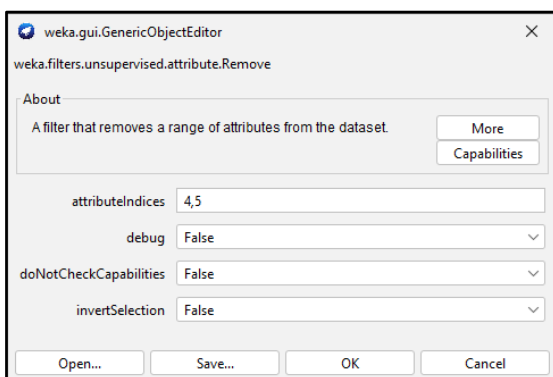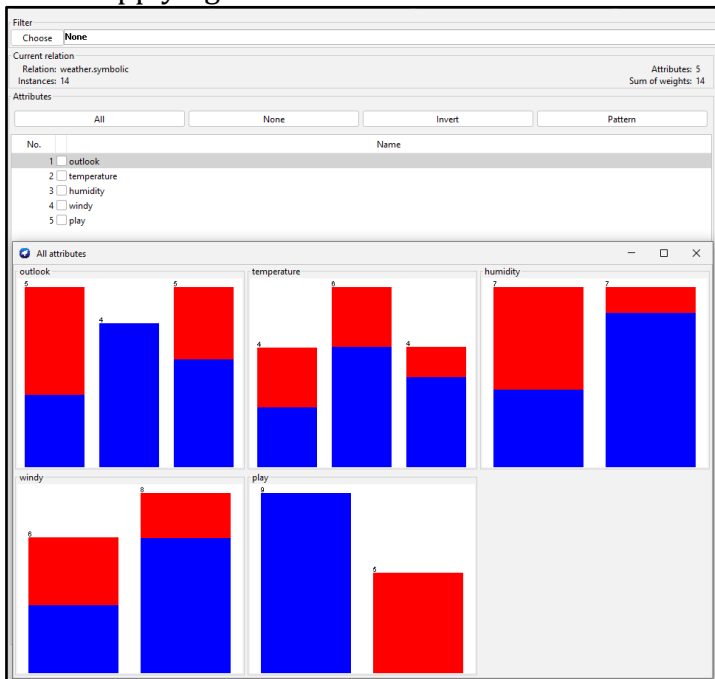## Experiment 4

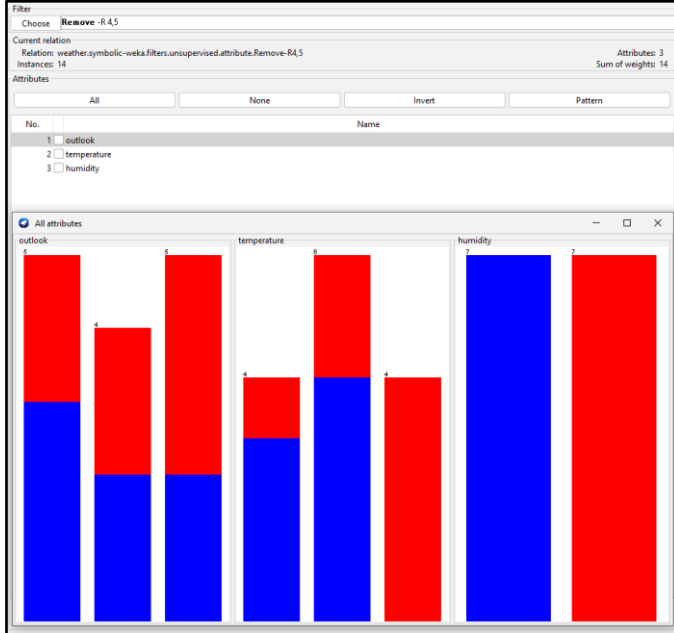**Title**: **Apply Preprocessing techniques on dataset using filters: Remove, ReplaceMissingValues, ReplaceMissingWithUserConstant, ReplaceWithMissingValue, Descritize. Also do the result analysis before and after preprocessing.**

1. Remove filter: The Remove filter is an unsupervised attribute filter that allows users to delete specific columns from a dataset. It is particularly useful when you want to exclude irrelevant, redundant or sensitive attributes from the data before applying machine learning algorithms.

Before applying the "Remove filter"

after applying the "Remove filter"



2. ReplaceMissingValues filter: The ReplaceMissingValues filter is an unsupervised attribute and instance filter used to automatically fill in missing values in dataset. It ensures that incomplete data does not negatively impact the performance of machine learning algorithms.
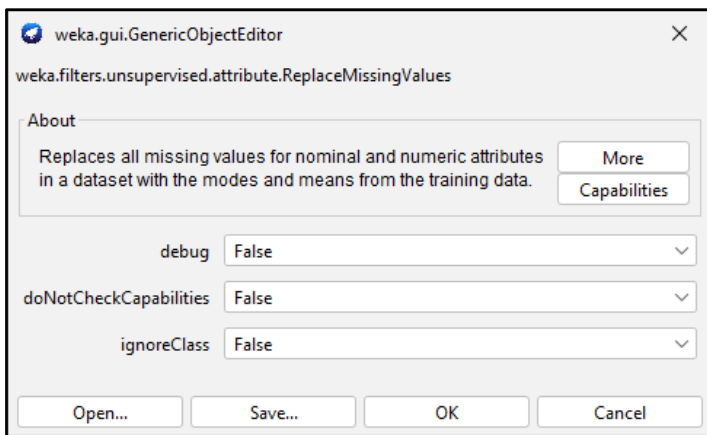
Before applying the "ReplaceMissingValues filter"

After applying the "ReplaceMissingValues filter"



| No. | 1: handicapped-infants Nominal | 2: water-project-cost-sharing Nominal | 3: adoption-of-the-budget-resolution Nominal | 4: physician-fee-freeze Nominal | 5: el-salvador-aid Nominal | 6: religious-groups-in-schools Nominal | 7: ant |
|---|---|---|---|---|---|---|---|
| 1 | n | y | n | y | y | y | n |
| 2 | n | y | n | y | y | y | n |
| 3 | n | y | y | n | y | y | n |
| 4 | n | y | y | n | y | y | n |
| 5 | y | y | y | n | y | y | n |
| 6 | n | y | y | n | y | y | n |
| 7 | n | y | n | y | y | y | n |
| 8 | n | y | n | y | y | y | n |
| 9 | n | y | n | y | y | y | n |
| 10 | y | y | y | n | n | n | y |
| 11 | n | y | n | y | y | n | n |
| 12 | n | y | n | y | y | y | n |
| 13 | n | y | y | n | n | n | y |
| 14 | y | y | y | n | n | y | y |
| 15 | n | y | n | y | y | y | n |
| 16 | n | y | n | y | y | y | n |
| 17 | y | n | y | n | n | y | n |
| 18 | y | y | y | n | n | n | y |
| 19 | n | y | n | y | y | y | n |
| 20 | y | y | y | n | n | n | y |
| 21 | y | y | y | n | n | y | y |
| 22 | y | y | y | n | n | n | y |
| 23 | y | y | y | n | n | n | y |

3. ReplaceMissingWithUserConstant filter: The <u>Replace Missing With User Constant</u> is an unsupervised attribute filter used to replace all missing values in a dataset with a constant value specified by the user. Allow users to define the custom constant to replace.

Before applying the "ReplaceMissingWithUserConstant filter"

After applying the "ReplaceMissingWithUserConstant filter"



4. Descritize filter: The Descritize filter is an unsupervised attribute into nominal attributes by dividing the numeric range into discrete intervals or bins. Descritization helps simplify data and can improve model performance in certain scenarios.

Before applying the "Descritize filter"

After applying the "Descritize filter"

**Viewer**

Relation: pima_diabetes-weka.filters.unsupervised.attribute.Discretize-B2-M-1.0-Rfirst-last-precision2

| No. | 1: preg<br>Nominal | 2: plas<br>Nominal | 3: pres<br>Nominal | 4: skin<br>Nominal | 5: insu<br>Nominal | 6: mass<br>Nominal | 7: pedi<br>Nominal | 8: age<br>Nominal | 9: class<br>Nominal |
|---|---|---|---|---|---|---|---|---|---|
| 1 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(33.55-inf)' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 2 | '(-inf-8.5]' | '(-inf-99.5]' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(-inf-51]' | tested_negative |
| 3 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 4 | '(-inf-8.5]' | '(-inf-99.5]' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(-inf-51]' | tested_negative |
| 5 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 6 | '(-inf-8.5]' | '(-inf-99.5]' | '(-inf-61]' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 7 | '(8.5-inf)' | '(99.5-inf)' | '(-inf-61]' | '(-inf-49.5]' | '(-inf-423]' | '(33.55-inf)' | '(-inf-1.25]' | '(-inf-51]' | tested_negative |
| 8 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(423-inf)' | '(-inf-33.55]' | '(-inf-1.25]' | '(51-inf)' | tested_positive |
| 9 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(51-inf)' | tested_positive |
| 10 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(33.55-inf)' | '(-inf-1.25]' | '(-inf-51]' | tested_negative |
| 11 | '(8.5-inf)' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(33.55-inf)' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 12 | '(-inf-8.5]' | '(99.5-inf)' | '(-inf-61]' | '(-inf-49.5]' | '(423-inf)' | '(-inf-33.55]' | '(-inf-1.25]' | '(51-inf)' | tested_positive |
| 13 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 14 | '(-inf-8.5]' | '(99.5-inf)' | '(-inf-61]' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 15 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(33.55-inf)' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 16 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 17 | '(-inf-8.5]' | '(99.5-inf)' | '(-inf-61]' | '(-inf-49.5]' | '(-inf-423]' | '(33.55-inf)' | '(-inf-1.25]' | '(-inf-51]' | tested_negative |
| 18 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(33.55-inf)' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 19 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(33.55-inf)' | '(-inf-1.25]' | '(-inf-51]' | tested_negative |
| 20 | '(-inf-8.5]' | '(-inf-99.5]' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(33.55-inf)' | '(-inf-1.25]' | '(-inf-51]' | tested_negative |
| 21 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(33.55-inf)' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 22 | '(8.5-inf)' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 23 | '(8.5-inf)' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 24 | '(8.5-inf)' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 25 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(33.55-inf)' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |
| 26 | '(-inf-8.5]' | '(-inf-99.5]' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(-inf-51]' | tested_negative |
| 27 | '(8.5-inf)' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(51-inf)' | tested_negative |
| 28 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(33.55-inf)' | '(-inf-1.25]' | '(-inf-51]' | tested_negative |
| 29 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(33.55-inf)' | '(-inf-1.25]' | '(51-inf)' | tested_negative |
| 30 | '(-inf-8.5]' | '(99.5-inf)' | '(61-inf)' | '(-inf-49.5]' | '(-inf-423]' | '(-inf-33.55]' | '(-inf-1.25]' | '(-inf-51]' | tested_positive |