

**BE540 - Introduction to Biostatistics
Computer Illustration**

**Topic 1 – Summarizing Data
Software: MINITAB**

A Visit to Yellowstone National Park, USA

Source:

Chatterjee, S; Handcock MS and Simonoff JS *A Casebook for a First Course in Statistics and Data Analysis*. New York, John Wiley, 1995.

Setting:

Upon completion of BE540, you decide to take a vacation to the United States. Of particular interest is seeing an eruption of the famous "Old Faithful" geyser at Yellowstone National Park. Unfortunately, your time is limited and you do not wish to miss seeing an eruption.

This worked example illustrates descriptive analysis of a data set of 222 interval times between eruptions of the Old Faithful Geyser, measured during August 1978 and 1979.

Data File:

GEYSER1.DAT - This is a data set in ASCII format.

Description of Data:

There are three variables, in the following order:

INDEX - An index of the date of the eruption. We will not be using this variable.

DURATION - The duration of the eruption in minutes.

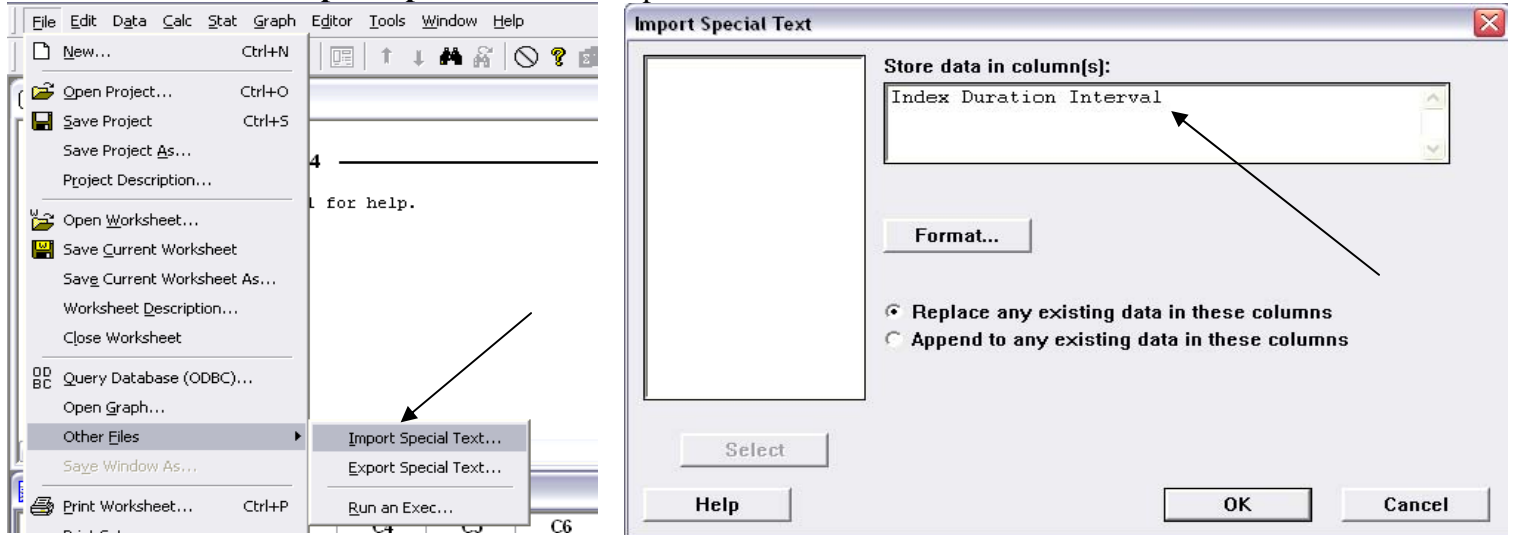
INTERVAL - The length of the interval between the current eruption and the next eruption.

Objective:

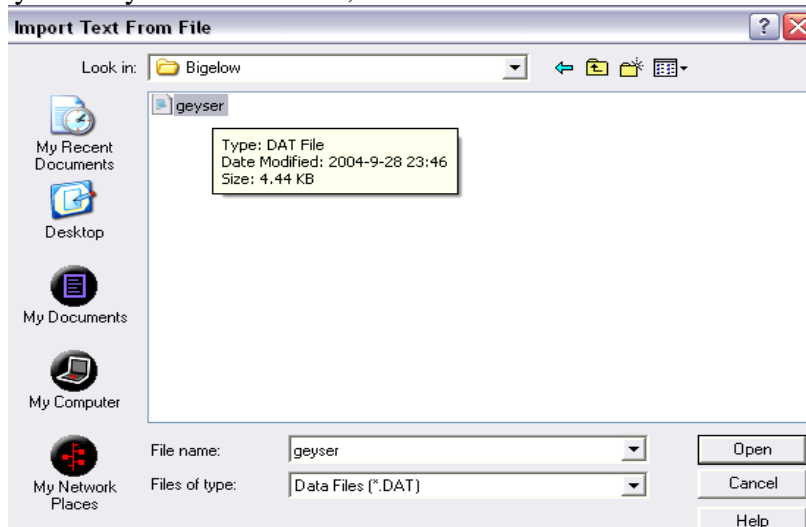
Describe the pattern of eruptions and predict the interval of time to the next eruption.

1. Read in the ASCII format data ‘GEYSER1.DAT’:

1. **File > Other Files > Import Special Text ...**
2. Input variable names: **Index, Duration, Interval** in the “Import Special Text” frame. Click **OK**.



3. Choose correct directory where you save the data, select the .DAT file “GEYSER”. Click **OPEN**.

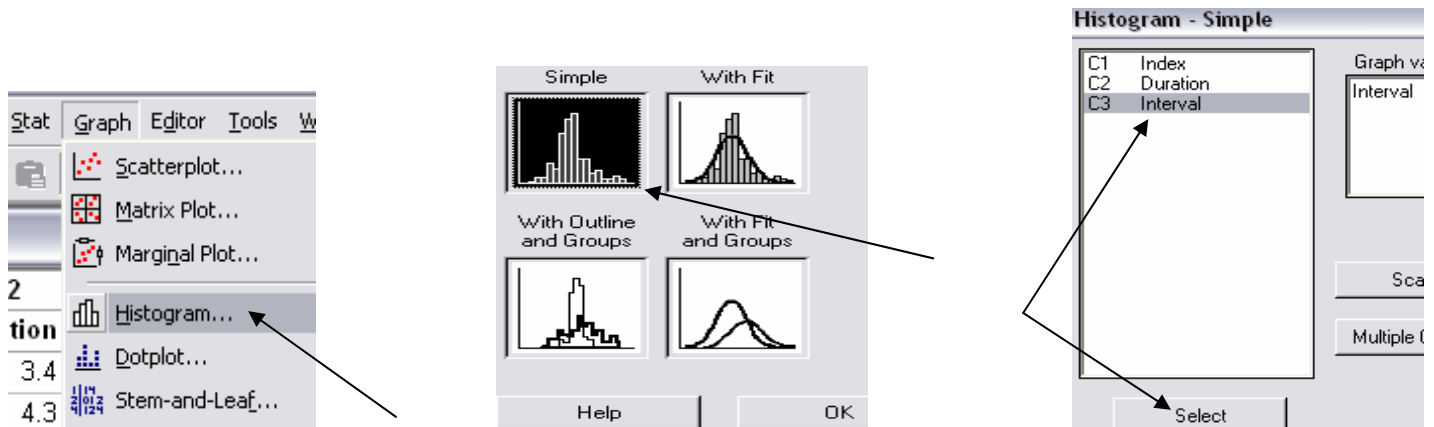


Data set (222 observations) is imported to MINITAB. You should see the following worksheet:

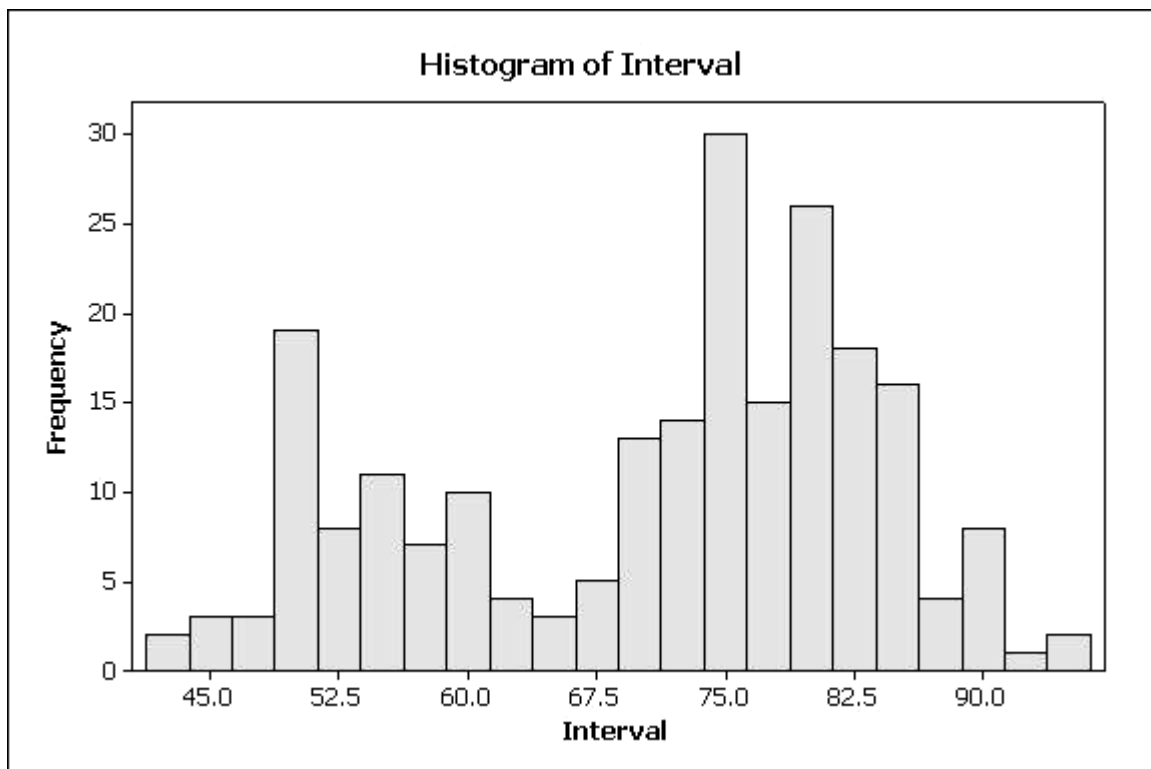
	C1	C2	C3	C4	C5	C6	C7	C8
	Index	Duration	Interval					
1	1	4.4	78					
2	1	3.9	74					
3	1	4.0	68					
4	1	4.0	76					
5	1	3.5	80					

2. Obtain a Histogram of Interval Times.

Menu > **GRAPH** > **HISTOGRAM...** > Choose histogram type, **OK** > Click the interested variable, **SELECT**, **OK**

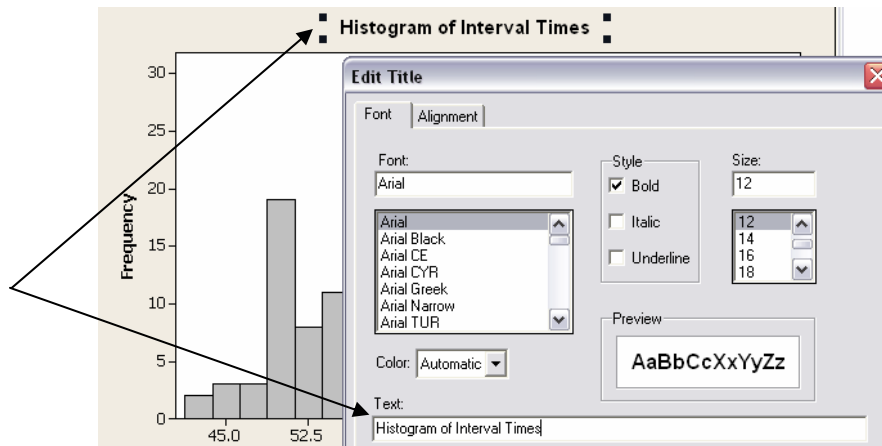


You should see:

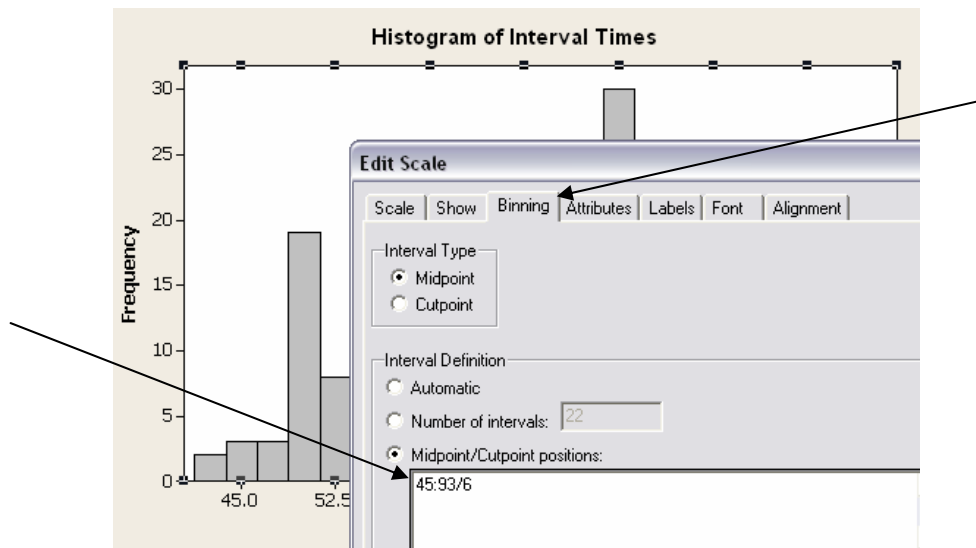


This histogram is generated with MINITAB default settings. If you like, you can revise the format by simply double clicking on the part you want to change. The following are some examples of how to change graph options.

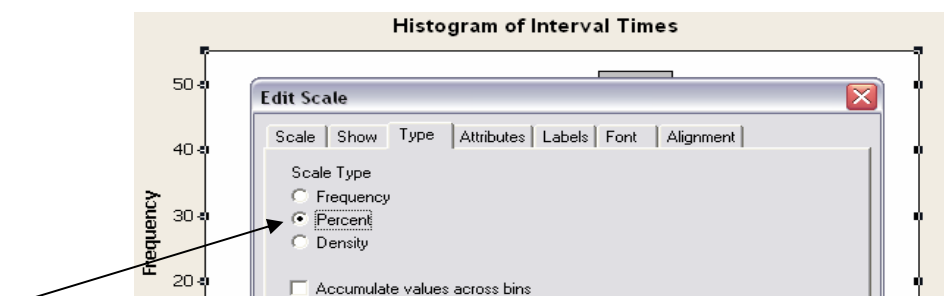
Double click title “Histogram of Interval”, change plot title:



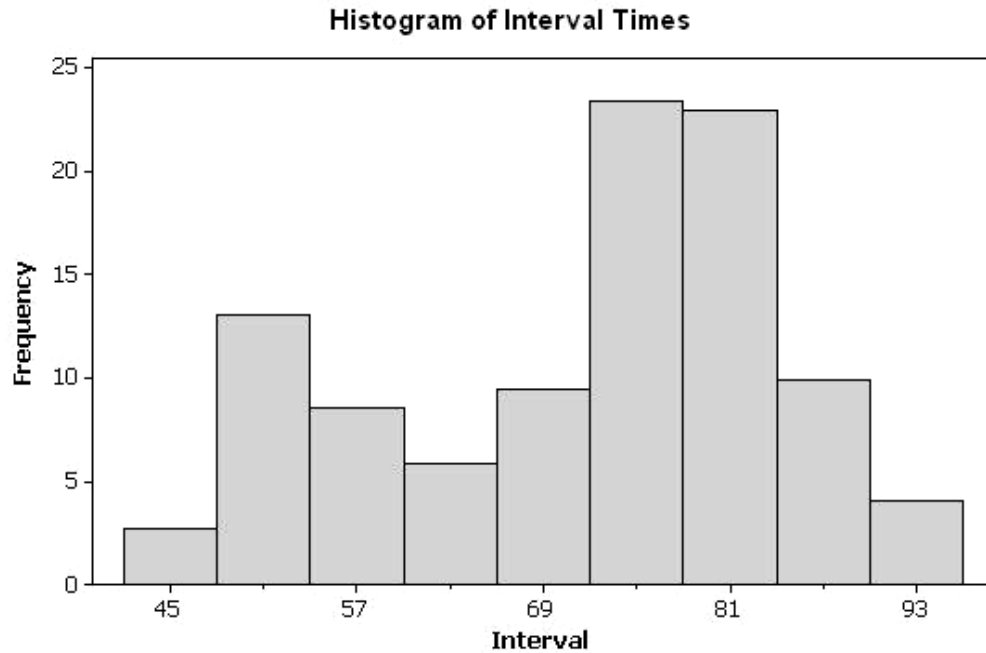
Double click on X-axis label, you can change X scale such as interval number, midpoints position, etc.



Change Y-axis setting by the same way. You can change the frequency plot to percentage histogram: double click Y-axis label, click “TYPE” in Edit Scale frame and check “Percent”. So Y-axis label represents percentage of all 222 observations.



Finally, the revised histogram looks like:



Remarks

The interval times are in the range of 40 to 100 minutes, approximately.

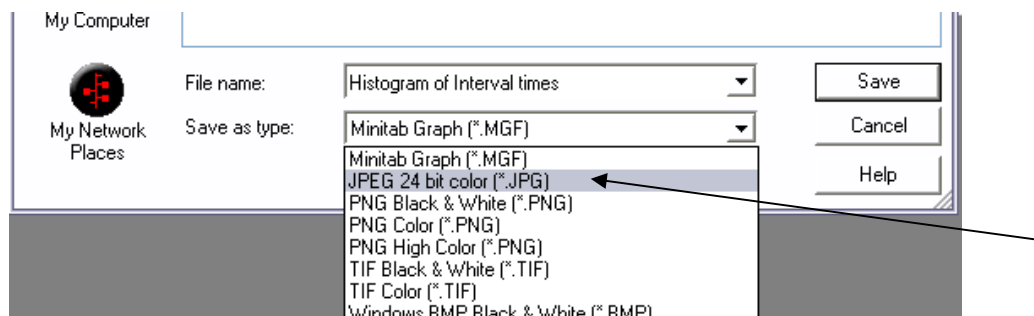
There appears to be two groupings of interval times.

They are centered at 55 and 80 minutes, approximately.

Interestingly, there is a gap in the middle.

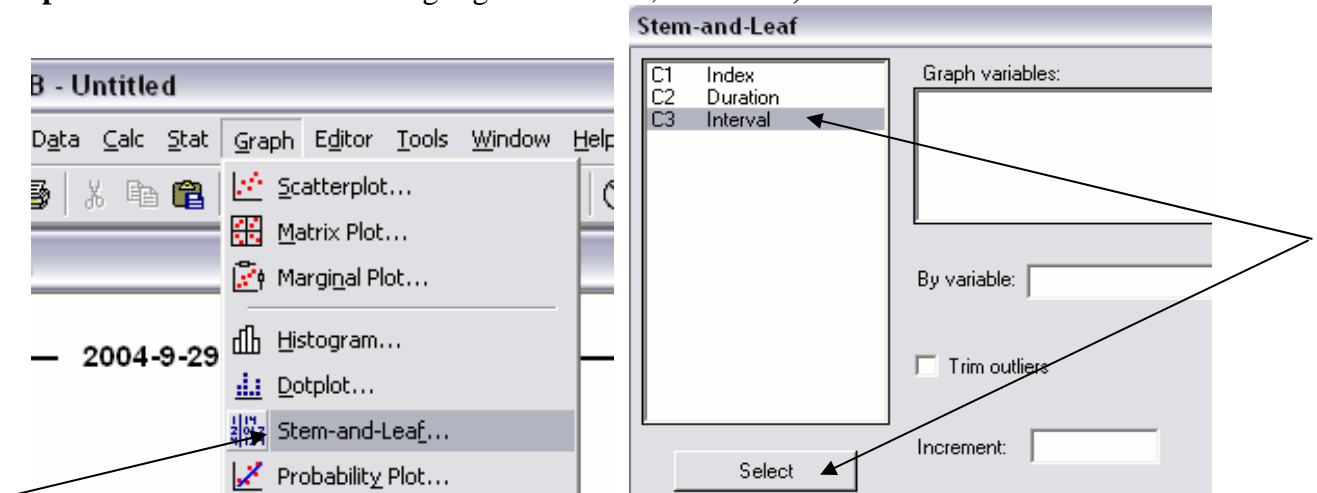
3. Save this histogram as a picture that you can print directly or that you can insert into a document such as this one.

Go to **Menu > FILE > SAVE GRAPH AS ... > Type file name and select picture format (.jpg) > SAVE**



4. Instead of a histogram, we might have constructed a stem-and-leaf diagram.

Menu > Graph > Stem-and-Leaf ... > Highlight "Interval", SELECT, OK.



You should see the following in the *SESSION* window:

Stem-and-Leaf Display: Interval

```
Stem-and-leaf of Interval  N = 222
Leaf Unit = 1.0
 3      4  234
 11     4  55788999
 39     5  0011111111111111222333334444
 54     5  555566677778889
 67     6  0000111112223
 78     6  66677788999
107     7  00000111112222233333333344444
(44)   7  5555555555555566666666667777778888889999
 71     8  0000000000001111111122222222333333334444444444
 22     8  5666666788899
  9     9  00011134
  1     9  5
```

Remarks.

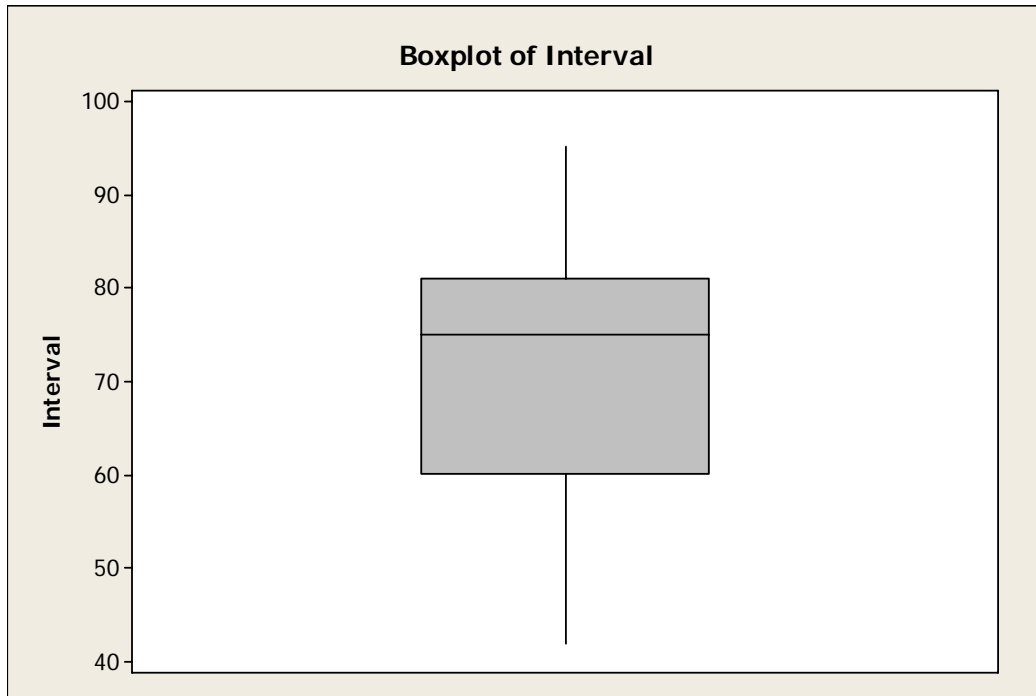
You can see that a stem and leaf diagram is very similar to a histogram. However, we can also see that the minimum and maximum interval times are 42 and 95 minutes, respectively, and that the median time is 75 minutes.

The column of numbers to the left of the stem and leaf diagram is
 from the bottom - a cumulative frequency from the bottom and up.
 from the top - a cumulative frequency from the top and down.

5. In this example, a Box and Whisker plot is not very informative. Let's see why.

Menu > Graph > Boxplot ... > 1-Y Simple , OK > Choose variable "Interval", SELECT, OK.

You should see



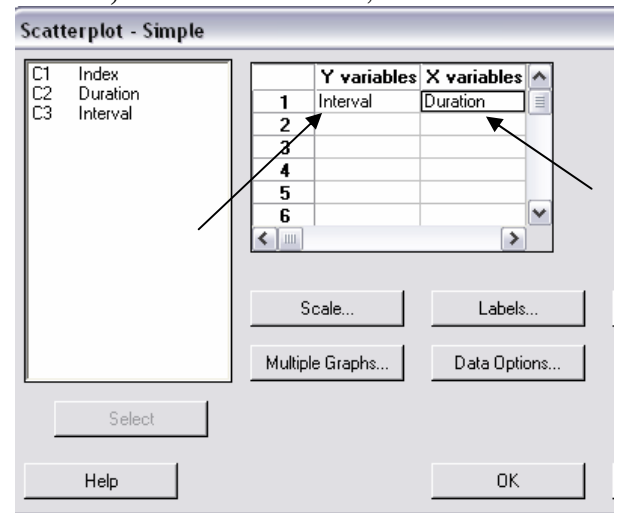
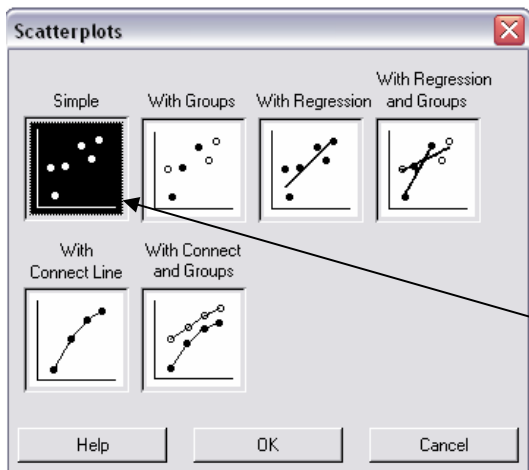
Remarks:

Both the histogram and stem and leaf summaries suggested that there are two groups of interval times. This cannot be seen in a Box and Whisker plot.

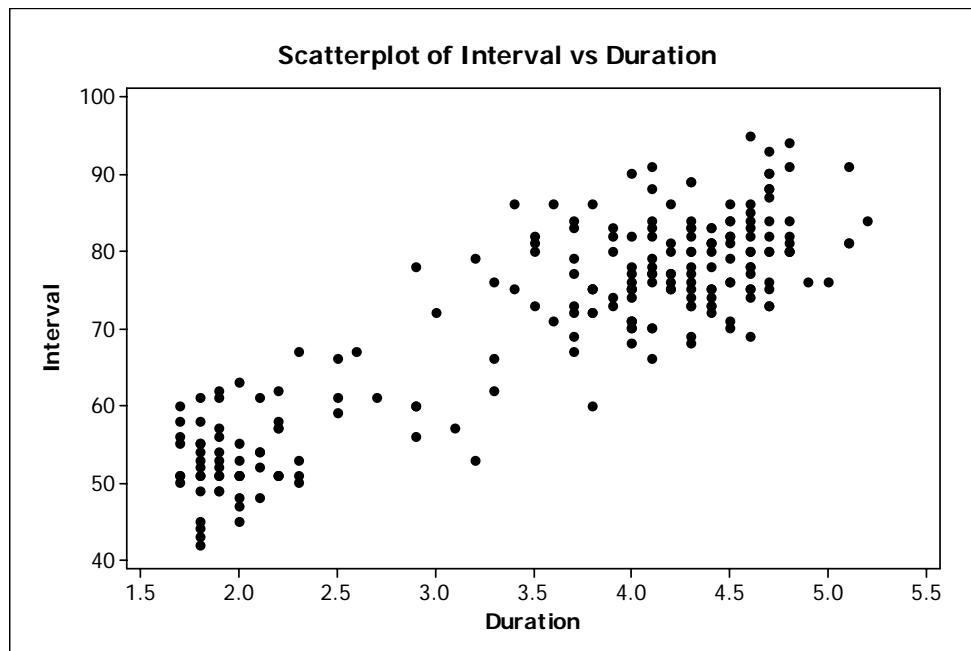
Box and Whisker plots are excellent for summarizing the distribution of ONE population. They are not informative when the sample being summarized actually represents MORE THAN ONE population.

6. We have information on duration of eruption also. One possibility is that the duration of the current eruption is a predictor of the interval time to the next eruption. To investigate this possibility, construct a scatter plot of interval time versus duration. Plot the predictor DURATION on the horizontal axis (X) and the outcome INTERVAL time to the next eruption on the vertical axis (Y).

Menu > Graph > Scatterplot ... > Simple, OK > Select INTERVAL in Y, DURATION in X, OK.



You should see

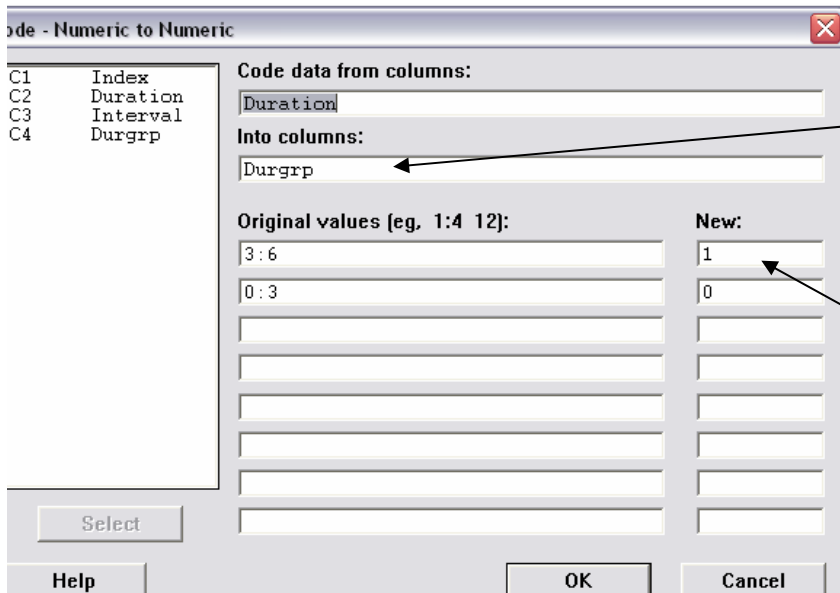


Remarks

The scatter plot confirms a suspected positive association. Longer duration times appear to predict longer intervals to the next eruption. Interestingly, the scatter plot still suggests that there are two distinct subgroups, distinguished by durations of less than versus greater than three minutes.

7. Create a grouped measure of duration and construct separate box and whisker plots of interval times for the interval times that follow eruptions less than 3 minutes in duration and the interval times that follow eruptions greater than 3 minutes in duration.

Menu > Data > Code > Numeric to Numeric ... > Code data: Duration, Into Col: Durgrp, set conditions, OK.

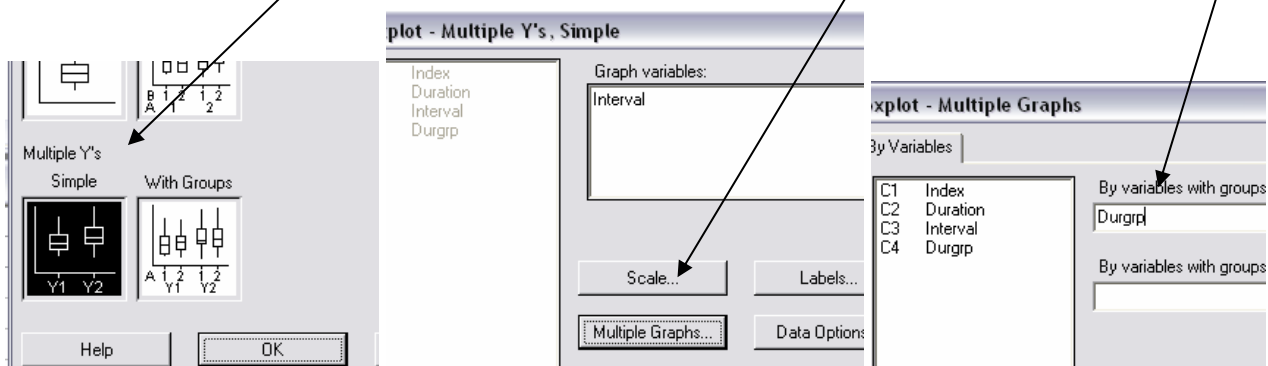


This function is to create new variables from existing variable. "Duration" is the existing, and "Durgrp" is the new variable.

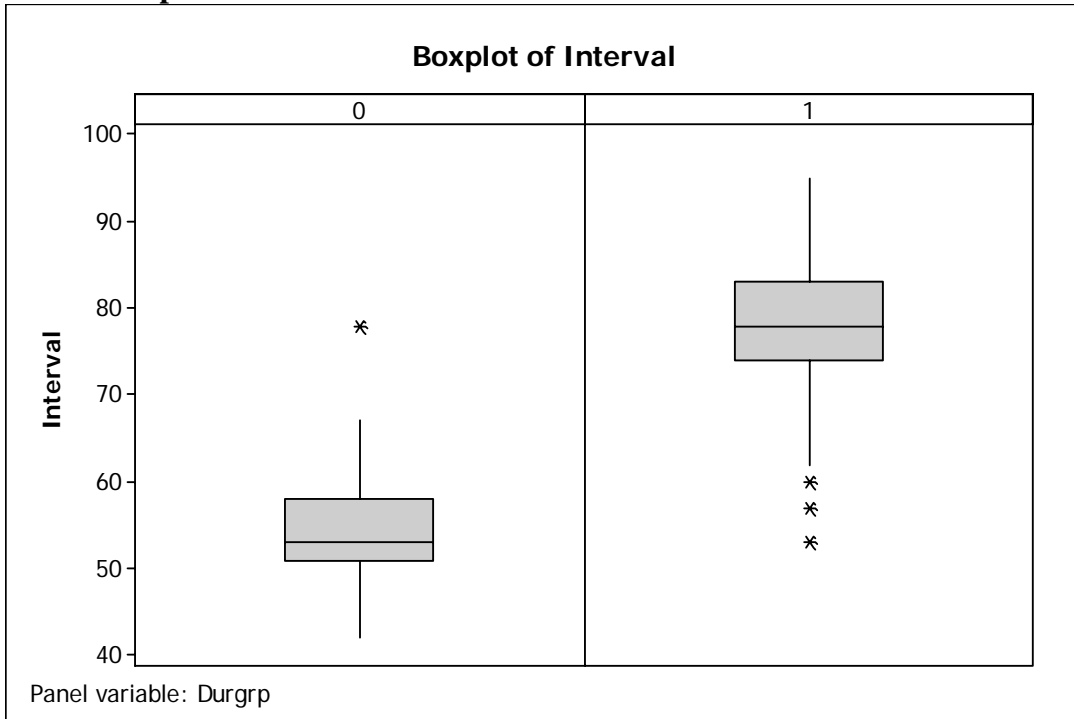
If $Duration > 3$ minutes, Duration group is 1, otherwise, Duration group is 0.

Note: You have just created what is called an indicator variable to indicate a duration time that is greater than 3 minutes. It is equal to 0 for all durations less than 3 minutes and is equal to 1 for all durations greater than 3 minutes. Indicator variables are also called dummy variables or design variables.

Menu>Graph > BoxPlot...>Multiple Y Simple, OK>Select Interval, Click Multiple Graph...>Input "Durgrp", OK.

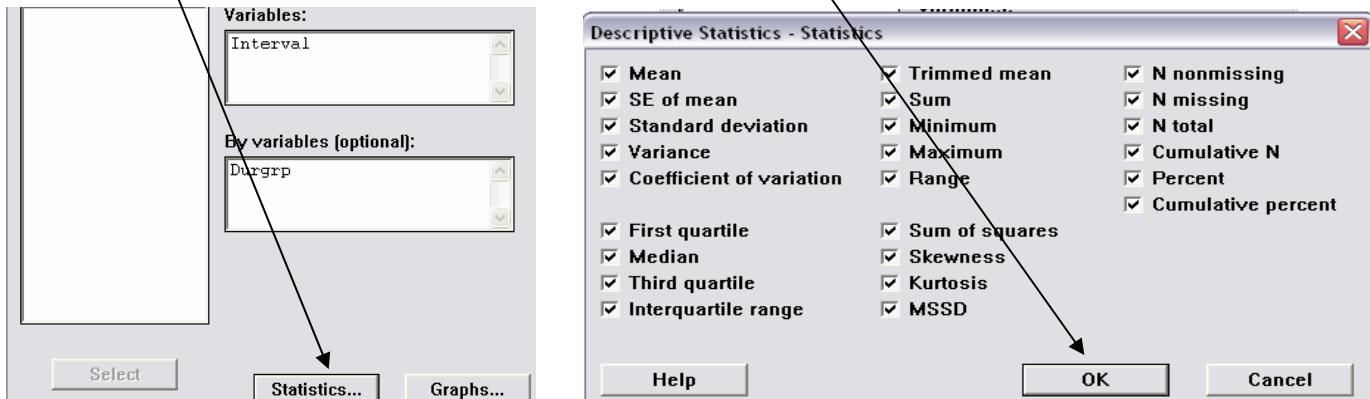


Separate box and whisker plots:



10. Finally, let's look at some numerical summaries, classified by the two groups.

Menu> Stat> Basic Statistics > Display Descriptive Statistics> Input INTERVAL, by DURGRP, STATISTICS. > Choose summary statistics in the list, OK, OK.



Descriptive Statistics: Interval									
		Total							
Variable	Durgrp	Count	N	N*	CumN	Percent	CumPct	Mean	SE Mean
Interval	0	67	67	0	67	30.1802	30.180	54.463	0.770
	1	155	155	0	222	69.8198	100.000	78.161	0.554
		Descriptive Statistics							
Variable	Durgrp	TrMean	StDev	Variance	CoefVar	Sum	Sum of Squares		
Interval	0	54.230	6.299	39.677	11.57	3649.000	201353.000		
	1	78.281	6.891	47.487	8.82	12115.000	954237.000		
		Descriptive Statistics							
Variable	Durgrp	Minimum	Q1	Median	Q3	Maximum	Range	IQR	
Interval	0	42.000	51.000	53.000	58.000	78.000	36.000	7.000	
	1	53.000	74.000	78.000	83.000	95.000	42.000	9.000	
		Descriptive Statistics							
Variable	Durgrp	Skewness	Kurtosis						
Interval	0	0.85	1.92						
	1	-0.34	1.04						

So, what should you do? If you arrive to Old Faithful just after an eruption of less than 3 minutes, with 95% confidence, your waiting time to the next eruption will be between 53 and 56 minutes. Alternatively, if you arrive just after an eruption of greater than 3 minutes, with 95% confidence, your waiting time to the next eruption will be between 77 and 79 minutes.