                                    f =     1.1139
Ho: ratio = 1                            degrees of freedom =     99, 99

   Ha: ratio < 1              Ha: ratio != 1              Ha: ratio > 1
 Pr(F < f) = 0.7037        2*Pr(F > f) = 0.5925        Pr(F > f) = 0.2963
```

The null hypothesis tested is that the variances of the groups are equal.

The alternative hypothesis (Ha) is that the variances of the groups are not equal or that their ratio is 1.

# T-Test

**General syntax for running a t-test in Stata**

Test that the mean of v1 is equal between two groups defined by `catvar`

```
ttest v1, by(catvar)
```

Use this when there are two distinct groups

As above, but assume unequal variances

```
ttest v1, by(catvar) unequal
```

Use this when there are two distinct groups with unequal variances (it is equivalent to Welch's t-test)

Paired *t* test of v2 and v3

```
ttest v2 == v3
```

Use this when there is only one single group with pretest scores and posttest scores (single group pretest-posttest design)

# Research Designs Suitable for a T-test

| Treatment | Posttest Measure |
|-----------|------------------|
| $X_T$ | $O_2$ |

**One-group posttest only design:**

This design does not require a t-test as there is only one group and a single score, posttest score.
Too weak, no control for internal validity threats, don't use it for an effectiveness study.

| Pretest Measure | Treatment | Posttest Measure |
|-----------------|-----------|------------------|
| $O_1$ | $X_T$ | $O_2$ |

Compare

**One-group pretest-posttest design:**

Paired t-test can be used to compare the pretest and posttest scores for the single group.
Better than one group posttest only design but still a weak design as there is no control group.

# Research Designs Suitable for a T-test

**These are quasi experiments: no random assignments**



Treatment | Posttest

$X_T$     $O_2$

Compare

$X_C$     $O_2$

Nonequivalent comparison-group design. The dashed line indicates nonrandom assignment to comparison groups.

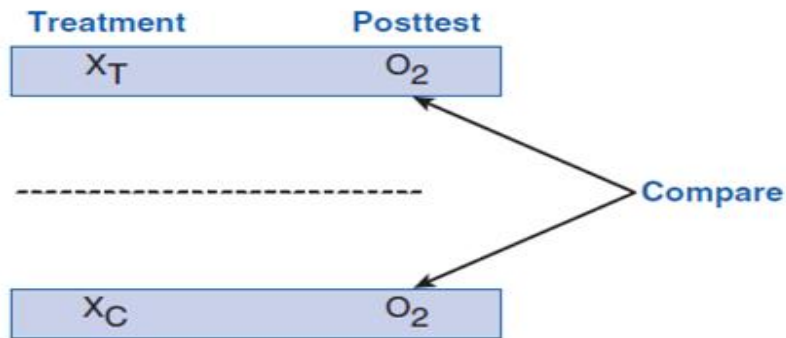| | Pretest Measure | Treatment | Posttest Measure |
|---|---|---|---|
| Experimental Group | $O_1$ | $X_1$ | $O_2$ |
| Control Group | $O_1$ | $X_2$ | $O_2$ |

Posttest-only design with nonequivalent groups:

There are two groups but since the groups are not randomly assigned, these groups are nonequivalent. Hence, this design is vulnerable to selection bias, a threat to interval validity.
Use a t-test to compare posttest scores of treatment and control group.

Nonequivalent comparison-group design:

There are two groups as well as pretest and posttest scores.
On a new column, compute the difference between pretest and posttest, then run a t-test to compare the groups on the difference score.

# Research Designs Suitable for a T-test

**These are randomized control trials (true experiments): participants were randomly assigned to groups**



**Posttest-only control group design**

Two groups are formed through random assignment. Only the posttest scores were obtained.
Use a t-test to compare the control and experimental group on the post test

**Pretest-posttest control group design**

This is a strong design and controls for threats to internal validity.
Use a t-test to compare the change or difference in pretest and posttest scores.

# Let's Run a T-test

```
. ttest diff_pre_post_math, by(teaching_method)

Two-sample t test with equal variances
```

| Group | Obs | Mean | Std. Err. | Std. Dev. | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| 0 | 100 | -2.38 | .8633714 | 8.633714 | -4.093116 | -.6668839 |
| 1 | 100 | 3.47 | .818036 | 8.18036 | 1.846839 | 5.093161 |
| combined | 200 | .545 | .6283822 | 8.886666 | -.6941424 | 1.784142 |
| diff | | -5.85 | 1.189367 | | -8.195452 | -3.504548 |

```
    diff = mean(0) - mean(1)                                  t =  -4.9186
Ho: diff = 0                              degrees of freedom =      198

    Ha: diff < 0              Ha: diff != 0              Ha: diff > 0
Pr(T < t) = 0.0000      Pr(|T| > |t|) = 0.0000      Pr(T > t) = 1.0000
```

The null hypothesis is that: the hybrid teaching method does not have significant effect on the math scores of college students.

The alternative hypothesis (Ha) is that, the hybrid teaching method has a statistical significant effect on the math scores of college students.

Since the p-value is less than 0.05, we reject the null hypothesis and support the alternative hypothesis

# Report the Results of the T-test

- Include the following as part of your report:
  - The purpose of the t-test as well as the categorical and outcome variable in the research question.
  - Report whether the assumptions are met or not

- Report whether the t-test is significant or not.
- If the t-test is significant, report the t-statistics and the degrees of freedom and also report the mean and standard deviation of the groups on the outcome variable.
- If the test is significant, draw a conclusion in a particular direction whether the treatment increased or decreased the outcome in question.

# Report the Results of the T-test (Example)

A t-test was conducted to investigate the effect of the hybrid teaching method on the math score of college students. The normality and homogeneity assumptions were tested and met. The independence assumption was assumed met.

The results of the t-test showed that the hybrid teaching method had a statistically significant effect on the math scores of college students, t (198) = -4.92, p < 0.05. On average, the hybrid class had a higher improvement in math scores, M=3.47 (0.82), compared to the traditional class, M = -2.38 (.86).

This is the mean value

The standard deviation goes into the parenthesis after reporting the mean.

# Analysis of Variance (ANOVA)

- ANOVA tries to find if three or more groups or treatment levels differ on a numeric outcome.
- Assumptions are same as those for a T-test. However, ANOVA is robust to normality.

- Run the Test and report the results following the steps already mentioned for reporting a t-test. However, note that the statistic is F-statistics and there are two degrees of freedom.
- Do a post hoc test if ANOVA results are significant and report the post hoc results as well.

# ANOVA

- There are different types of ANOVA:

  ☐ One-way ANOVA: this is an ANOVA where there is only one categorical variable and one outcome variable

  ☐ Factorial ANOVA: This is an ANOVA with two categorical variables and one outcome variable. In this type of ANOVA, we are trying to investigate if there is a significant interaction effect of the two categorical variables on the outcome as well as if there is a significant main effect of the individual categorical variables on the outcome.

In this lesson, we will focus on one-way ANOVA

# ANOVA

- Research question:
  - Do different versions of an intervention (A, B and C) have a significant difference on depression level?

- Null Hypothesis:
  - There is no significant difference between interventions A, B and C on depression level.

- Alternative Hypothesis:
  - At least, there is a significant difference between one pair of interventions on depression level.

# ANOVA

## Data used

| | id | gender | treatment | post_depre~1 |
|---|---|---|---|---|
| 1 | 1 | Male | A | 12 |
| 2 | 2 | Female | A | 21 |
| 3 | 3 | Male | A | 15 |
| 4 | 4 | Female | A | 19 |
| 5 | 5 | Male | A | 16 |
| 6 | 6 | Female | A | 18 |
| 7 | 7 | Male | A | 17 |
| 8 | 8 | Female | A | 24 |
| 9 | 9 | Male | A | 14 |
| 10 | 10 | Female | A | 25 |
| 11 | 11 | Male | B | 14 |
| 12 | 12 | Female | B | 21 |
| 13 | 13 | Male | B | 17 |
| 14 | 14 | Female | B | 20 |
| 15 | 15 | Male | B | 19 |

| | | | | |
|---|---|---|---|---|
| 16 | 16 | Female | B | 23 |
| 17 | 17 | Male | B | 20 |
| 18 | 18 | Female | B | 27 |
| 19 | 19 | Male | B | 17 |
| 20 | 20 | Female | B | 25 |
| 21 | 21 | Male | C | 25 |
| 22 | 22 | Female | C | 37 |
| 23 | 23 | Male | C | 27 |
| 24 | 24 | Female | C | 34 |
| 25 | 25 | Male | C | 29 |
| 26 | 26 | Female | C | 36 |
| 27 | 27 | Male | C | 24 |
| 28 | 28 | Female | C | 26 |
| 29 | 29 | Male | C | 22 |
| 30 | 30 | Female | C | 29 |

# ANOVA

Test for normality of outcome for each group

```
. tabstat post_depression_level, by(treatment) statistic(count mean median var sd min max skew kurt)

Summary for variables: post_depression_level
    by categories of: treatment
```
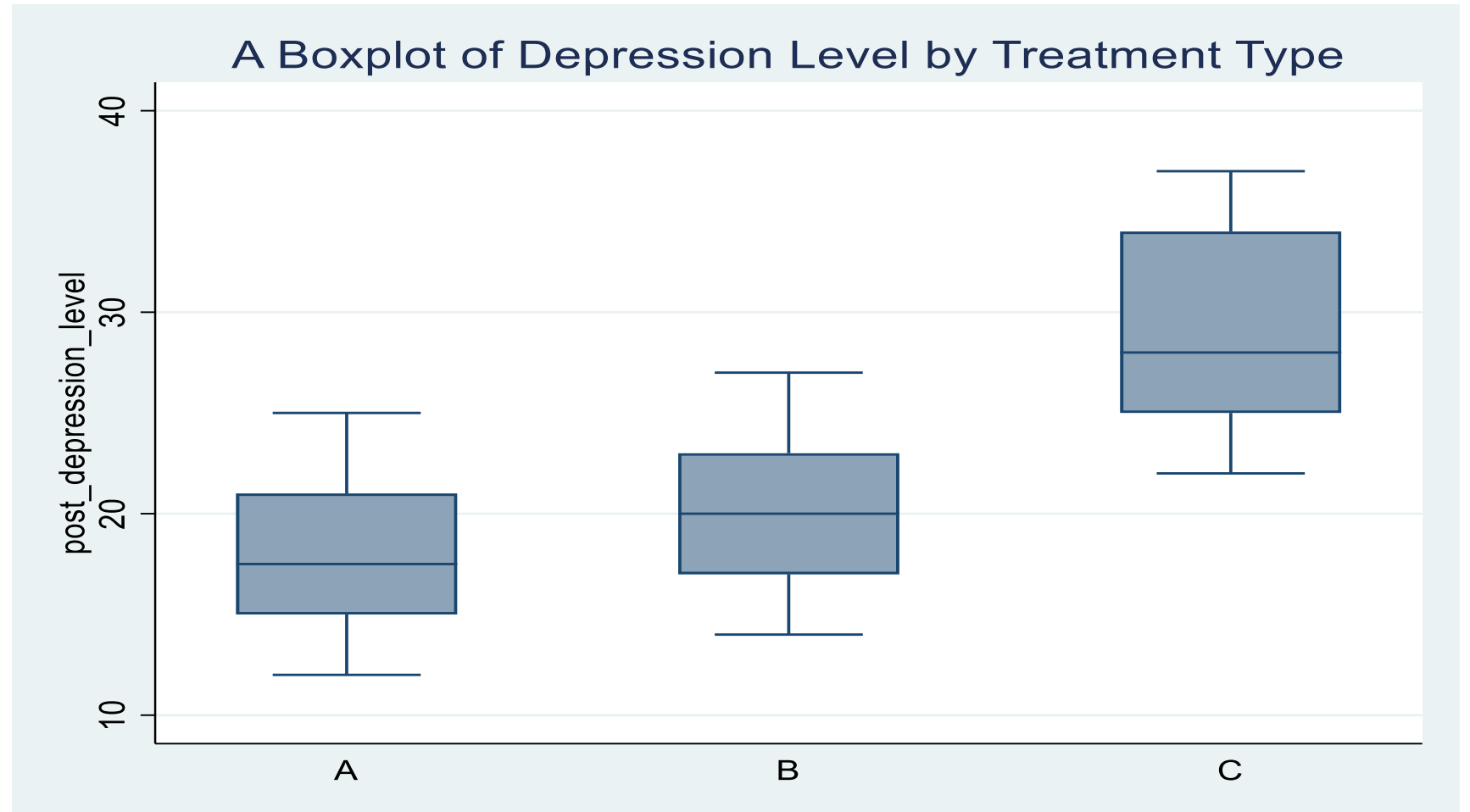
| treatment | N | mean | p50 | variance | sd | min | max | skewness | kurtosis |
|---|---|---|---|---|---|---|---|---|---|
| A | 10 | 18.1 | 17.5 | 17.87778 | 4.228212 | 12 | 25 | .3450843 | 2.058904 |
| B | 10 | 20.3 | 20 | 15.34444 | 3.917199 | 14 | 27 | .1957114 | 2.292317 |
| C | 10 | 28.9 | 28 | 26.76667 | 5.173651 | 22 | 37 | .4104282 | 1.840451 |
| Total | 30 | 22.43333 | 21.5 | 41.08161 | 6.409494 | 12 | 37 | .5727155 | 2.785366 |

Treatments A, B, and C appear to be normality distributed since their skewness values are within -1 and 1. The mean and median of each treatment group on depression level are also nearly the same indicating the depression level of each group approximately normally distributed.

# ANOVA

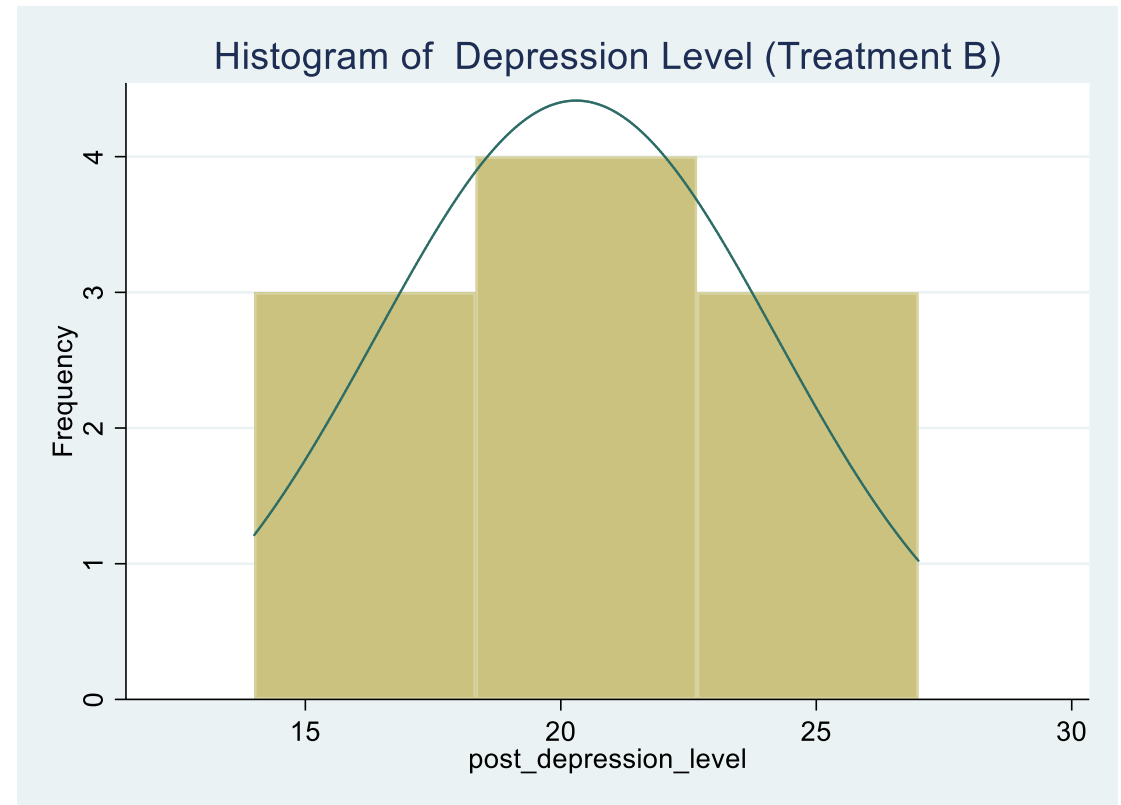The boxplots also indicate that the treatment groups are approximately normally distributed.



A Boxplot of Depression Level by Treatment Type

```
graph box post_depression_level, over(treatment) title(A Boxplot of Depression Level by Treatment Type)
```

# ANOVA



Histogram of Depression Level (Treatment A)



Histogram of Depression Level (Treatment B)

```
. histogram post_depression_level if treatment=="A", frequency normal title(Histogram of  Depression Level (Treatment A))

. histogram post_depression_level if treatment=="B", frequency normal title(Histogram of  Depression Level (Treatment B))
```
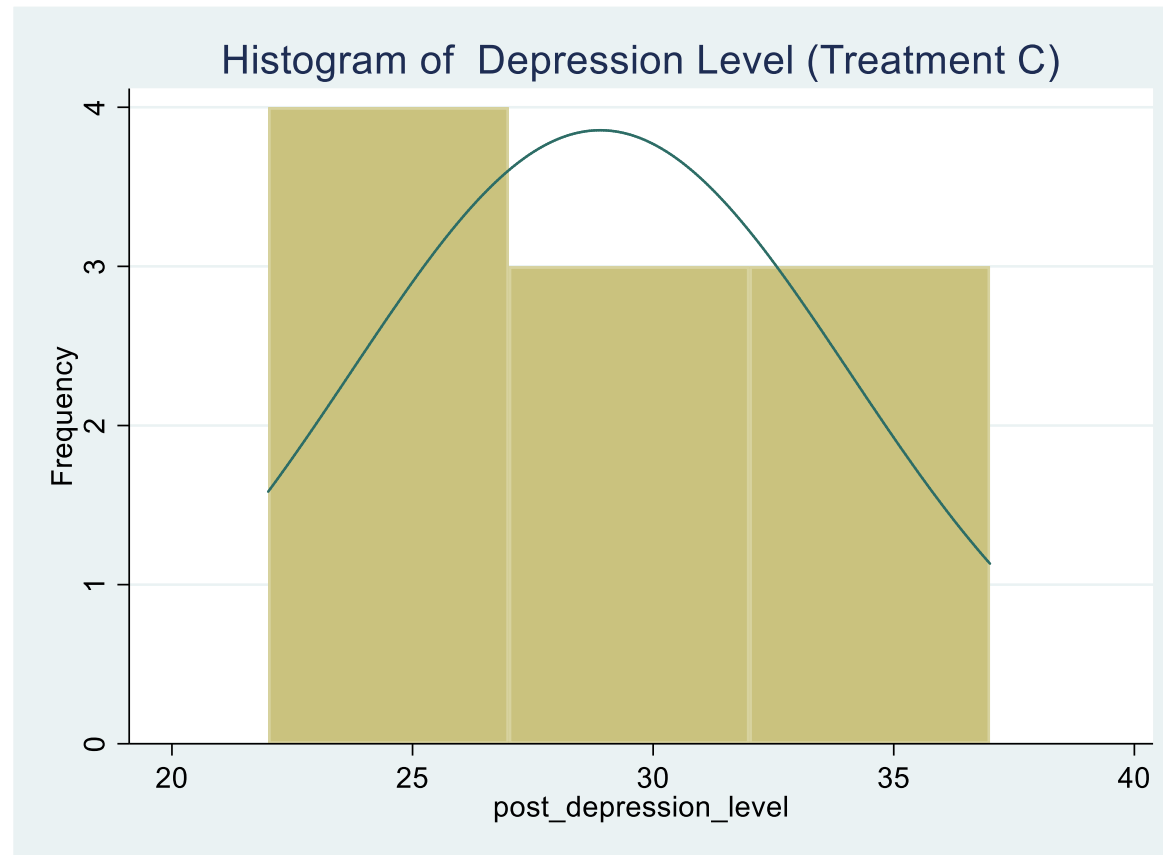
# ANOVA

Visual inspection of the histograms also indicate that the treatment groups are approximately normally distributed.

# ANOVA    Syntax for running the one-way anova in Stata

One-way ANOVA model of y for factor a

    oneway y a     Syntax for running a one way ANOVA: y is the outcome, a is the categorical variable

Report the mean and std. dev. of y and number of observations for each level of a

    oneway y a, tabulate     Include **tabulate** to generate descriptive statistics

Report all pairwise comparisons of the means of y across levels of a with $p$-values adjusted using Bonferroni's procedure     If you ANOVA test is significant, run it again and add **bonferroni** to get post hoc results with Bonferroni adjustment to p-value

    oneway y a, bonferroni

As above, but adjust $p$-values for multiple comparisons using Scheffé's method

    oneway y a, scheffe     Instead of using **bonferroni,** you could alternatively use **scheffe** for the post hoc multiple-comparison post hoc test. The Bonferroni and scheffe methods correct for type 1 error when doing pairwise multiple comparison.

# ANOVA

## ANOVA Results

The p-value of the F statistics is less than 0.05 showing that we need to reject the null hypothesis and support the alternative hypothesis. Therefore, there is a significant different between at least two treatment types.

```
. oneway post_depression_level treatment
```

Degrees of freedom

                    Analysis of Variance
| Source | SS | df | MS | F | Prob > F |
|---|---|---|---|---|---|
| Between groups | 651.466667 | 2 | 325.733333 | 16.29 | 0.0000 |
| Within groups | 539.9 | 27 | 19.9962963 | | |
| Total | 1191.36667 | 29 | 41.0816092 | | |

Bartlett's test for equal variances:   chi2(2) =   0.7305   Prob>chi2 = 0.694

In Stata, we do not have a separate test for equality or homogeneity of variance for ANOVA. The test ANOVA results have the homogeneity of variance test results included. A Bartlett's test is used to test for equal variances.

# ANOVA — ANOVA Results including Post Hoc Test Results

```
. oneway post_depression_level treatment, bonferroni

                        Analysis of Variance
    Source              SS          df      MS              F      Prob > F
─────────────────────────────────────────────────────────────────────────
Between groups       651.466667      2   325.733333       16.29     0.0000
 Within groups          539.9       27   19.9962963
─────────────────────────────────────────────────────────────────────────
    Total            1191.36667      29   41.0816092

Bartlett's test for equal variances:   chi2(2) =    0.7305  Prob>chi2 = 0.694

                Comparison of post_depre~l by treatment
                            (Bonferroni)
Row Mean-
Col Mean          A               B

    B            2.2
               0.843

    C           10.8             8.6
               0.000            0.001

.
```

Test statistics

P-value

Since the one-way ANOVA test results were significant, we need to test which pair of groups are significantly different using the post hoc test.
We run the post hoc test by adding **bonferroni** to the one-way ANOVA syntax.

There is a significant difference between treatment C and A. There is also a significant difference between treatment C and A since the p-value is less than 0.05

# ANOVA Test Results

An analysis of variance (ANOVA) test was conducted to examine whether there was a statistically significant difference between treatment types (A, B, and C) on depression level. The normality and homogeneity of variance assumptions were tested and met. The independence assumption was assumed met.

The results of the ANOVA test indicated that there was a statistically significant difference between the treatment types, $F(2, 27) = 16.29$, $p < 0.05$. A pairwise multiple-comparison test was conducted to investigate which pairs of treatment types were significantly different. The results of the post hoc test indicated that the group with treatment B had lower depression level, $M = 20.3$ (3.91) compared to the group that received treatment C, $M = 28.9$ (5.17), and the group with treatment A had a significantly lower depression level, $M = 18.1$ (4.23), than group with treatment C.