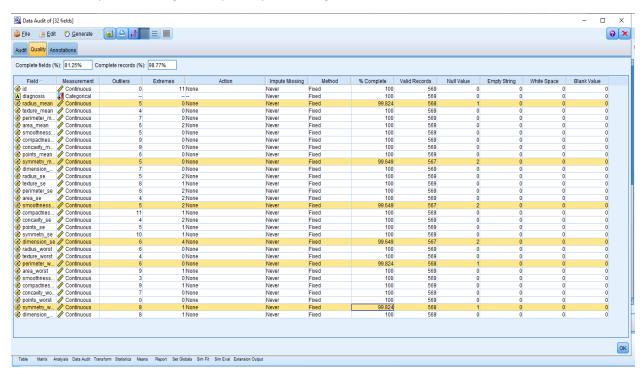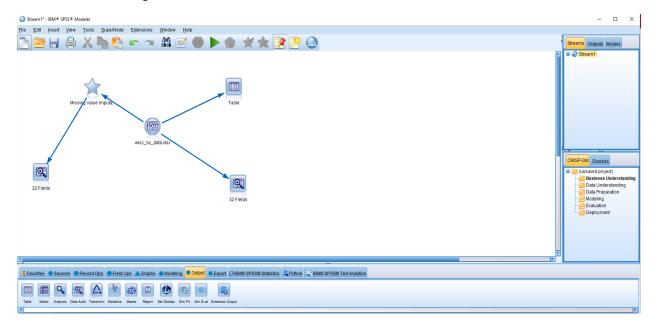**1. Explore whether there are missing values for any of the variables. Show how you did so in SPSS Modeler**
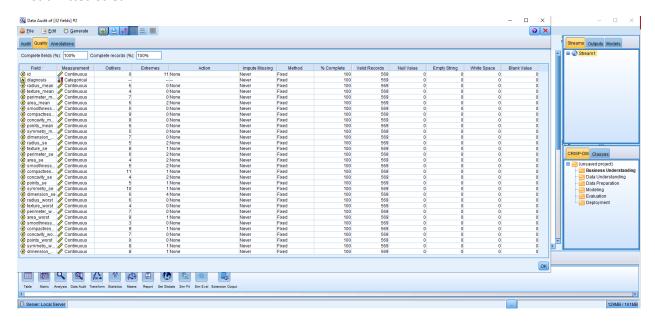
Below is the snapshot for all the missing variables from the given data set, we have used the features of Data Audit output table to get the quantity of missing variables.



Once the missing variables has been handled, then we have generated the missing value node to address these variables from given data set.

Below is the snapshot of data after handling the missing variables, complete detail of variable is available in submitted stream.



**2. Use a graph to determine visually whether there are any outliers in the average concavity field (concavity_mean) that measures the severity of concave portions of the contour. What kind of graph would provide better visualization for this task.**

Below is visual representation of the average concavity attribute (concavity_mean), from the below histogram graph we can say the there is no outliers available, we will also check with box plot graph.

Below is the box plot graph representation of the average concavity attribute (concavity_mean).



From a above box plot, we can confirm that there are no outliers available.



The above is the snapshot of stream till this task.

**3. Transform the radius_mean attribute using Z-score standardization. Using a graph, describe the range of the standardized values. Hint: consider adding a Derive node to create the Z-score. You can find the required summary statistics (mean and std deviation) to calculate the Z-score using the statistics node.**
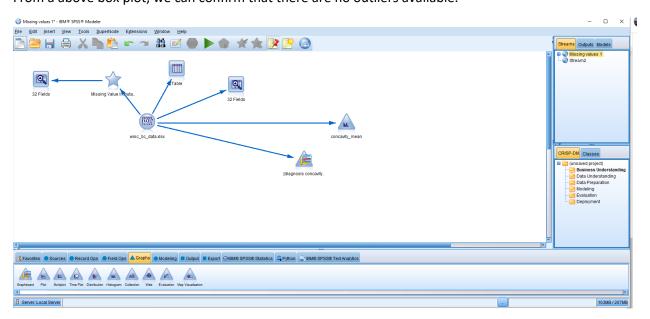
In order to calculate the Z-score standardization, we first have to calculate the mean and standard deviation which can be shown below.

| Statistics of [radius_mean] | | — □ ✕ |
|---|---|---|

| radius_mean Statistics | |
|---|---|
| Count | 568 |
| Mean | 14.125 |
| Min | 6.981 |
| Max | 28.110 |
| Range | 21.129 |
| Variance | 12.438 |
| Standard Deviation | 3.527 |
| Standard Error of Mean | 0.148 |

The above is mean and standard deviation for given data using statics, we will be using these values to calculate the asked z-score.



Histogram of Z-Score-Radius_mean #5

The above is the histogram graph representation of Z-score standardization, which mostly varies from -2 to 2.

We have transformed the radius_mean attribute using Z-score standardization (mean $\mu$ = 14.125) and standard deviation ($\sigma$) = 3.527. The Z-score ranges from about - 2.026 to 3.965 indicating the potential existence for several outliers.

**4.How would you transform all 30 metrics to Z-scores without having to create derived nodes for each of them? Hint: Investigate the Auto Data Prep node. We did not cover this in class, but it should be very simple with a little experimentation**.

In order to calculate the Z-score for all the 30 metrics, we first need to add and connect the Auto Data prep node, which is shown in the snap shot above.



From the edit option in Auto Data prep node, we have to select the custom analysis option, as shown above.



Under the fields option we can select the target and inputs field to calculate the required Z-score, we can select any attribute, to valid our calculation we have taken the radius_mean attribute whose Z-score we have already calculate in previous part.

Now under setting tab and in prepare Inputs and target section, we have to insert mean and standard deviation of the selected attribute, for this case we have, mean = 14.125 and standard deviation = 3.527.





Once we click on Analyze Data from the above tab, we get the required results of selected attribute, this case the Z-score is varies in the same range which we have calculated earlier.

Marist College | Fall 2020

**5. Investigate whether there are any correlated variables among the first 10 numeric fields (the ones with _mean prefix)**

| | Corelations | — □ × |
| --- | --- | --- |

File       Edit       Generate

Statistics   Annotations

Collapse All       Expand All

**radius_mean**
  Statistics
  Pearson Correlations

| texture_mean | 0.324 | Strong |
| --- | --- | --- |
| perimeter_mean | 0.998 | Strong |
| area_mean | 0.987 | Strong |
| smoothness_mean | 0.170 | Strong |
| compactness_mean | 0.506 | Strong |
| concavity_mean | 0.677 | Strong |
| points_mean | 0.823 | Strong |
| symmetry_mean | 0.150 | Strong |
| dimension_mean | -0.312 | Strong |

**texture_mean**
  Statistics
  Pearson Correlations

| radius_mean | 0.324 | Strong |
| --- | --- | --- |
| perimeter_mean | 0.330 | Strong |
| area_mean | 0.321 | Strong |
| smoothness_mean | -0.023 | Medium |
| compactness_mean | 0.237 | Strong |
| concavity_mean | 0.302 | Strong |
| points_mean | 0.293 | Strong |
| symmetry_mean | 0.072 | Strong |
| dimension_mean | -0.076 | Strong |

**perimeter_mean**
  Statistics
  Pearson Correlations

| radius_mean | 0.998 | Strong |
| --- | --- | --- |
| texture_mean | 0.330 | Strong |
| area_mean | 0.987 | Strong |
| smoothness_mean | 0.207 | Strong |
| compactness_mean | 0.557 | Strong |
| concavity_mean | 0.716 | Strong |
| points_mean | 0.851 | Strong |
| symmetry_mean | 0.185 | Strong |
| dimension_mean | -0.261 | Strong |

OK

## Corelations

File | Edit | Generate

**Statistics** | Annotations

Collapse All | Expand All

### area_mean
#### Statistics
#### Pearson Correlations

| | | |
|---|---|---|
| radius_mean | 0.987 | Strong |
| texture_mean | 0.321 | Strong |
| perimeter_mean | 0.987 | Strong |
| smoothness_mean | 0.177 | Strong |
| compactness_mean | 0.499 | Strong |
| concavity_mean | 0.686 | Strong |
| points_mean | 0.823 | Strong |
| symmetry_mean | 0.153 | Strong |
| dimension_mean | -0.283 | Strong |

### smoothness_mean
#### Statistics
#### Pearson Correlations

| | | |
|---|---|---|
| radius_mean | 0.170 | Strong |
| texture_mean | -0.023 | Medium |
| perimeter_mean | 0.207 | Strong |
| area_mean | 0.177 | Strong |
| compactness_mean | 0.659 | Strong |
| concavity_mean | 0.522 | Strong |
| points_mean | 0.554 | Strong |
| symmetry_mean | 0.558 | Strong |
| dimension_mean | 0.585 | Strong |

### compactness_mean
#### Statistics
#### Pearson Correlations

| | | |
|---|---|---|
| radius_mean | 0.506 | Strong |
| texture_mean | 0.237 | Strong |
| perimeter_mean | 0.557 | Strong |
| area_mean | 0.499 | Strong |
| smoothness_mean | 0.659 | Strong |
| concavity_mean | 0.883 | Strong |
| points_mean | 0.831 | Strong |
| symmetry_mean | 0.603 | Strong |
| dimension_mean | 0.565 | Strong |

OK

**Corelations** — □ ×

File    Edit    Generate

Statistics    Annotations

Collapse All    Expand All

concavity_mean
- Statistics
- Pearson Correlations

| radius_mean | 0.677 | Strong |
|---|---|---|
| texture_mean | 0.302 | Strong |
| perimeter_mean | 0.716 | Strong |
| area_mean | 0.686 | Strong |
| smoothness_mean | 0.522 | Strong |
| compactness_mean | 0.883 | Strong |
| points_mean | 0.921 | Strong |
| symmetry_mean | 0.501 | Strong |
| dimension_mean | 0.337 | Strong |

points_mean
- Statistics
- Pearson Correlations

| radius_mean | 0.823 | Strong |
|---|---|---|
| texture_mean | 0.293 | Strong |
| perimeter_mean | 0.851 | Strong |
| area_mean | 0.823 | Strong |
| smoothness_mean | 0.554 | Strong |
| compactness_mean | 0.831 | Strong |
| concavity_mean | 0.921 | Strong |
| symmetry_mean | 0.464 | Strong |
| dimension_mean | 0.167 | Strong |

symmetry_mean
- Statistics
- Pearson Correlations

| radius_mean | 0.150 | Strong |
|---|---|---|
| texture_mean | 0.072 | Strong |
| perimeter_mean | 0.185 | Strong |
| area_mean | 0.153 | Strong |
| smoothness_mean | 0.558 | Strong |
| compactness_mean | 0.603 | Strong |
| concavity_mean | 0.501 | Strong |
| points_mean | 0.464 | Strong |
| dimension_mean | 0.479 | Strong |

dimension_mean
- Statistics
- Pearson Correlations

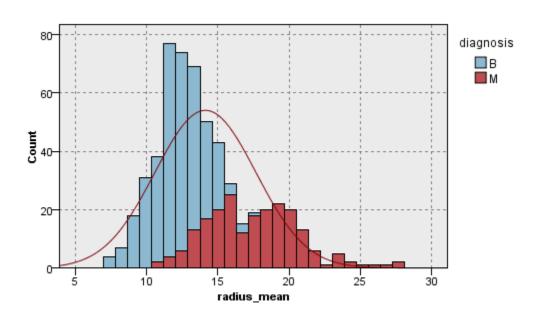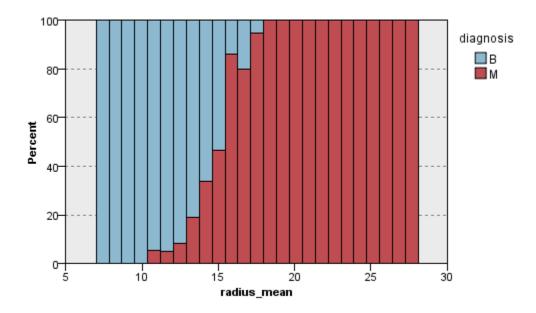| radius_mean | -0.312 | Strong |
|---|---|---|
| texture_mean | -0.076 | Strong |
| perimeter_mean | -0.261 | Strong |
| area_mean | -0.283 | Strong |
| smoothness_mean | 0.585 | Strong |
| compactness_mean | 0.565 | Strong |
| concavity_mean | 0.337 | Strong |
| points_mean | 0.167 | Strong |
| symmetry_mean | 0.479 | Strong |

OK

From the above we can say that there is no corelated variables.

**6. Construct a normalized histogram of radius_mean, with an overlay of the target variable diagnosis. Explain the results.**





The above normalized histogram graph of the required attribute (radius_mean), the diagnosis is coded as "M" to indicate malignant, or "B" to indicate benign. The graph shows that radius_mean for (M) malignant are higher then as compared to (B) benign, these features relate to the shape and size of the cell nuclei.