Distribution and Correlation Analysis of a CSV Imported File Using SAS Studio

Didem Bulut Aykurt

MIS500-1 – Foundations of Data Analytics

Colorado State University-Global Campus
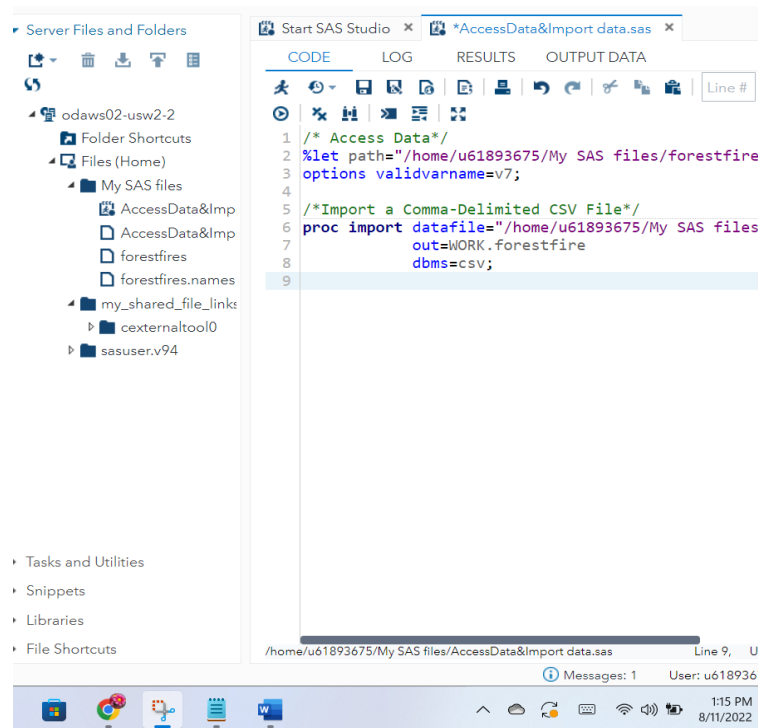
Dr. Steve Chung

August 14, 2022

**Introduction**

The forest fire dataset contains information that helps the most affected forest fire damage and puts human and animal life at risk. This dataset is publicly available from the UCI Machine Learning Repository. The dataset contains 517 observations with 13 attributes. Before applying and executing the model, we should perform a distribution analysis for each variable and check the correlation relation to individual variables. In this assignment, I will use the forest fire dataset's variables 'temp' and 'RH' to analyze the distribution and correlation with temperature in Celsius degrees from 2.2 to 33.30 and RH is relative to the humidity in percentage from a percentage from 15.0 to 100.

**Uploading and Import the forestfire CSV(Comma-delimited) file to SAS Studio**

The first step to uploading CSV data from my computer into SAS Studio is right click on the new icon and click Folder. Then, a new page will pop up, giving a new folder name and save. Then select an existing folder and right click new folder to upload a file from my laptop. Now is forestfire.csv online on SAS Studio; change the file name right, click then rename forestfire. Screenshot 2 has details.

## Screenshot 1;

*SAS code to import forestfire CSV file into Studio on SAS studio*



## Screenshot 2

*That shows the dataset table on OUTPUT DATA on SAS studio.*

**Analyzing the distribution of temp 'temp' variable**

Analyze the distribution of the temp variable with SAS Studio;
click Tasks and Utilities from in the navigation menu, then
Distribution Analyst under the Statistic as it opens a new code
page. The upper ribbon has the list of the Data and Options menu.
Add forest fire dataset and analysis variable is 'temp' following
select 'Histogram' and 'Normal quantile-quantile plot' adds inset
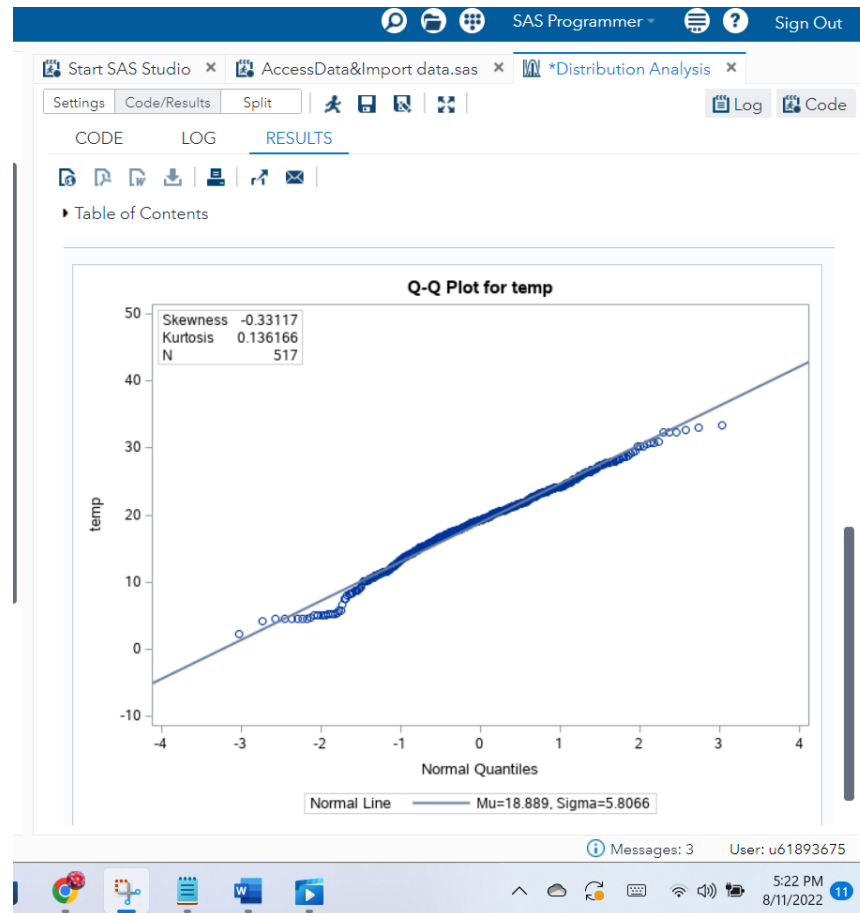statistics—the result on screenshots 3 and 4.

**Screenshot 3**

*The resulting distribution analysis is a histogram of the 'temp' variable
on SAS studio.*

**Screenshot 4**

*The resulting distribution analysis is the Norma quantile-quantile plot of the 'temp' variable on SAS studio.*



Skewness tells dataset has an asymmetric distribution or not that measures three different distributions. One zero skew means the distribution is symmetrical. Others negative skew when the number is negative and positive skew when the number is positive. The skewness of temp shows a negative number of -0.33, which tells the left skew.

Kurtosis shows the variable's probability or frequency, which also helps to compare which variable has a heavy distribution tail with three kurtosis types. I found so many different ranges people use, and I would like to use several zero for the normal distribution of kurtosis because I check outliers and best explanate for the case outliers with the number of zero for kurtosis. Medium tails are **mesokurtic(kurtosis=0)**, low kurtosis is **platykurtic(kurtosis<0),** and high kurtosis is **leptokurtic(kurtosis>0).** The temp variable has leptokurtic because the temp's kurtosis is 0.13 as positive kurtosis more than zero as a normal distribution kurtosis's number is zero or a close zero. Thus, the temp distribution should have outliers; the skewness result is close to zero, which can be accepted for normal distribution if outliers exceed what we expect.
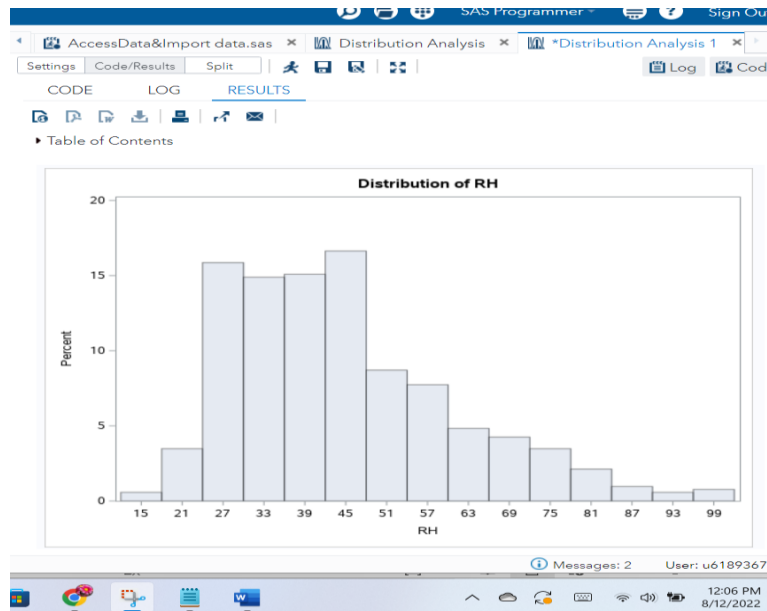
**Analyzing the distribution of temp 'RH' variable**

First, I change the variable from the Data section on the top ribbon Option has a checklist as I use Histogram under the exploring data, a Normal quantile-quantile plot with inset statistic under checking for Normality, and a number of observations, skewness, and kurtosis. Screenshot 5 has a result for histogram distribution analysis, and 6 shows a Q-Q plot for the 'RH' variable.
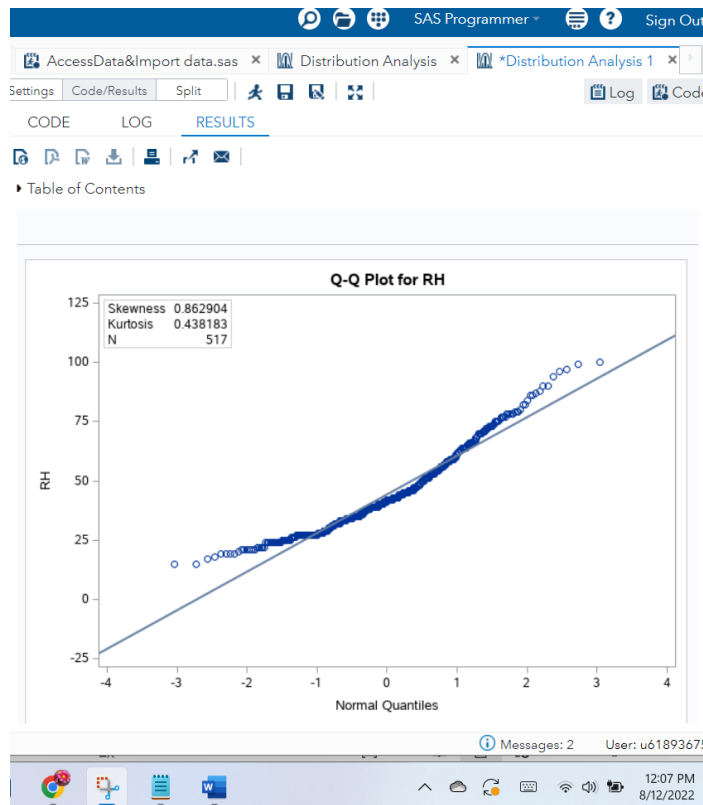
**Screenshot 5**
*The resulting distribution analysis is a histogram of the 'RH' variable on SAS studio on SAS studio.*

## Screenshot 6

*The resulting distribution analysis is the Norma quantile-quantile plot of the 'temp' variable on SAS studio.*

RH's histogram(screenshot5) shows more variables on the right side which means the distribution has a strong positive skew in the dataset. That also indicates that the' RH' distribution is not normal.

Q-Q plot for 'RH's results skewness is 0.86 that also say same what we found histogram strongly positive skewness.

Kurtosis's result is 0.43, as positive kurtosis would have more outliers. Thus, 'RH' has a positive skew and leptokurtic distribution.
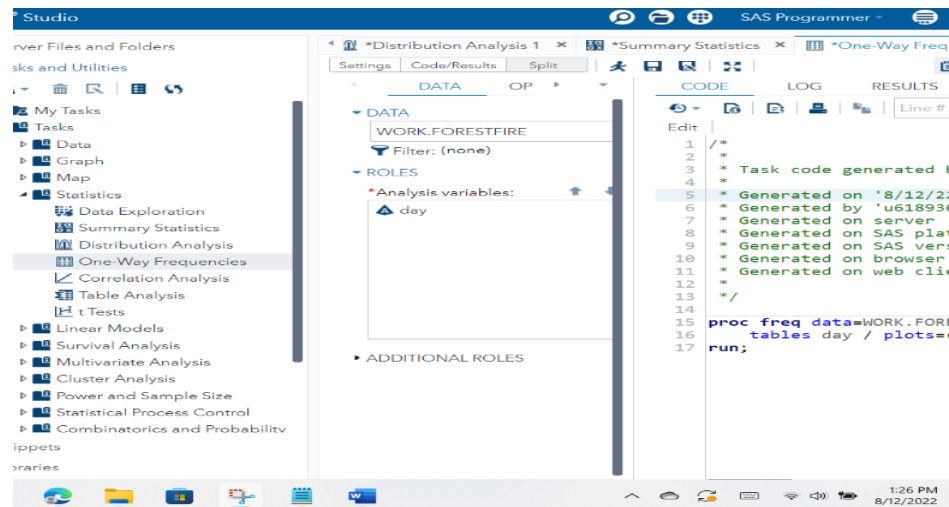
**One-Way Frequencies analysis by 'day.'**

What tools are ideally used for nominal variables to count frequency by each category? I worked on SAS studio's main menu, Tasks and Utilities, then One-way Frequencies. Add variable 'day' to show each day's frequency on the histogram. Screenshot 7 shows the code detail, and screenshot 8 has the result page.
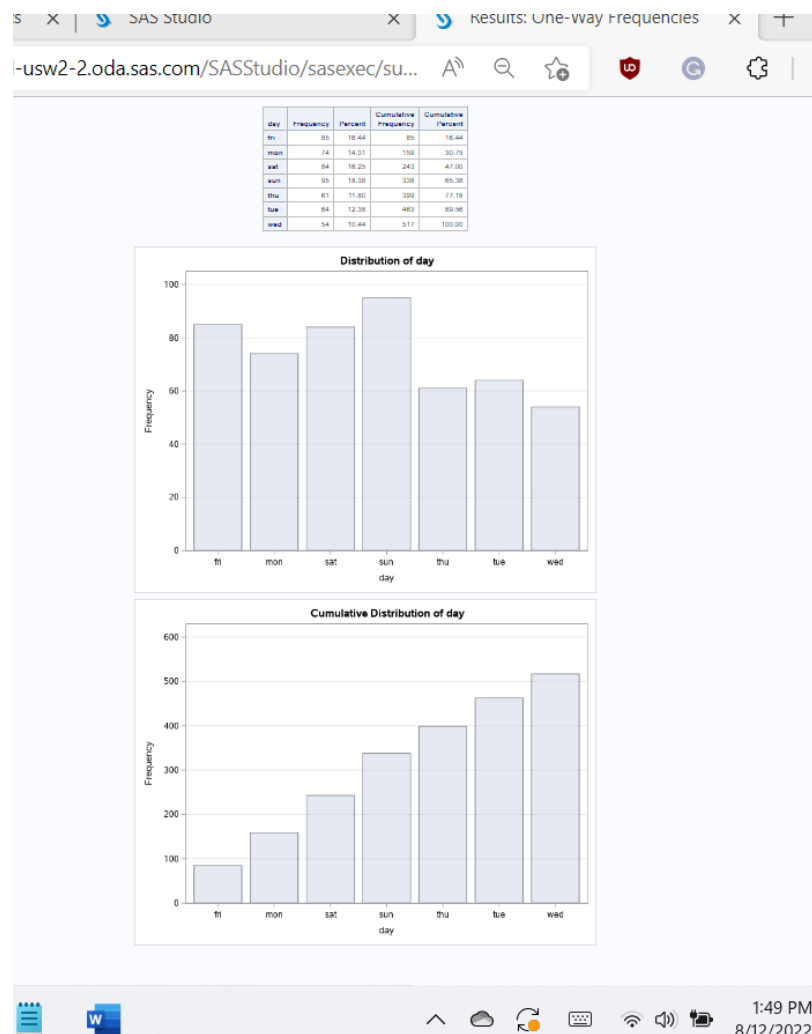
**Screenshot 7**
*The code of the one-way frequencies by 'day' on SAS studio.*

**Screenshot 8**

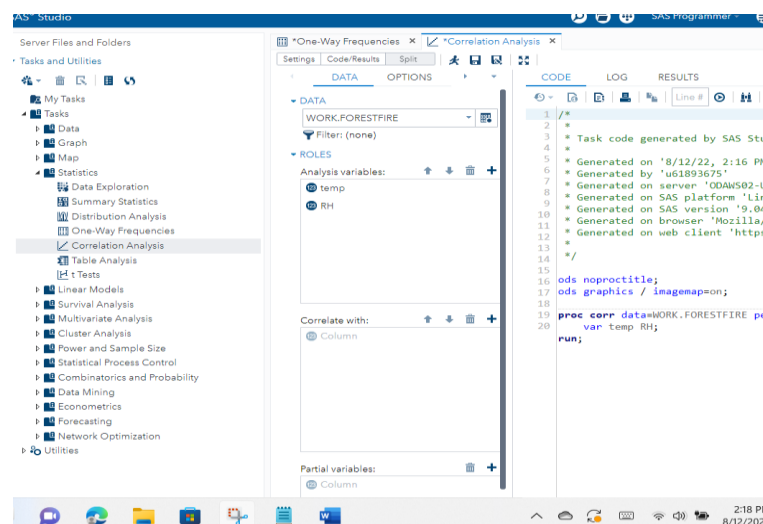*The result of the distribution of 'day' by One-way Frequencies on SAS studio.*

As we can easily compare, which day has more frequency than other days with frequency distribution is an outline of all prominent values. Screenshot 8, the top on the graph, shows the most frequent day on Sunday. The bottom graph is a cumulative frequency distribution that helps if there is a specific value and to see how probability is.

**Correlation Analysis of the Variables of 'temp' and 'RH.'**

Correlation analysis helps us to know how to relation between two variables. I use the Tasks and Utilities main menu, double-click Correlation Analysis, add variables 'temp' and 'RH' on the DATA page, and move to the OPTIONS page select Matrix of scatter plots. All steps are on screenshots 9 and 10. The result page is in screenshot 11.

**Screenshot 9**

*The code of the correlation analysis of 'temp' and 'RH' on SAS Studio.*
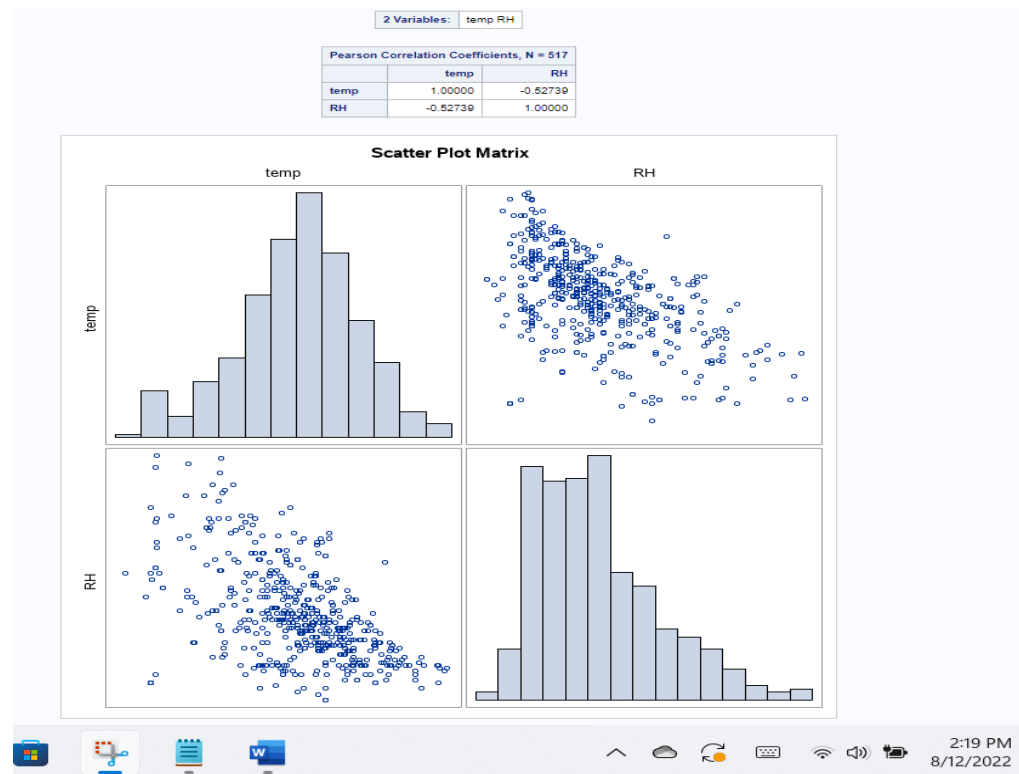
**Screenshot 10**

*The code of correlation analysis Option page for 'temp' and "RH' variables on SAS Studio.*



**Screenshot 11**

*The result of correlation analysis for 'temp' and 'RH' variables on SAS Studio.*

Thus, 'temp' and 'RH' have a negative correlation of -0.52, which means when the temperature increases, the humidity decreases (see screenshot 11).

**Conclusion**

That great to learn how to upload complex files online on SAS Studio and Import the file. I use Distribution Analysis for each variable to see how the data behave. The one-way frequencies excellent tools that I use for nominal variables; in this case, I used the 'day' variable. The correlation analysis with the matrix plot menu shows the 'temp' and 'RH' variable relation.

**References**

UCI Machine Learning Repository: Forest Fires Data Set