## Program – 4

## AIM: Explain missing value treatment and outlier treatment in WEKA on sample dataset.
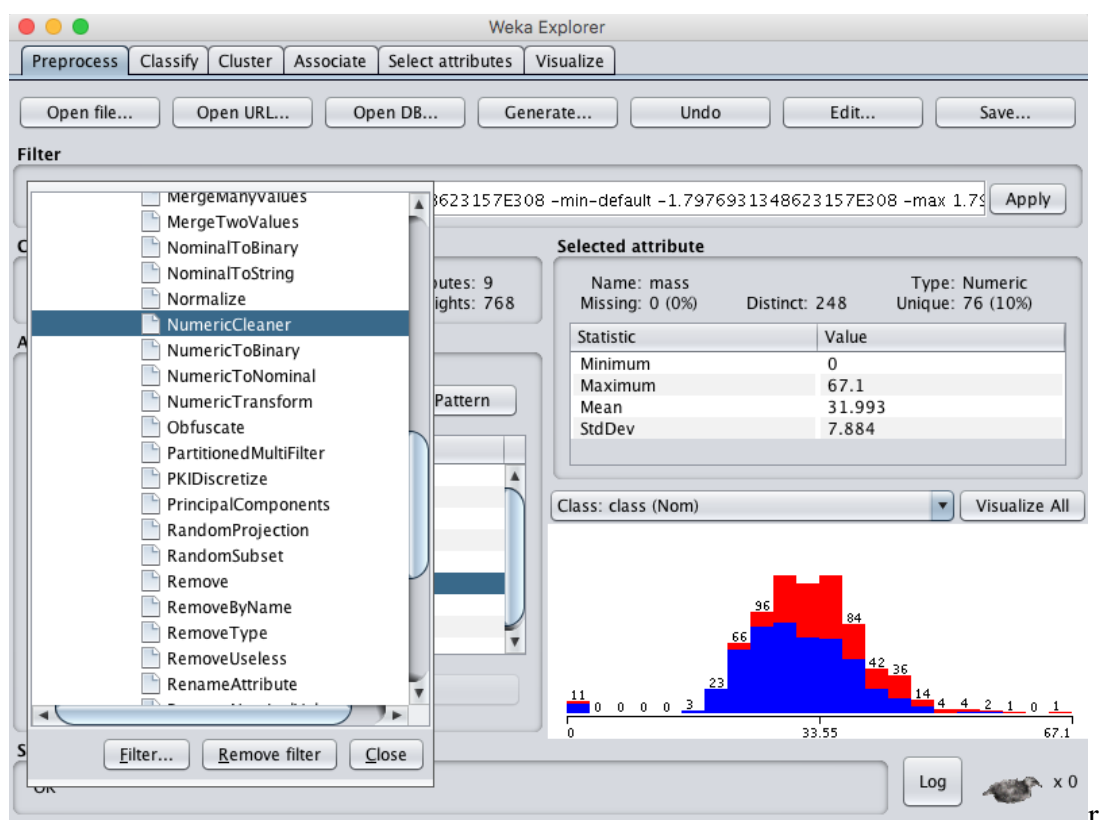
### Introduction and Theory

Data is rarely clean and often you can have corrupt or missing values. It is important to identify, mark and handle missing data when developing machine learning models in order to get the very best performance.

In most cases it is normal to find some values missing in the data, in other cases there may be some data elements whose values don't match with any other data points, such values are called outliers, normally any value outside of 1.5 quantiles of the data mean is taken as an outlier.
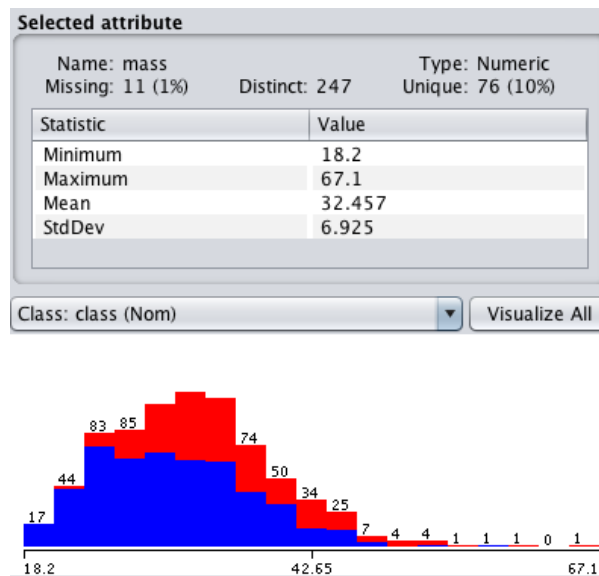
## Procedure
### Mark Missing Values
1. Open the Weka Explorer.

2. Load the Pima Indians onset of diabetes dataset.

3. Click the "Choose" button for the Filter and select NumericalCleaner, it us under

   unsupervised.attribute.NumericalCleaner



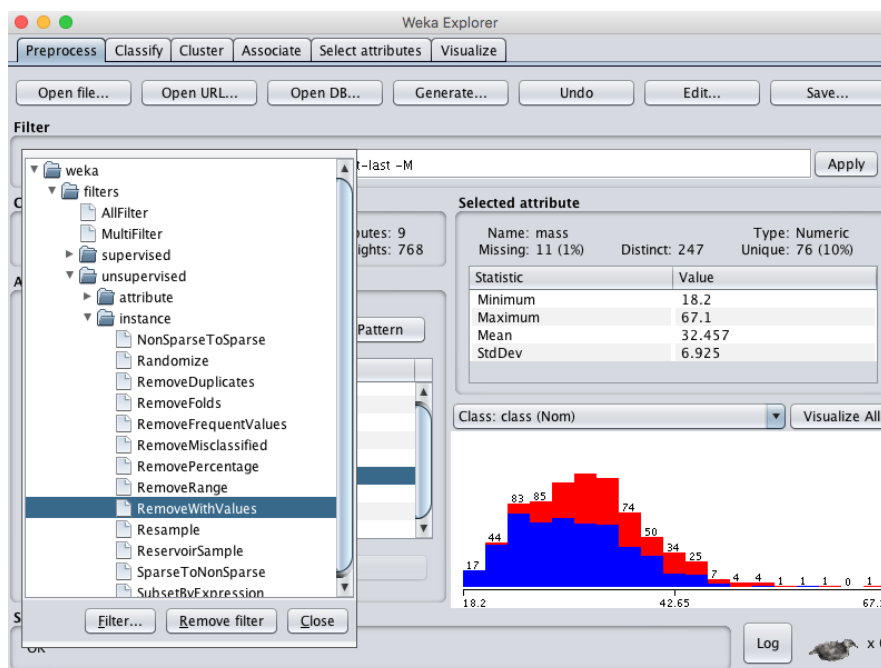4. Click on the filter to configure it.

# Program – 4

5. Set the attributeIndicies to 6, the index of the mass attribute.

6. Set minThreshold to 0.1E-8 (close to zero), which is the minimum value allowed for the attribute.

7. Set minDefault to NaN, which is unknown and will replace values below the threshold.

8. Click the "OK" button on the filter configuration.

9. Click the "Apply" button to apply the filter.



## Removing missing values

1. Click the "Choose" button for the Filter and select RemoveWithValues, it us under unsupervized.instance.RemoveWithValues.



2. Click on the filter to configure it.

# Program – 4

3. Set the *attributeIndicies* to 6, the index of the mass attribute.

4. Set *matchMissingValues* to "True".

5. Click the "OK" button to use the configuration for the filter.

6. Click the "Apply" button to apply the filter.

**Selected attribute**

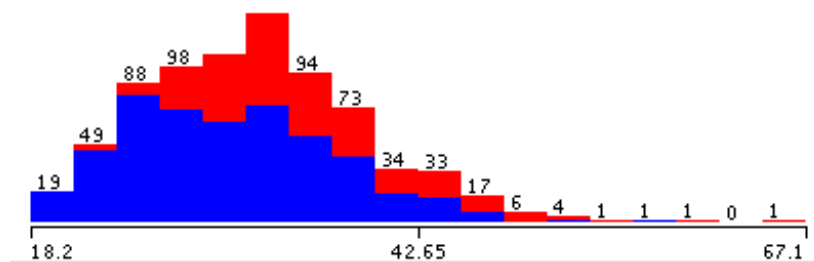| | |
|---|---|
| Name: mass | Type: Numeric |
| Missing: 0 (0%)  Distinct: 247 | Unique: 76 (10%) |

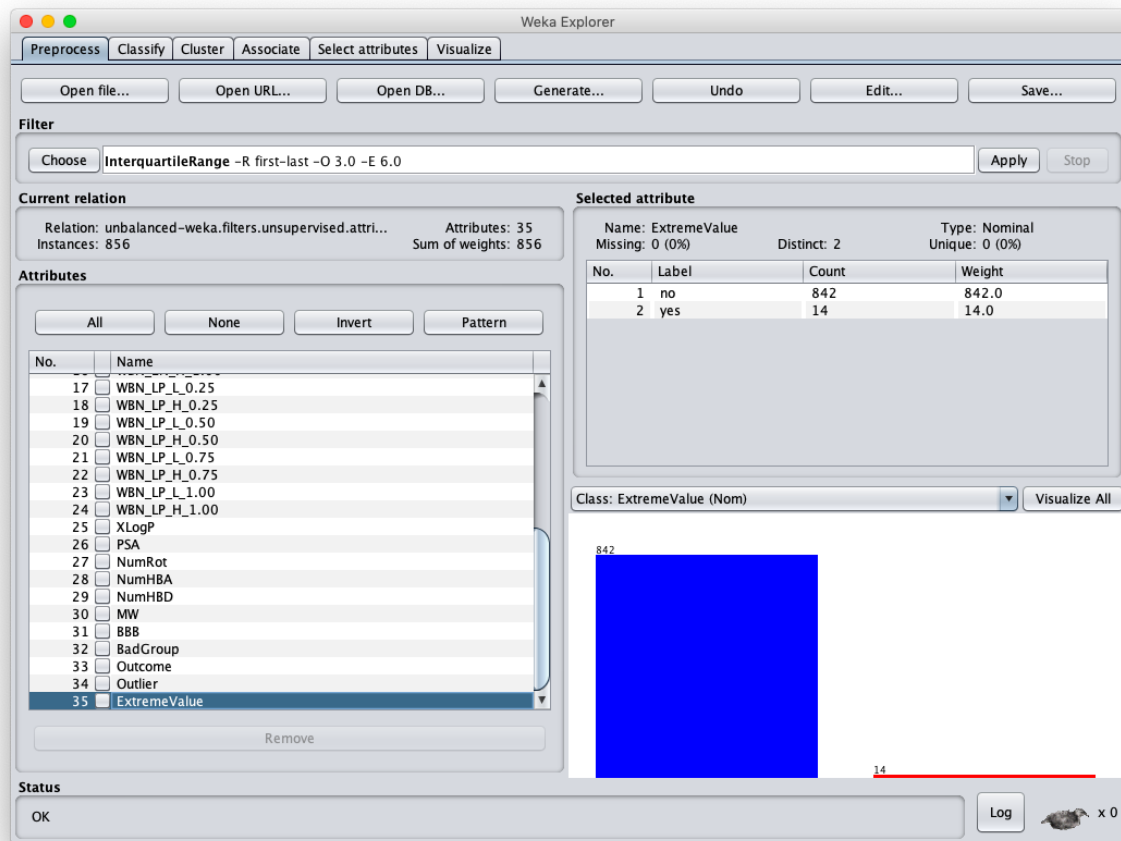| Statistic | Value |
|---|---|
| Minimum | 18.2 |
| Maximum | 67.1 |
| Mean | 32.457 |
| StdDev | 6.925 |

Class: class (Nom) ▼   Visualize All



## Outlier treatment

1. Open the Weka Explorer.

2. Load the Unbalanced dataset.

3. Click the "Choose" button for the Filter and select InterquantileRange, it us under unsupervized.attribute.InterquantileRange

4. Click on the filter to configure it.

5. Click the "OK" button on the filter configuration.

6. Click the "Apply" button to apply the filter.

# Program – 4



## Findings and Learnings:

1. Missing value treatment and outlier removal are very important steps to be performed before data analysis and statistical modeling.
2. Weka Software makes it easy and simple to visualize data and apply appropriate treatments to remove the anomalies.