Samin Naji - 9812399057

1. Import csv data:

2. Scale variables are usually measured at the interval or ratio level of measurement. Therefore, the suitable descriptive statistics table for scale variables in SPSS is the one that includes the mean, standard deviation, variance, range, minimum, maximum, and optionally the geometric mean, harmonic mean, and coefficient of variation.

Interval: This level of measurement measures the amount or degree of something, such as temperature, IQ, or income. The numbers have a meaningful order and equal intervals, but there is no true zero point. The suitable descriptive statistics for interval variables are mean, standard deviation, variance, range, minimum, maximum, and parametric tests.

Ratio: This level of measurement is similar to interval, but it has a true zero point, such as height, weight, or age. The numbers have a meaningful order, equal intervals, and a meaningful ratio. The suitable descriptive statistics for ratio variables are the same as interval variables, plus geometric mean, harmonic mean, and coefficient of variation.

Here we can see our measures and we can recognize which descriptive statistics table is suitable for our variables:

for scale measures, we have:

1.



Nominal: This level of measurement assigns numbers to categories or groups, such as gender, race, or religion. The numbers do not have any inherent order or meaning. The suitable descriptive statistics for nominal variables are frequency, percentage, mode, and chi-square test. So for nominal measures, we have:

2.



cyl

| | | Frequency | Percent | Valid Percent | Cumulative Percent |
|---|---|---|---|---|---|
| Valid | 4 | 11 | 34.4 | 34.4 | 34.4 |
| | 6 | 7 | 21.9 | 21.9 | 56.3 |
| | 8 | 14 | 43.8 | 43.8 | 100.0 |
| | Total | 32 | 100.0 | 100.0 | |

vs

| | | Frequency | Percent | Valid Percent | Cumulative Percent |
|---|---|---|---|---|---|
| Valid | 0 | 18 | 56.3 | 56.3 | 56.3 |
| | 1 | 14 | 43.8 | 43.8 | 100.0 |
| | Total | 32 | 100.0 | 100.0 | |

am

| | | Frequency | Percent | Valid Percent | Cumulative Percent |
|---|---|---|---|---|---|
| Valid | 0 | 19 | 59.4 | 59.4 | 59.4 |
| | 1 | 13 | 40.6 | 40.6 | 100.0 |
| | Total | 32 | 100.0 | 100.0 | |

gear

3.



| | | Frequency | Percent | Valid Percent | Cumulative Percent |
|---|---|---|---|---|---|
| Valid | 0 | 19 | 59.4 | 59.4 | 59.4 |
| | 1 | 13 | 40.6 | 40.6 | 100.0 |
| | Total | 32 | 100.0 | 100.0 | |

gear

| | | Frequency | Percent | Valid Percent | Cumulative Percent |
|---|---|---|---|---|---|
| Valid | 3 | 15 | 46.9 | 46.9 | 46.9 |
| | 4 | 12 | 37.5 | 37.5 | 84.4 |
| | 5 | 5 | 15.6 | 15.6 | 100.0 |
| | Total | 32 | 100.0 | 100.0 | |

carb

| | | Frequency | Percent | Valid Percent | Cumulative Percent |
|---|---|---|---|---|---|
| Valid | 1 | 7 | 21.9 | 21.9 | 21.9 |
| | 2 | 10 | 31.3 | 31.3 | 53.1 |
| | 3 | 3 | 9.4 | 9.4 | 62.5 |
| | 4 | 10 | 31.3 | 31.3 | 93.8 |
| | 6 | 1 | 3.1 | 3.1 | 96.9 |
| | 8 | 1 | 3.1 | 3.1 | 100.0 |
| | Total | 32 | 100.0 | 100.0 | |

3. Scale variables are usually measured at the interval or ratio level of measurement. Therefore, the suitable graphs for scale variables in SPSS are the ones that can show the distribution, outliers, or relationship of your variables, such as histograms, boxplots, or scatterplots.

   For interval or ratio variables, we can use histograms, boxplots, or scatterplots to show the distribution, outliers, or relationship of your variables.

   Histograms show the frequency of values in intervals or bins, and you can also overlay a normal curve to check the normality of your data.

   Boxplots show the median, quartiles, and extreme values of your variables, and they can help you identify outliers or skewness.

   Scatterplots show the relationship between two variables by plotting them as dots on a coordinate plane.

   1.

2.



3.

4.



For nominal or ordinal variables, we can use bar graphs, pie charts, or histograms to show the frequency or percentage of each category.
Bar graphs are useful for comparing multiple categories or groups, while
pie charts are good for showing the proportion of each category in the whole.
Histograms are similar to bar graphs, but they show the distribution of a single variable in intervals or bins.

5.



6.

7.



8.

9.

4.

This table shows the frequencies of different variables in a dataset. The table is divided into columns for the variable name, the number of valid cases, the number of missing cases, and the mean, median, mode, standard deviation, variance, skewness, and kurtosis for each variable. This table can help researchers understand the distribution of their data and identify any potential outliers or patterns.
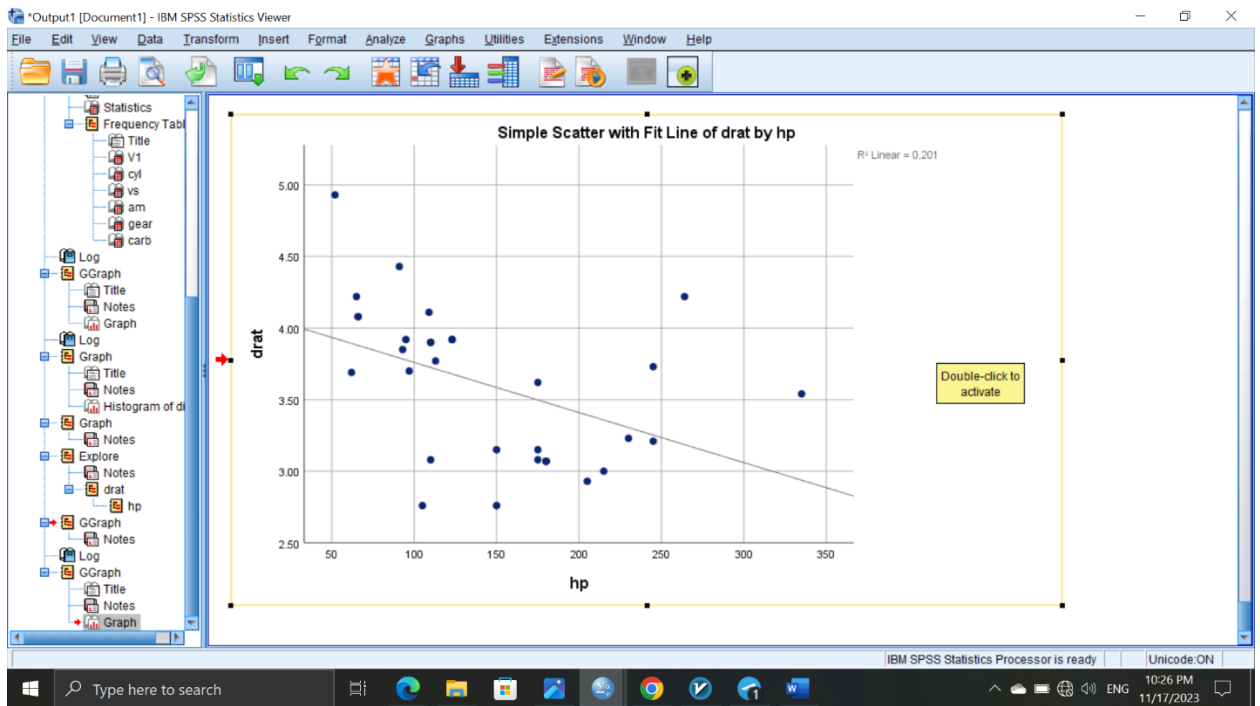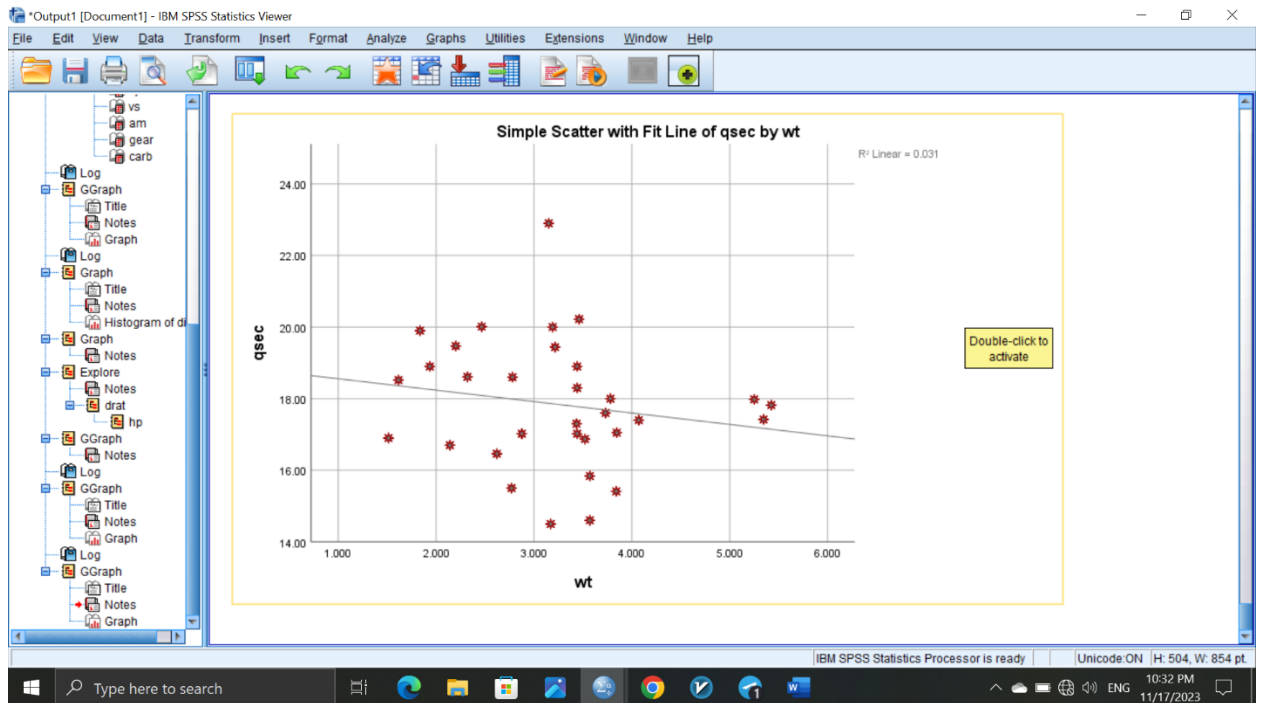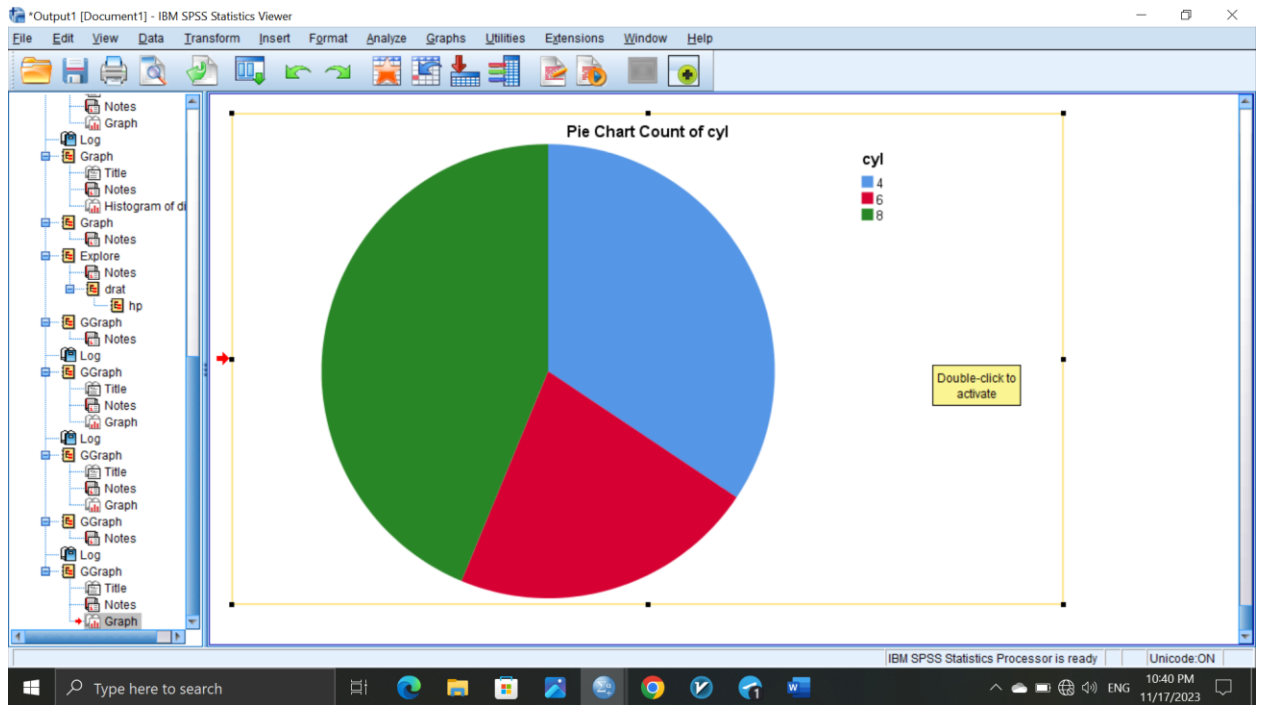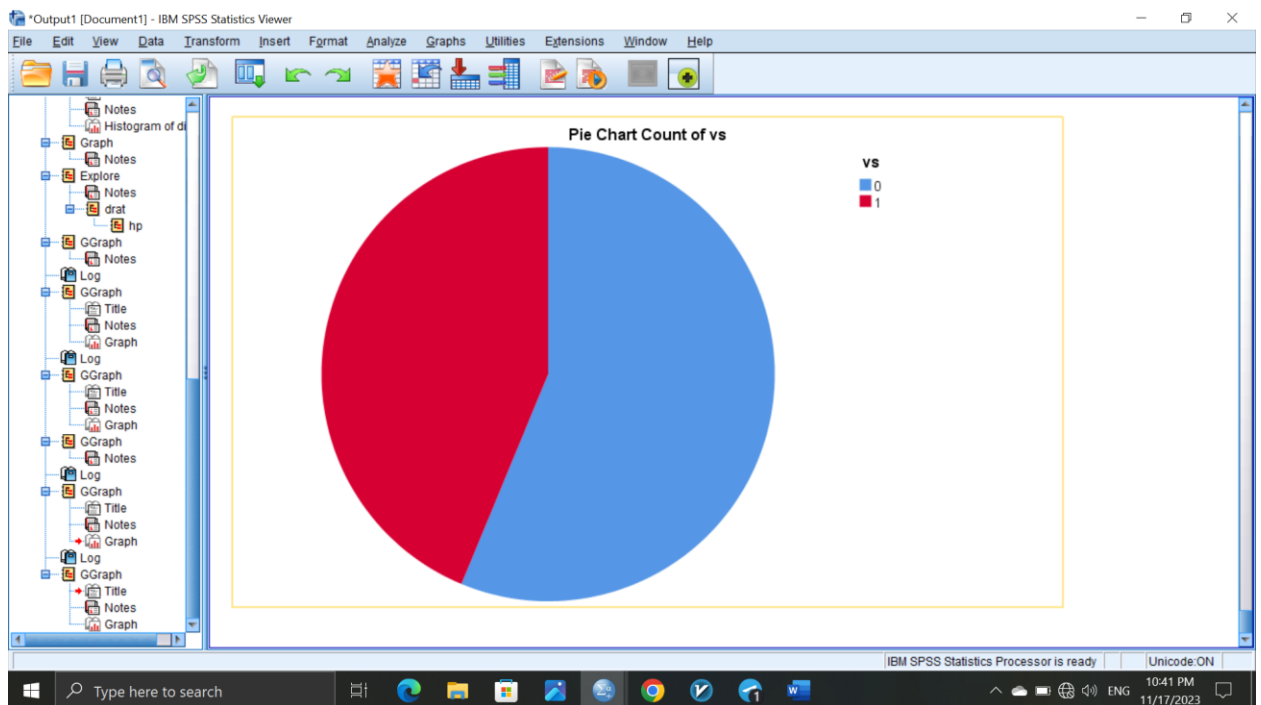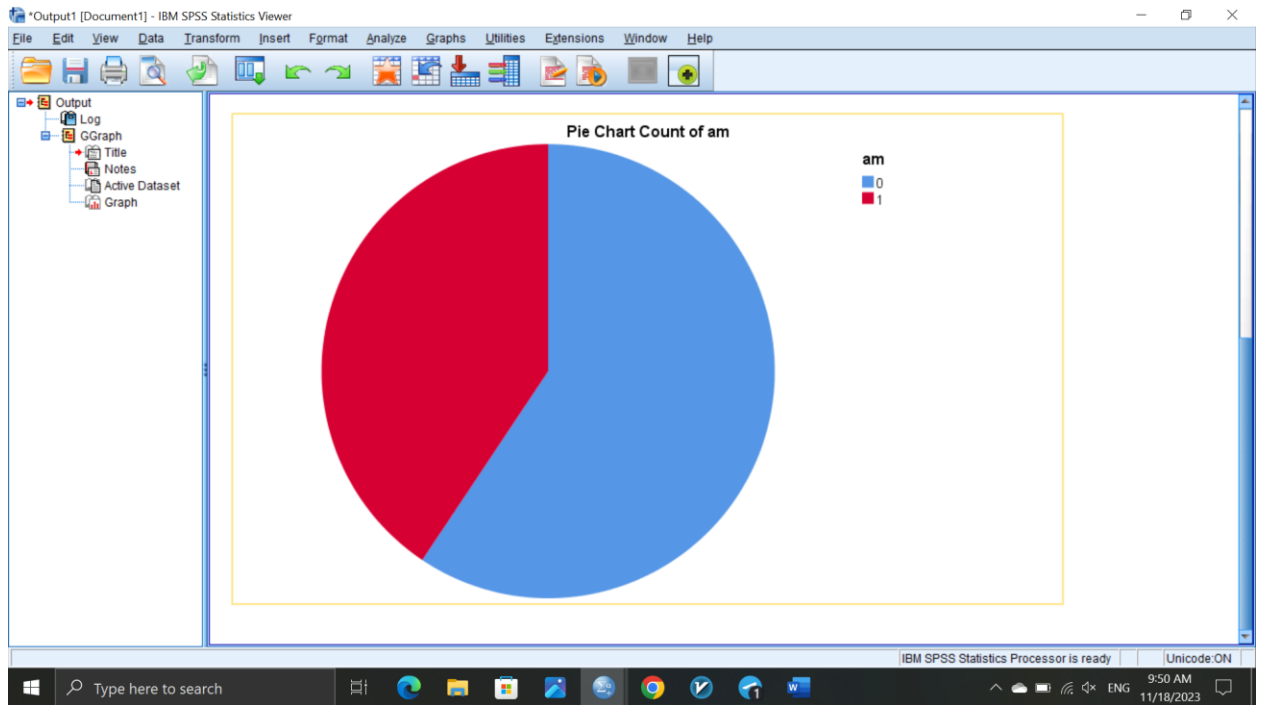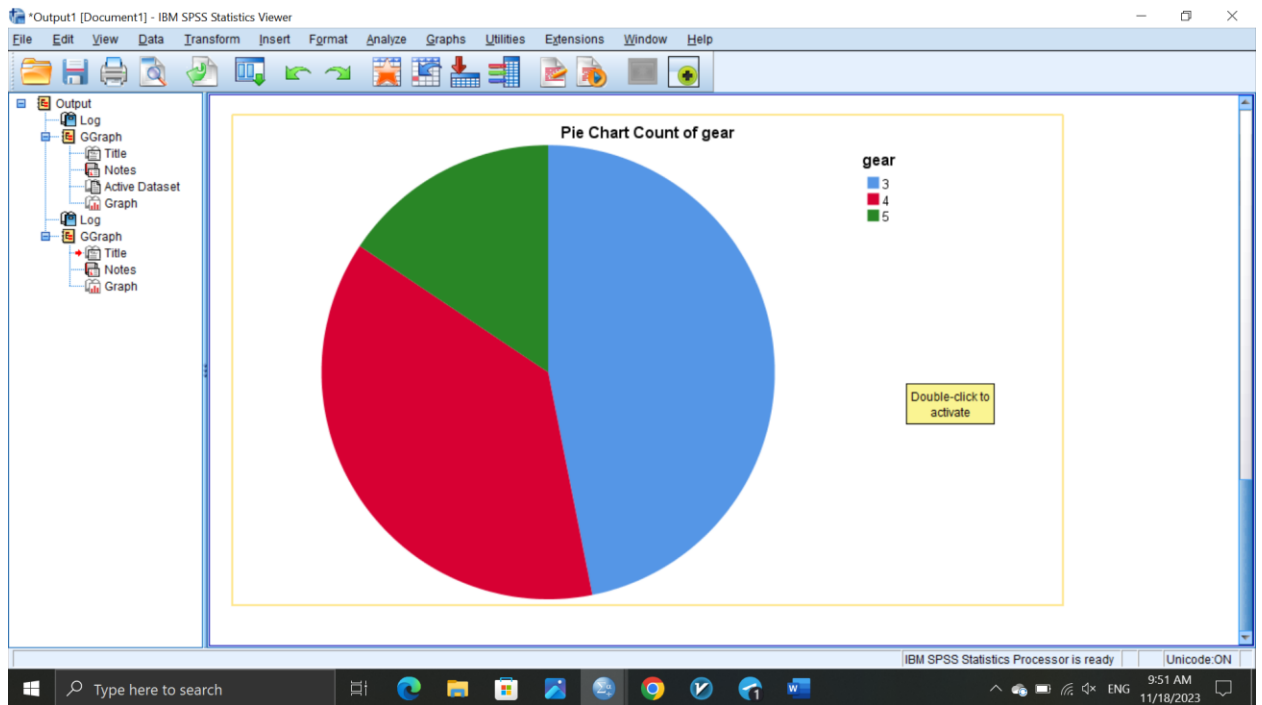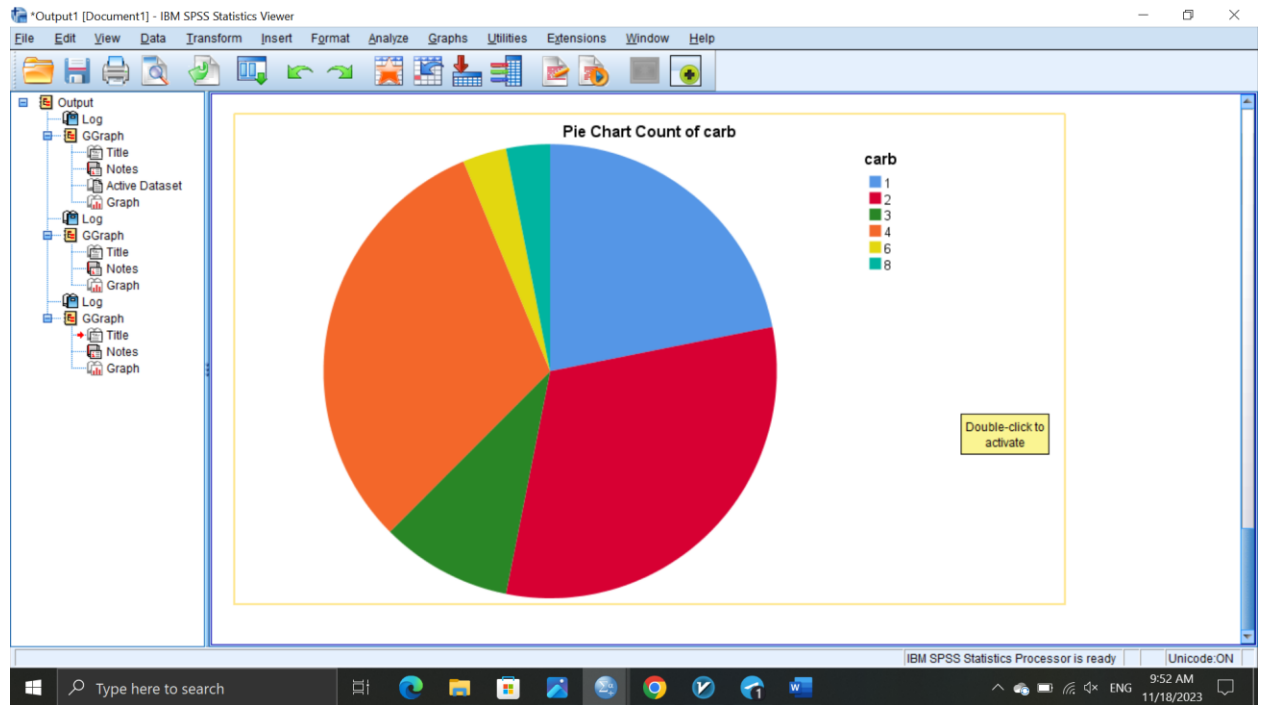
Let's explain what each column means:

- **Variable**: This is the name of the variable in the dataset.
- **Valid**: This is the number of cases that have a valid value for the variable. For example, if there are 100 cases in the dataset, and 10 of them have a missing value for age, then the valid number for age is 90.
- **Missing**: This is the number of cases that have a missing value for the variable. For example, if there are 100 cases in the dataset, and 10 of them have a missing value for age, then the missing number for age is 10.
- **Mean**: This is the average value of the variable. It is calculated by adding up all the valid values and dividing by the number of valid cases. For example, the mean of the hp variable is 146.69, which means that the average horsepower of the cases in the dataset is 146.69.
- **Median**: This is the middle value of the variable. It is calculated by sorting the valid values in ascending order and finding the value that is in the middle. If there is an even number of valid values, then the median is the average of the two middle values.
- **Mode**: This is the most frequent value of the variable. It is calculated by counting how many times each valid value occurs and finding the value that occurs the most. If there is more than one value that has the same frequency, then the mode is the lowest of those values.
- **Std. Deviation**: This is the standard deviation of the variable. It is a measure of how much the valid values vary from the mean. It is calculated by finding the difference between each valid value and the mean, squaring those differences, adding them up, dividing by the number of valid cases, and taking the square root.
- **Variance**: This is the variance of the variable. It is the square of the standard deviation. It is a measure of how much the valid values vary from the mean. It is calculated by finding the difference between each valid value and the mean, squaring those differences, adding them up, and dividing by the number of valid cases.
- **Skewness**: This is the skewness of the variable. It is a measure of how symmetric or asymmetric the distribution of the valid values is. It is calculated by finding the difference between each valid value and the mean, cubing those differences, adding them up, dividing by the number of valid cases, and dividing by the cube of the standard deviation.
  A positive skewness means that the distribution is skewed to the right, meaning that there are more values above the mean than below. A negative skewness means that the distribution is skewed to the left, meaning that there are more values below the mean than above. A zero skewness means that the distribution is symmetric, meaning that there are equal values above and below the mean.

- **Kurtosis**: This is the kurtosis of the variable. It is a measure of how peaked or flat the distribution of the valid values is. It is calculated by finding the difference between each valid value and the mean, raising those differences to the fourth power, adding them up, dividing by the number of valid cases, and dividing by the fourth power of the standard deviation.
A positive kurtosis means that the distribution is more peaked than a normal distribution, meaning that there are more values close to the mean and fewer values far from the mean. A negative kurtosis means that the distribution is more flat than a normal distribution, meaning that there are fewer values close to the mean and more values far from the mean. A zero kurtosis means that the distribution is similar to a normal distribution, meaning that there is a balance between values close to and far from the mean.

For the second table:

This table shows the frequency distribution of two categorical variables in the dataset: "cyl" and "am". The "cyl" variable stands for the number of cylinders in a car engine, and the "am" variable stands for the type of transmission: 0 for automatic and 1 for manual. The table is divided into columns for the valid number of cases, the frequency of each category, and the percentage of each category. The table can help researchers understand the distribution of their data and compare the proportions of different categories.

Let's explain what each column means:

- **Valid**: This is the number of cases that have a valid value for the variable. For example, if there are 30 cases in the dataset, and none of them have a missing value for "cyl", then the valid number for "cyl" is 30.
- **Frequency**: This is the number of cases that belong to each category of the variable. For example, if there are 30 cases in the dataset, and 19 of them have a value of 6 for "cyl", then the frequency of 6 for "cyl" is 19.
- **Percent**: This is the percentage of cases that belong to each category of the variable. It is calculated by dividing the frequency by the valid number and multiplying by 100. For example, if there are 30 cases in the dataset, and 19 of them have a value of 6 for "cyl", then the percentage of 6 for "cyl" is (19 / 30) * 100 = 63.3.
- **Cumulative Percent**: This is the percentage of cases that have a value less than or equal to a specific value for the variable. It is calculated by adding up the percentages of all the values that are less than or equal to the current value.

For the third table:

This table is similar to the second table. The table is divided into columns for the valid number of cases, the percent of each value, and the cumulative percent of each value. The table can help researchers understand the distribution of their data and compare the proportions of different values. I explained valid, percent and cumulative percent cases above.

This plot is a simple bar graph that shows the frequency of different values of mpg (miles per gallon) in a dataset. The x-axis represents the mpg values and the y-axis represents the frequency. From the plot, we can see that the most frequent mpg value is around 150, and the frequency decreases as the mpg value increases. This means that most of the cars in the dataset have a low fuel efficiency, and only a few have a high fuel efficiency. The plot can help us understand the distribution of the mpg variable and compare the performance of different cars.

For the second plot:

This plot is a histogram in SPSS, which is a statistical software. A histogram is a graphical display of the frequency distribution of a numerical variable. The x-axis represents the values of the variable and the y-axis represents the frequency of each value. The plot shows the frequency distribution of the variable "disp" in the dataset. The "disp" variable stands for the displacement of a car engine, which is a measure of its size and power. The mean of the variable "disp" is 230.72 and the standard deviation is 123.93. These statistics are shown on the plot as well.

From the plot, we can see that the distribution of the variable "disp" is skewed to the right, meaning that there are more values below the mean than above. The most frequent value of "disp" is around 120, and the frequency decreases as the value increases. There are also some outliers or extreme values of "disp" that are above 400. The plot can help us understand the distribution of the variable "disp" in the dataset and see if it is normally distributed or not. A normal distribution is a symmetric and bell-shaped distribution that has many applications in statistics. If the variable "disp" is not normally distributed, we may need to transform it or use nonparametric tests for further analysis.

For the third plot:

The plot is a scatter plot with a line of best fit in SPSS, which is a statistical software. A scatter plot is a graphical display of the relationship between two numerical variables. The x-axis represents the values of one variable and the y-axis represents the values of another variable. The line of best fit is a straight line that shows the trend of the relationship between the two variables. It is calculated by minimizing the sum of squared errors between the observed data points and the predicted values on the line.

From the plot, we can see that there is a negative relationship between the two variables, hp and drat. The hp variable stands for horsepower, which is a unit of power that measures how much work a car engine can do in a given time. The drat variable stands for the rear axle ratio, which is a measure of how fast a car can accelerate. The line of best fit suggests that as hp increases, drat decreases. This means that cars with higher horsepower tend to have lower rear axle ratios, and

vice versa. The plot can help us understand the correlation and the regression between the two variables and see how well the line of best fit fits the data.

The R2 value also appears in the top right corner of the plot. This represents the coefficient of determination, which is a measure of how much variation in the response variable "drat" can be explained by the predictor variable "hp". It ranges from 0 to 1, with higher values indicating a better fit. In this case, the R2 value is 0.201, which means that 20.1% of the variation in drat can be explained by hp. This indicates a weak linear relationship between the two variables.

### For the fourth plot:

From the plot, we can see that there is a negative relationship between the two variables, "wt" and "qsec". The "wt" variable stands for the weight of a car in thousands of pounds, and the "qsec" variable stands for the time in seconds that a car takes to run a quarter mile. The line of best fit suggests that as "wt" increases, "qsec" decreases. This means that heavier cars tend to have lower quarter mile times, and vice versa. The plot can help us understand the correlation and the regression between the two variables and see how well the line of best fit fits the data.

The R2 value also appears in the top right corner of the plot. This represents the coefficient of determination, which is a measure of how much variation in the response variable "qsec" can be explained by the predictor variable "wt". It ranges from 0 to 1, with higher values indicating a better fit. In this case, the R2 value is 0.031, which means that 3.1% of the variation in "qsec" can be explained by wt. This indicates a very weak linear relationship between the two variables.

### For the fifth plot:

This plot is a pie chart in SPSS, which is a statistical software. A pie chart is a circular chart that uses "pie slices" to display the relative sizes of data. The plot shows the count of a variable "cyl" in the dataset. The chart is divided into three sections, each representing a different value of "cyl". The green section represents the highest count, the blue section represents the second highest count, and the red section represents the lowest count. This chart helps us understand the distribution of the variable "cyl" and how frequently each value occurs.

From the plot, we can see that the value of "cyl" that has the highest count is 8, with 14 cases. The value of "cyl" that has the second highest count is 4, with 11 cases. The value of "cyl" that has the lowest count is 6, with 7 cases. The table below the pie chart also shows these numbers in percentage forms: 46.7% of the cases have a value of 8 for "cyl", 36.7% of the cases have a value of 4 for "cyl", and 23.3% of the cases have a value of 6 for "cyl". The pie chart helps us easily see that the most common value of "cyl" is 8, as it occupies the largest slice of the chart.

 The plot shows the count of a variable "vs" in the dataset. The chart is divided into two sections, one red and one blue. The blue section represents the count of cases that have a value of 0 for "vs", and the red section represents the count of cases that have a value of 1 for "vs". This chart helps us understand the distribution of the variable "vs" and how frequently each value occurs.

From the plot, we can see that the value of "vs" that has the highest count is 0, with 18 cases. The value of "vs" that has the second highest count is 1, with 14 cases. The table below the pie chart also shows these numbers in percentage forms: 56.3% of the cases have a value of 0 for "vs", and 43.8% of the cases have a value of 1 for "vs". The pie chart helps us easily see that the most common value of "vs" is 0, as it occupies the larger slice of the chart.

The "vs" variable stands for the engine shape, which is either V-shaped or straight. A V-shaped engine has two rows of cylinders that form a V shape, while a straight engine has one row of cylinders that form a straight line. The "vs" variable is a binary variable, meaning that it can only take two values: 0 or 1. A value of 0 means that the engine is V-shaped, and a value of 1 means that the engine is straight. The plot shows that most of the cars in the dataset have a V-shaped engine, and only a few have a straight engine.

For the seventh plot:

The plot shows the count of a variable "am" in the dataset. The chart is divided into two sections, one red and one blue. The blue section represents the count of cases that have a value of 0 for "am", and the red section represents the count of cases that have a value of 1 for "am". This chart helps us understand the distribution of the variable "am" and how frequently each value occurs.

From the plot, we can see that the value of "am" that has the highest count is 0, with 19 cases. The value of "am" that has the second highest count is 1, with 13 cases. The table below the pie chart also shows these numbers in percentage forms: 59.4% of the cases have a value of 0 for "am", and 40.6% of the cases have a value of 1 for "am". The pie chart helps us easily see that the majority of the cases have a value of 0 for "am", as it occupies more than half of the chart.

The "am" variable stands for the type of transmission, which is either automatic or manual. A value of 0 means that the transmission is automatic, and a value of 1 means that the transmission is manual. The plot shows that most of the cars in the dataset have an automatic transmission, and only a few have a manual transmission.

For the eighth plot:

The plot shows the count of a variable "gear" in the dataset. The chart is divided into three sections, each representing a different value of "gear". The blue section represents the count of cases that have a value of 3 for "gear", the red section represents the count of cases that have a value of 4 for "gear", and the green section represents the count of cases that have a value of 5 for "gear". This chart helps us understand the distribution of the variable "gear" and how frequently each value occurs.

From the plot, we can see that the value of "gear" that has the highest count is 3, with 15 cases. The value of "gear" that has the second highest count is 4, with 12 cases. The value of "gear" that has the lowest count is 5, with 5 cases. The table below the pie chart also shows these numbers in percentage forms: 50% of the cases have a value of 3 for "gear", 40% of the cases have a value of 4 for "gear", and 16.7% of the cases have a value of 5 for "gear". The pie chart helps us easily see that the most common value of "gear" is 3, as it occupies half of the chart.

The "gear" variable stands for the number of forward gears in a car, which is a measure of its speed and performance. The "gear" variable is a discrete variable, meaning that it can only take a finite number of values. The plot shows that most of the cars in the dataset have either 6 or 4 forward gears, and only a few have 8 forward gears.

For the ninth plot:

The plot shows the count of a variable "carb" in the dataset. The chart is divided into six sections, each representing a different value of "carb". The different colors represent different categories of "carb".

From the plot, we can see that the value of "carb" that has the highest count is 2, with 10 cases. The value of "carb" that has the second highest count is 4, with 10 cases. The value of "carb" that has the third highest count is 1, with 7 cases. The value of "carb" that has the fourth highest count is 3, with 3 cases. The value of "carb" that has the fifth highest count is 6, with 1 case. The value of "carb" that has the lowest count is 8, with 1 case. The table below the pie chart also shows these numbers in percentage forms: 33.3% of the cases have a value of 2 for "carb", 33.3% of the cases have a value of 4 for "carb", 23.3% of the cases have a value of 1 for "carb", 10% of the cases have a value of 3 for "carb", 3.3% of the cases have a value of 6 for "carb", and 3.3% of the cases have a value of 8 for "carb". The pie chart helps us easily see that the most common values of "carb" are 2 and 4, as they occupy the largest slices of the chart.

The "carb" variable stands for the number of carburetors in a car, which is a device that mixes air and fuel for internal combustion engines. The "carb" variable is a discrete variable, meaning that it can only take a finite number of values. The plot shows that most of the cars in the dataset have either 2 or 4 carburetors, and only a few have 1, 3, 6, or 8 carburetors.