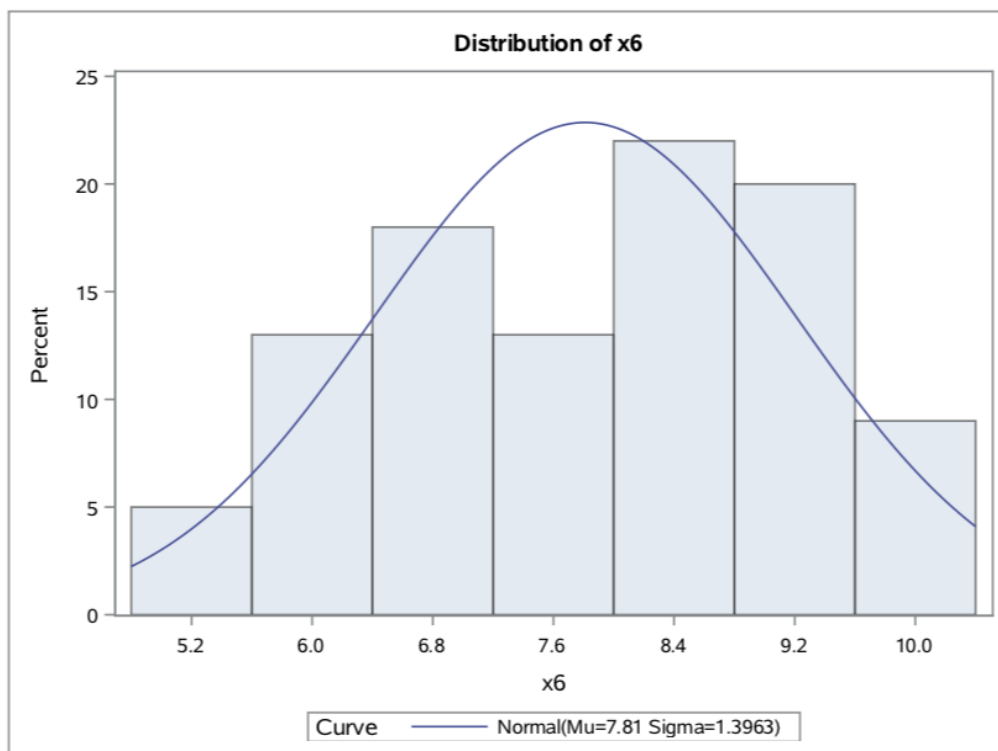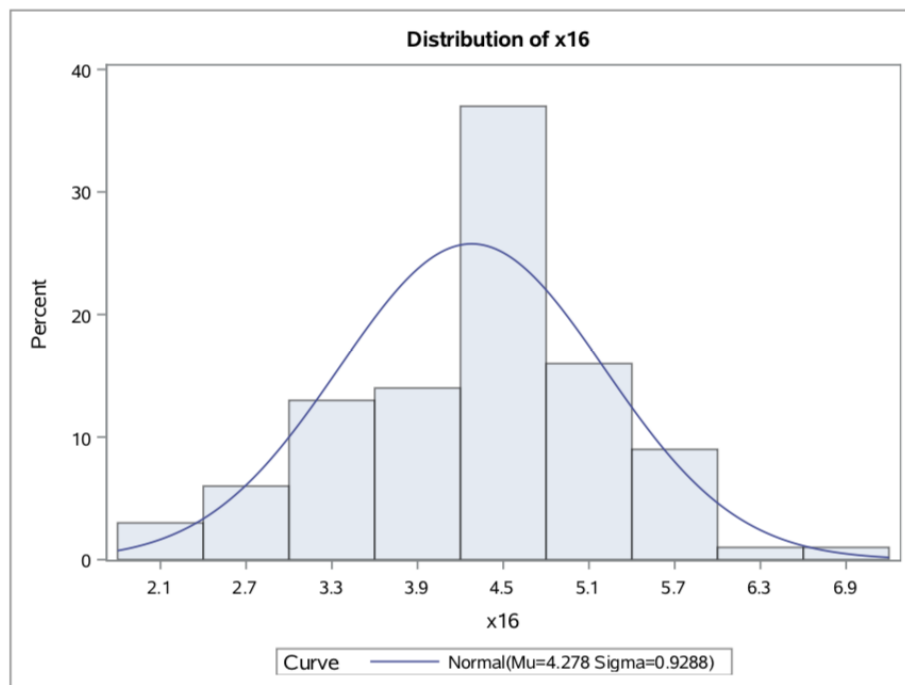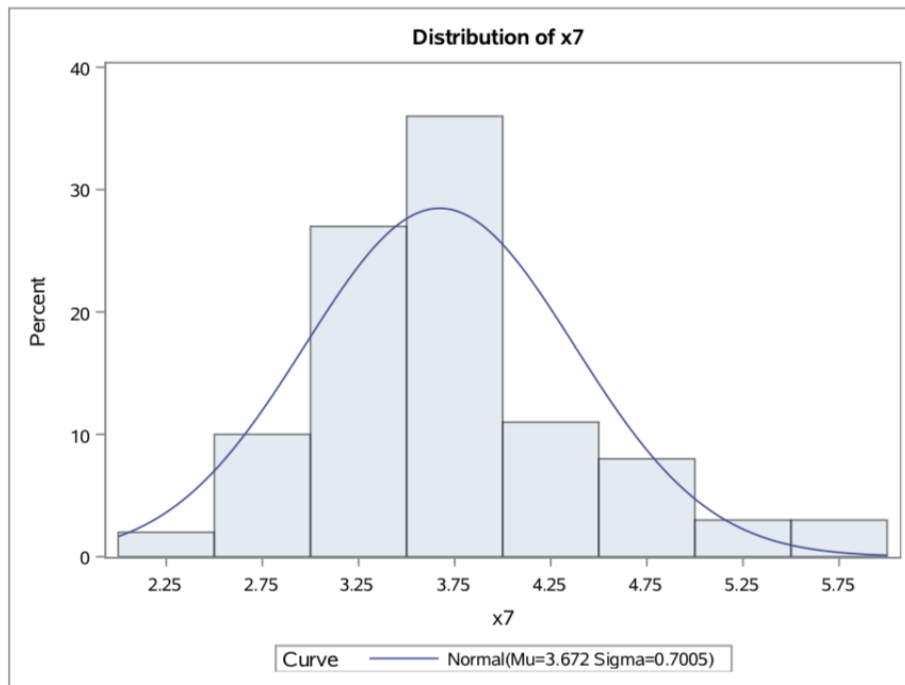1. Missing Data

Missing data is an important topic because data sets often will not be perfect when you consider the real world. Missing data is defined as information that is not available for a particular subject or case when there is other information about that subject/case available. This is a relatively common phenomenon.
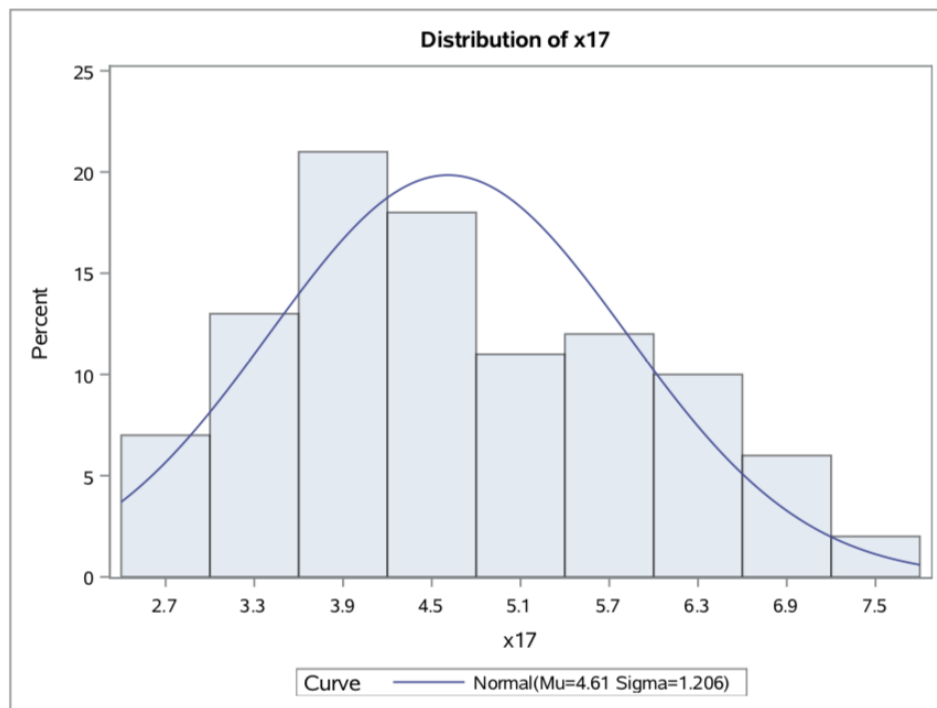
There are several causes of missing data. First, nonresponse can occur in surveys, when the subject answers some of the questions but not all of them. Many times the omitted responses come on personal questions. Next, attrition can occur when a subject begins an experiment but does not complete it. Another example is improper data collection. This is a more generic category. It can occur when the researcher loses data, mislabels data, improperly records data, improperly conducts the experiment, etc. Lastly, governments and private companies frequently publish data on many different topics, but sometimes there is data missing. This can be due to the entities choosing not to publish certain data if it does not look favorable for them, or it can simply be due to them not having the data or having missing data themselves.

After identifying missing data, one should consider any patterns. Is the missing data completely random, random, or not random at all? The answer to this question will help determine how one should treat the missing data. Another task after identifying missing data is quantifying it. If under 10% of the data for an individual observation is missing, this can be ignored. However, if the data is missing in a manner that is not random, this must be reconsidered. If the data is still going to be used, the missing data needs to be handled, and there are several strategies for this. One can only use observations with complete data, certain cases or variables with the missing data can be deleted, or the missing values can be estimated.

2a.



Distribution of x6

Curve —— Normal(Mu=7.81 Sigma=1.3963)

## Distribution of x7



Curve —— Normal(Mu=3.672 Sigma=0.7005)

## Distribution of x16



Curve —— Normal(Mu=4.278 Sigma=0.9288)

Distribution of x17

2b.

| Analysis Variable : x6 x6 | | | | | | |
|---|---|---|---|---|---|---|
| Mean | Std Dev | Minimum | Maximum | N | Skewness | Kurtosis |
| 7.8100000 | 1.3962793 | 5.0000000 | 10.0000000 | 100 | -0.2445019 | -1.1318375 |

| Analysis Variable : x7 x7 | | | | | | |
|---|---|---|---|---|---|---|
| Mean | Std Dev | Minimum | Maximum | N | Skewness | Kurtosis |
| 3.6720000 | 0.7005164 | 2.2000000 | 5.7000000 | 100 | 0.6603903 | 0.7353470 |

| Analysis Variable : x16 x16 | | | | | | |
|---|---|---|---|---|---|---|
| Mean | Std Dev | Minimum | Maximum | N | Skewness | Kurtosis |
| 4.2780000 | 0.9288398 | 2.0000000 | 6.7000000 | 100 | -0.3335404 | 0.2441491 |

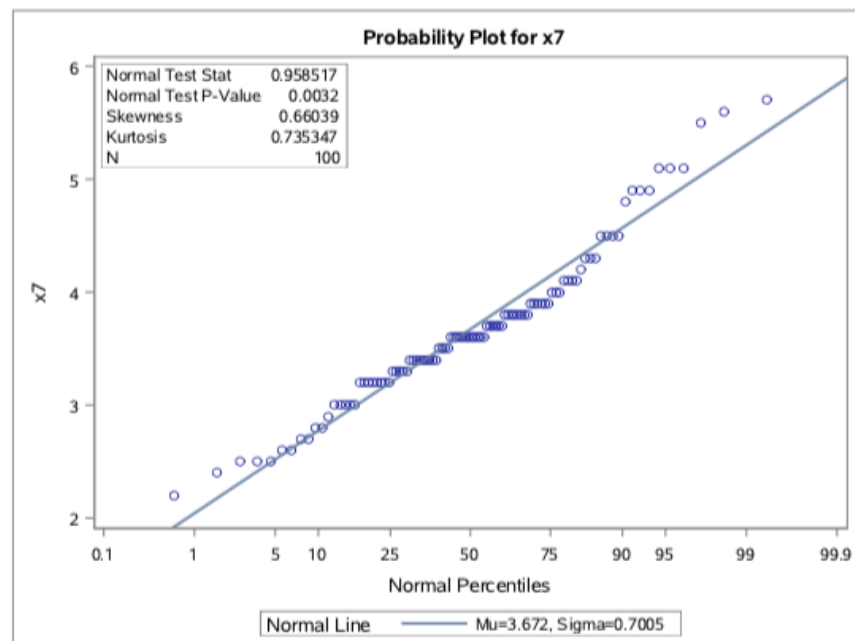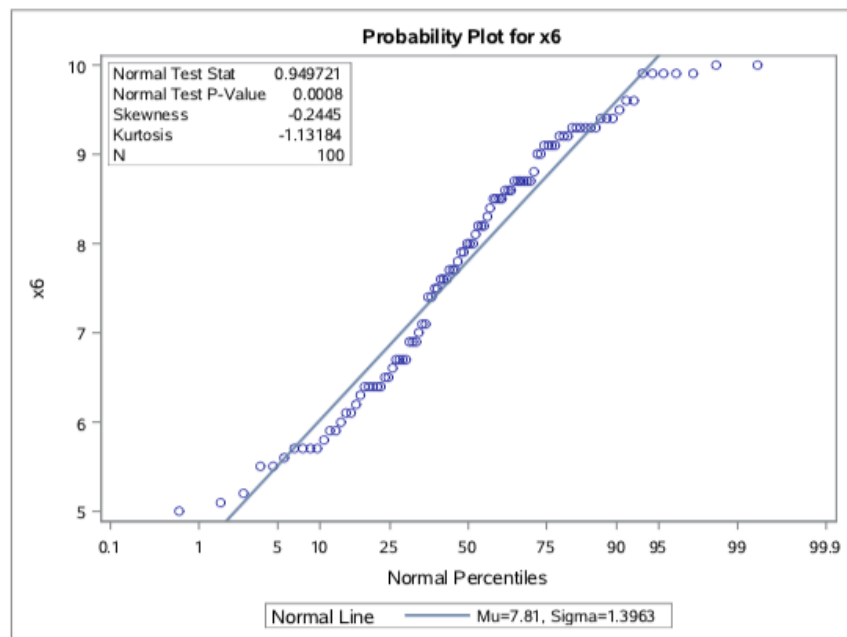| Analysis Variable : x17 x17 | | | | | | |
|---|---|---|---|---|---|---|
| Mean | Std Dev | Minimum | Maximum | N | Skewness | Kurtosis |
| 4.6100000 | 1.2060035 | 2.6000000 | 7.3000000 | 100 | 0.3227665 | -0.8158665 |

Using the formula, the z for skewness for x6 is ~1.0 and the z for kurtosis is ~2.3. The z for kurtosis may suggest a lack of normality depending on the significance level.
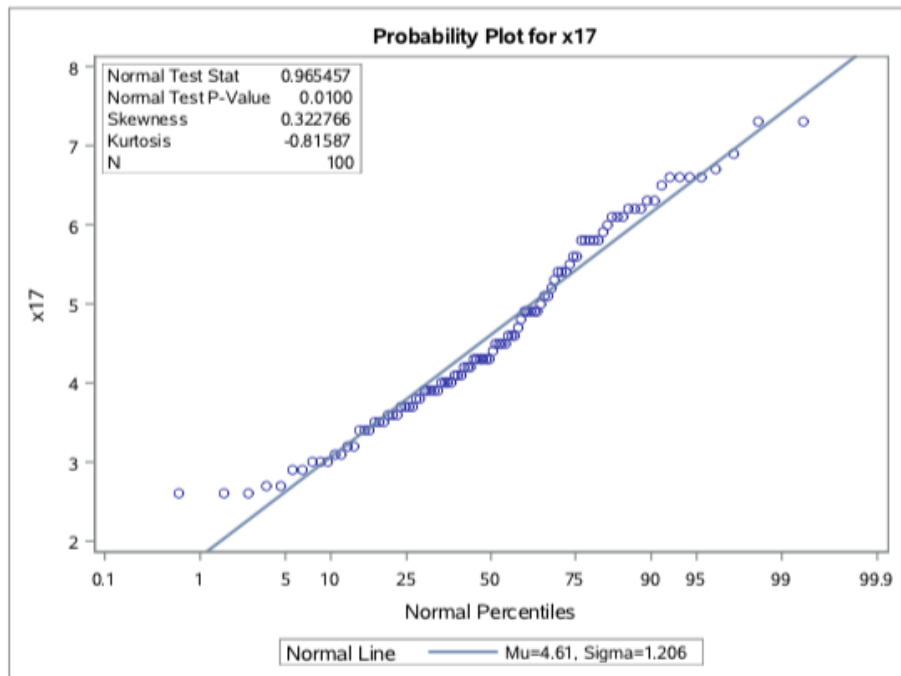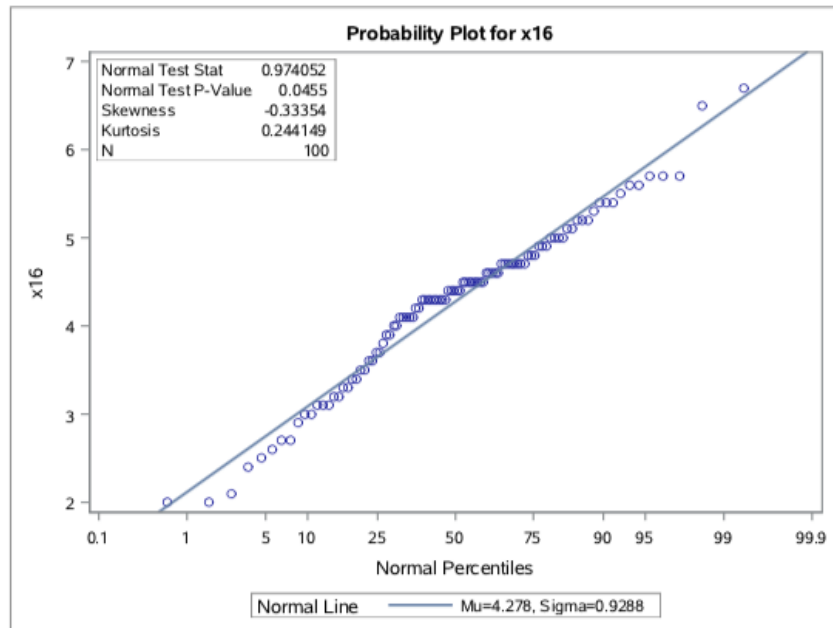
For x7, the z for skewness is ~2.7 and the z for kurtosis is ~1.5. The skewness z definitely suggests a lack of normality.

For x16, the z for skewness is ~1.4 and the z for kurtosis is ~.5. These z values do not suggest a lack of normality.

For x17, the z for skewness is ~1.3 and the z for kurtosis is ~1.7. These z values do not suggest a lack of normality.

2c.



Probability Plot for x6

| Normal Test Stat | 0.949721 |
| Normal Test P-Value | 0.0008 |
| Skewness | -0.2445 |
| Kurtosis | -1.13184 |
| N | 100 |

Normal Percentiles

Normal Line — Mu=7.81, Sigma=1.3963



Probability Plot for x7

| Normal Test Stat | 0.958517 |
| Normal Test P-Value | 0.0032 |
| Skewness | 0.66039 |
| Kurtosis | 0.735347 |
| N | 100 |

Normal Percentiles

Normal Line — Mu=3.672, Sigma=0.7005

## Probability Plot for x16

| Normal Test Stat | 0.974052 |
| Normal Test P-Value | 0.0455 |
| Skewness | -0.33354 |
| Kurtosis | 0.244149 |
| N | 100 |

Normal Line — Mu=4.278, Sigma=0.9288

Normal Percentiles

## Probability Plot for x17

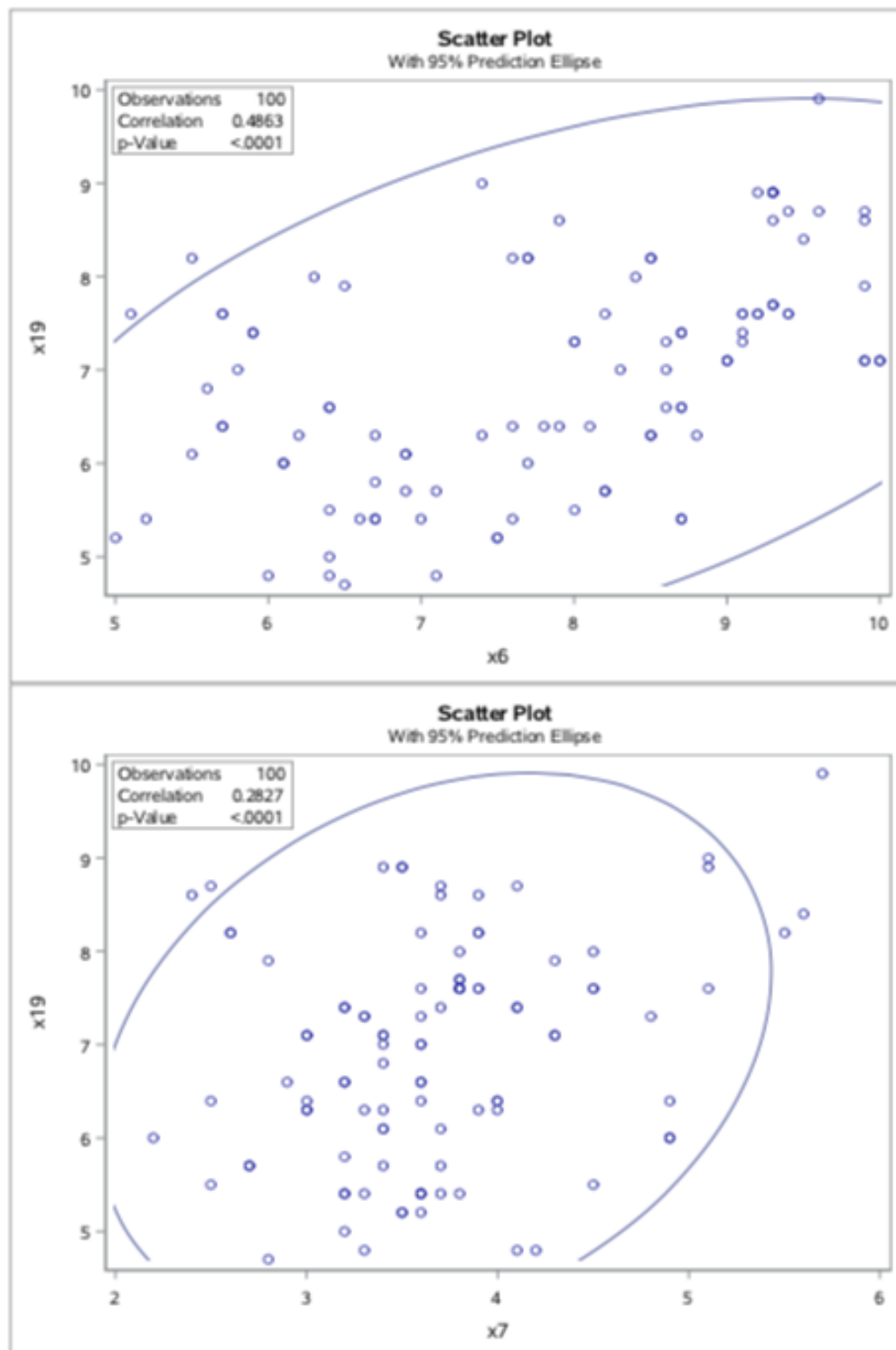| Normal Test Stat | 0.965457 |
| Normal Test P-Value | 0.0100 |
| Skewness | 0.322766 |
| Kurtosis | -0.81587 |
| N | 100 |

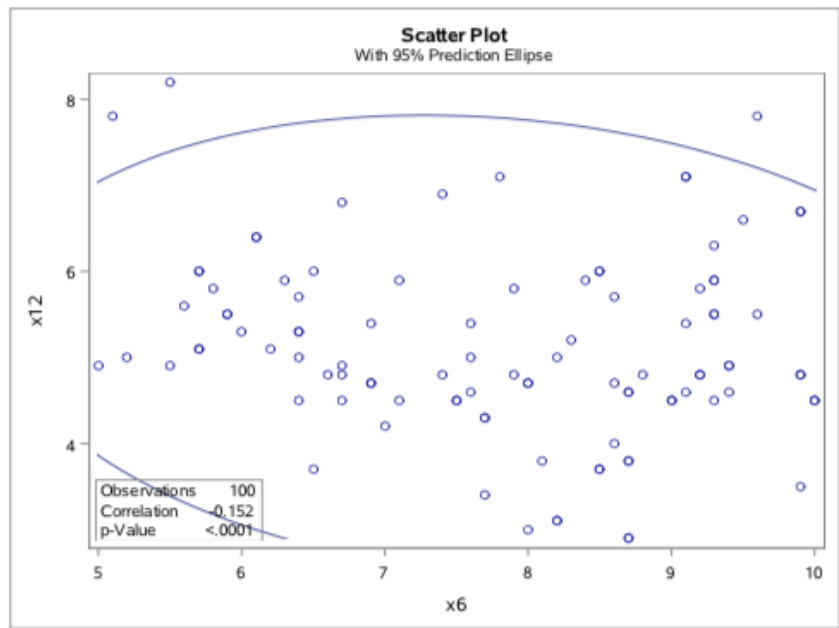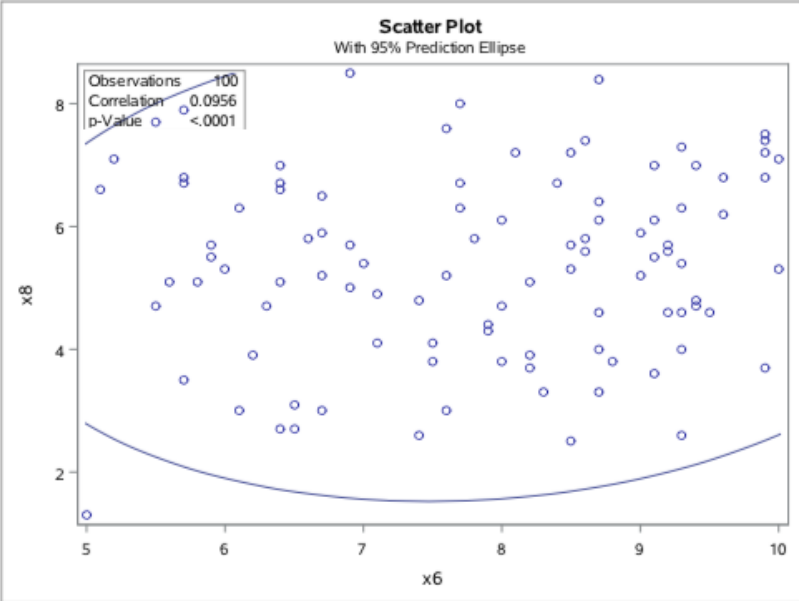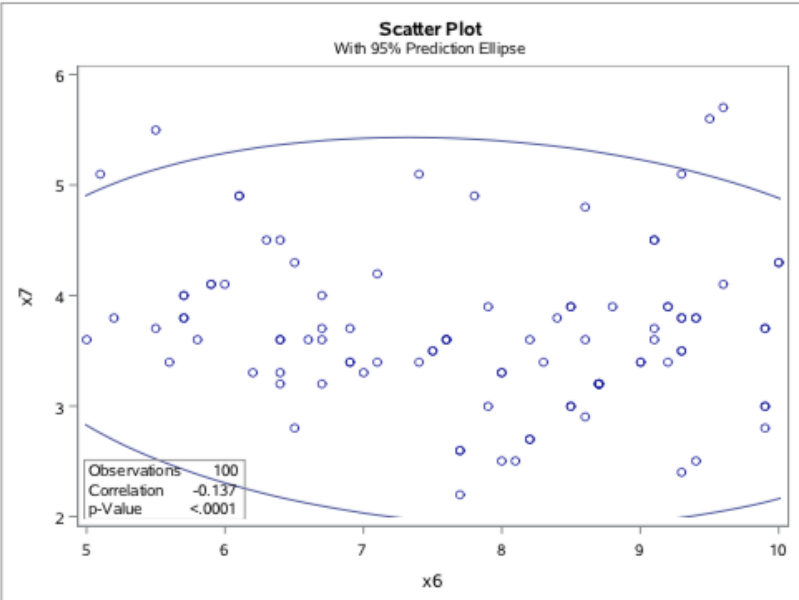Normal Line — Mu=4.61, Sigma=1.206

Normal Percentiles

All the p-values are below 0.05. However, x16 is not below 0.01. So depending on the significance level, this is the only variable that may not be normal.
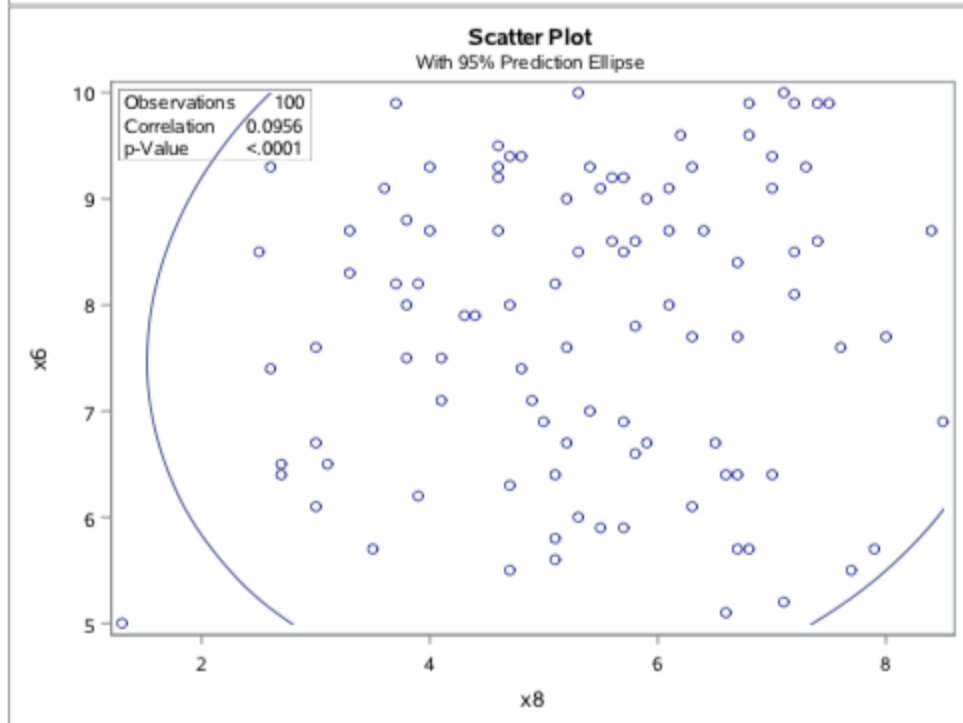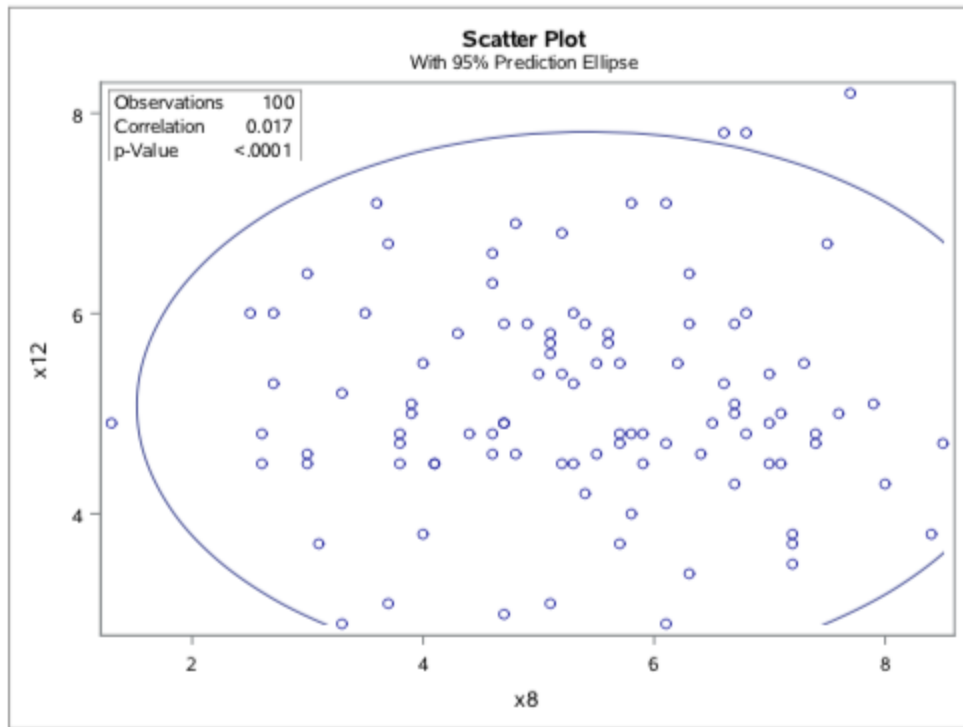
2d.



**Scatter Plot**
With 95% Prediction Ellipse

| Observations | 100 |
| Correlation | 0.4863 |
| p-Value | <.0001 |

**Scatter Plot**
With 95% Prediction Ellipse

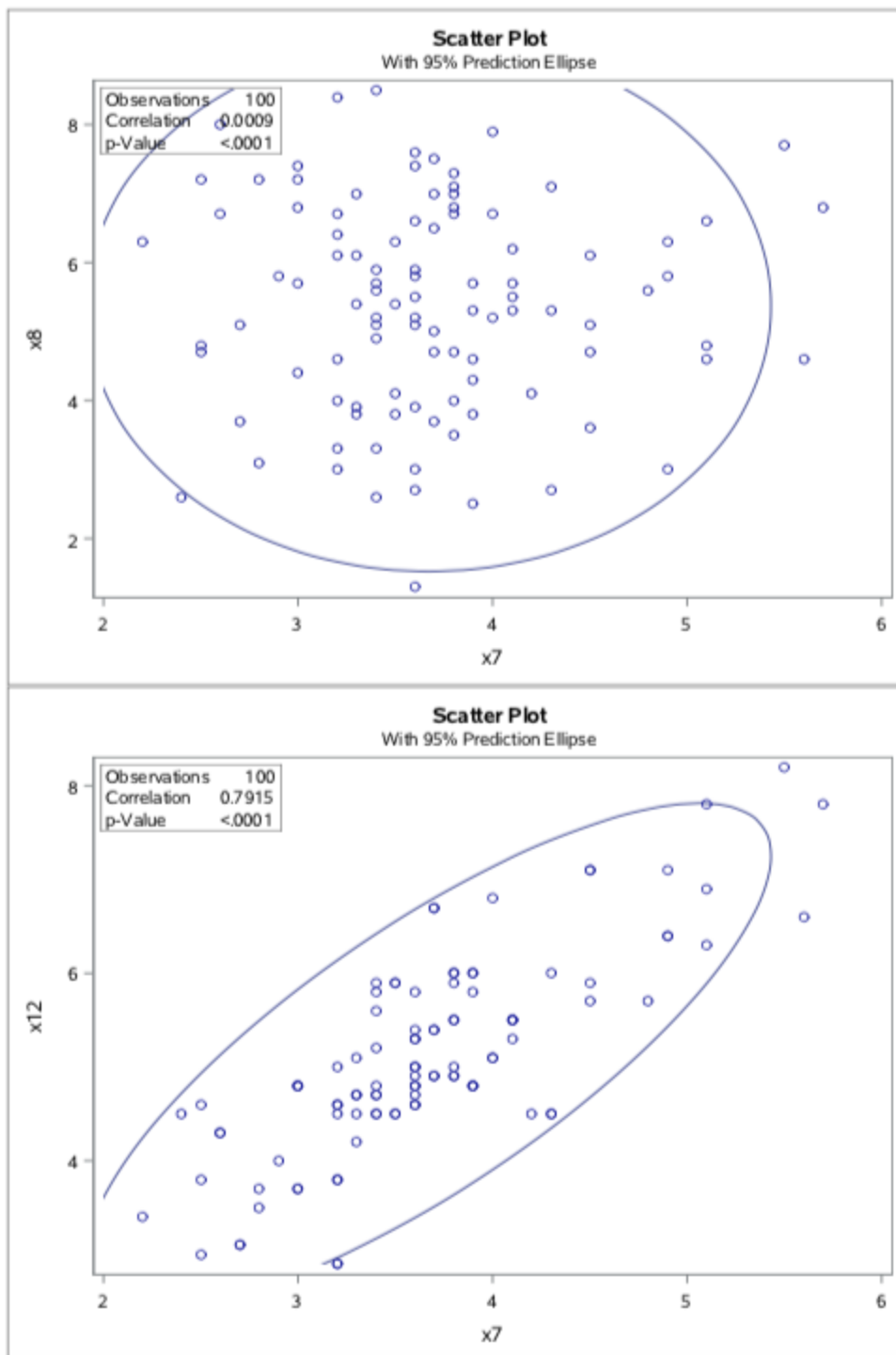| Observations | 100 |
| Correlation | 0.2827 |
| p-Value | <.0001 |

In the top plot, there are two outliers plus one that is right on the border. In the bottom plot, there are five outliers.

2e.



**Scatter Plot**
With 95% Prediction Ellipse

| Observations | 100 |
| Correlation | -0.137 |
| p-Value | <.0001 |

**Scatter Plot**
With 95% Prediction Ellipse

| Observations | 100 |
| Correlation | 0.0956 |
| p-Value | <.0001 |

**Scatter Plot**
With 95% Prediction Ellipse

| Observations | 100 |
| Correlation | 0.152 |
| p-Value | <.0001 |

# Scatter Plot
With 95% Prediction Ellipse

| Observations | 100 |
|---|---|
| Correlation | 0.017 |
| p-Value | <.0001 |

x12

x8

# Scatter Plot
With 95% Prediction Ellipse

| Observations | 100 |
|---|---|
| Correlation | 0.0956 |
| p-Value | <.0001 |

x6

x8

**Scatter Plot**
With 95% Prediction Ellipse

| Observations | 100 |
| Correlation | 0.0009 |
| p-Value | <.0001 |

**Scatter Plot**
With 95% Prediction Ellipse

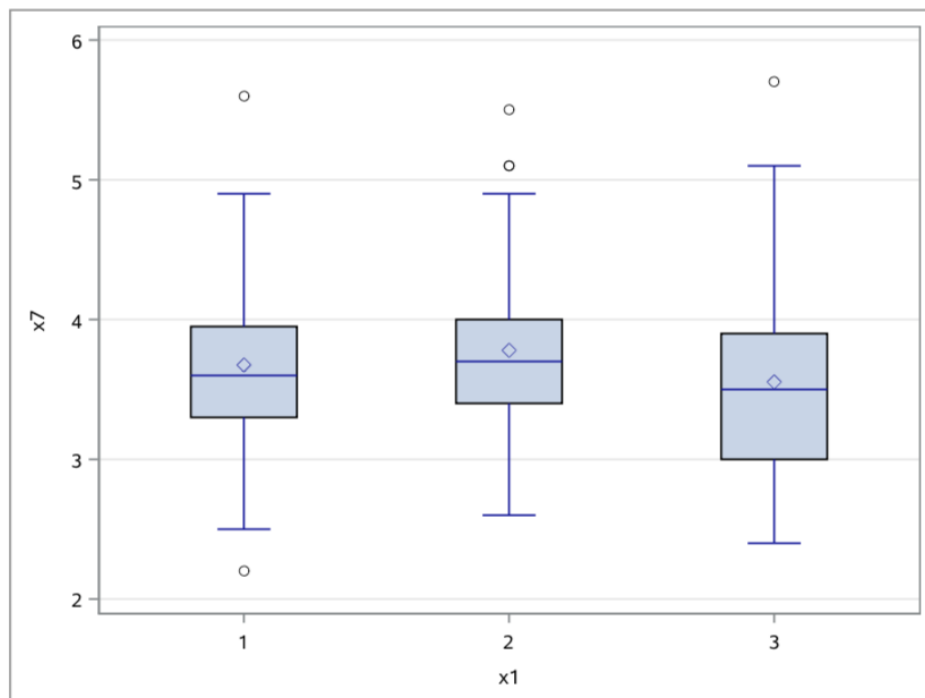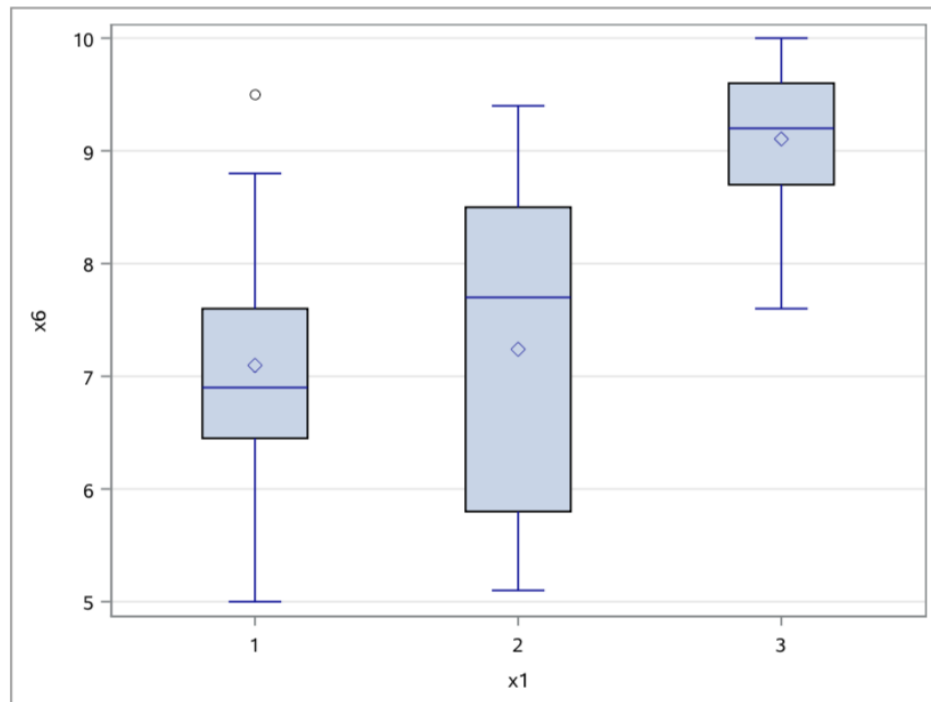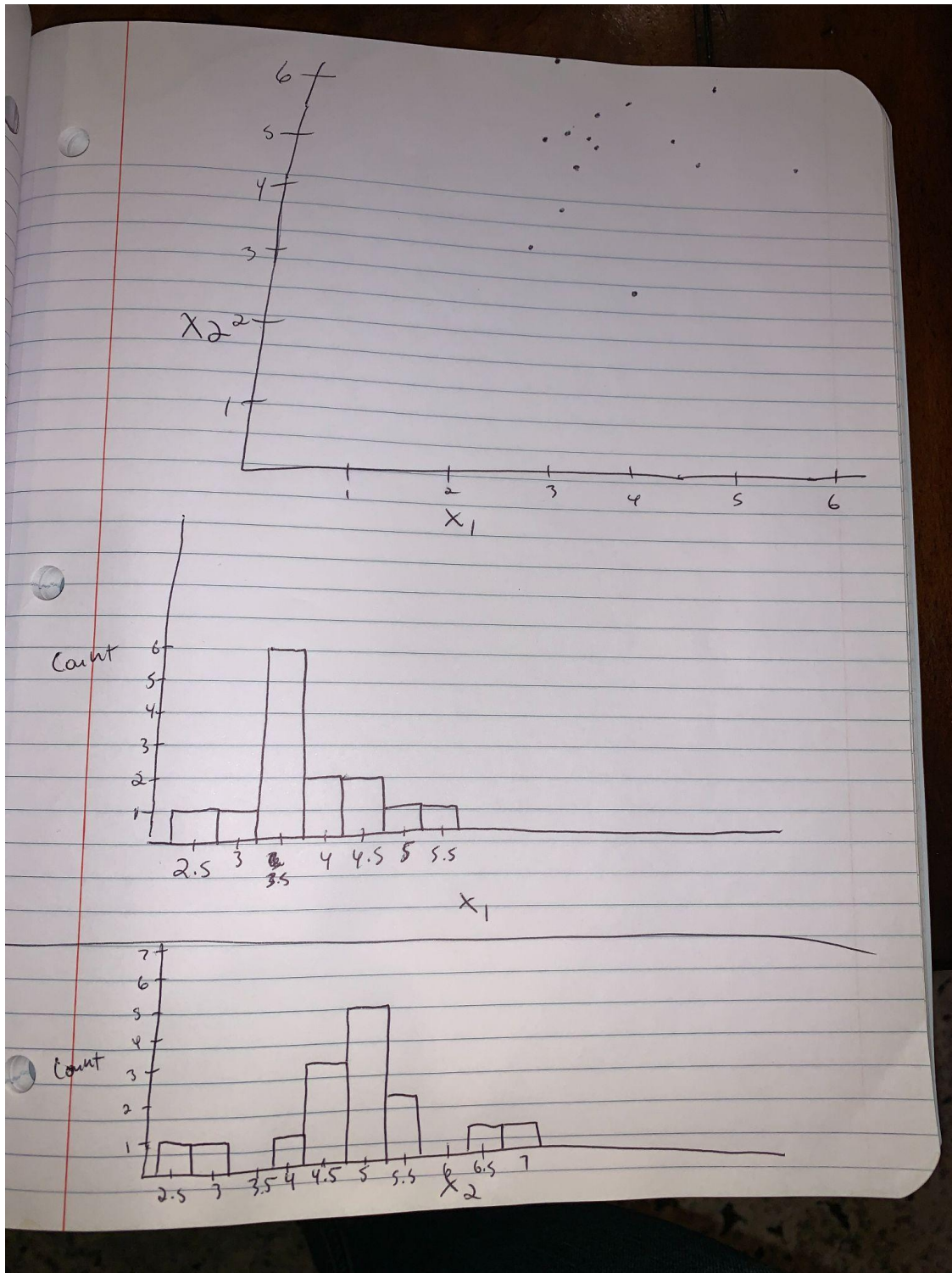| Observations | 100 |
| Correlation | 0.7915 |
| p-Value | <.0001 |

Based on the plots and the correlation numbers, x7 and x12 are the only ones that are correlated.

2f.





Yes, box plots can detect outliers. The 'whiskers' or error bars show the range of the data. If there are outliers, they will appear outside of the whiskers on the plot. In the top plot, you can see one outlier, and in the bottom there are five.

3.



X2 looks much more normal than X1 as there is much more natural distribution of numbers that fit the traditional bell curve shape.