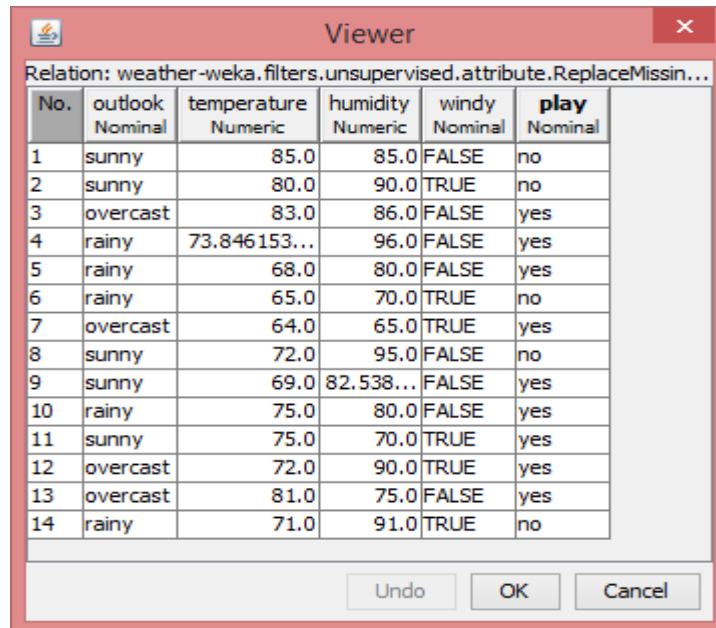# Lab Exercise One
# Data Preprocessing with WEKA Explorer

Binning With Filter

Unsupervised Attribute Filter – Discretize: This filter converts numeric attributes to nominal use equal-width (default) or equal-depth (frequency) binning.

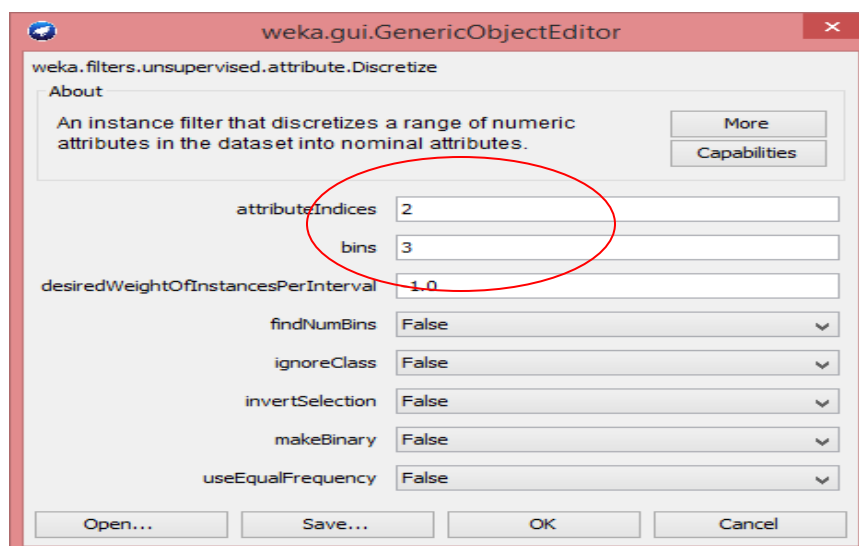1. Open the dataset **weather.numeric**. First, replace missing values with the filter.



2. Choose filter **Discretize** from the drop-down list of **unsupervised attribute** filters and then left-click to open its properties window. We want to perform **equal-width** binning on 2nd attribute – **temperature** with three bins.

3. Click **Apply** button. Then select **temperature** attribute to check the results.

4. To perform equal-depth (frequency) binning on the 3rd attribute **humidity**, we choose filter **PKIDiscretize** from the drop-down list of **unsupervised attribute** filters and then left-click to open its properties window. This filter use the square root of the number of values as the number of bins.
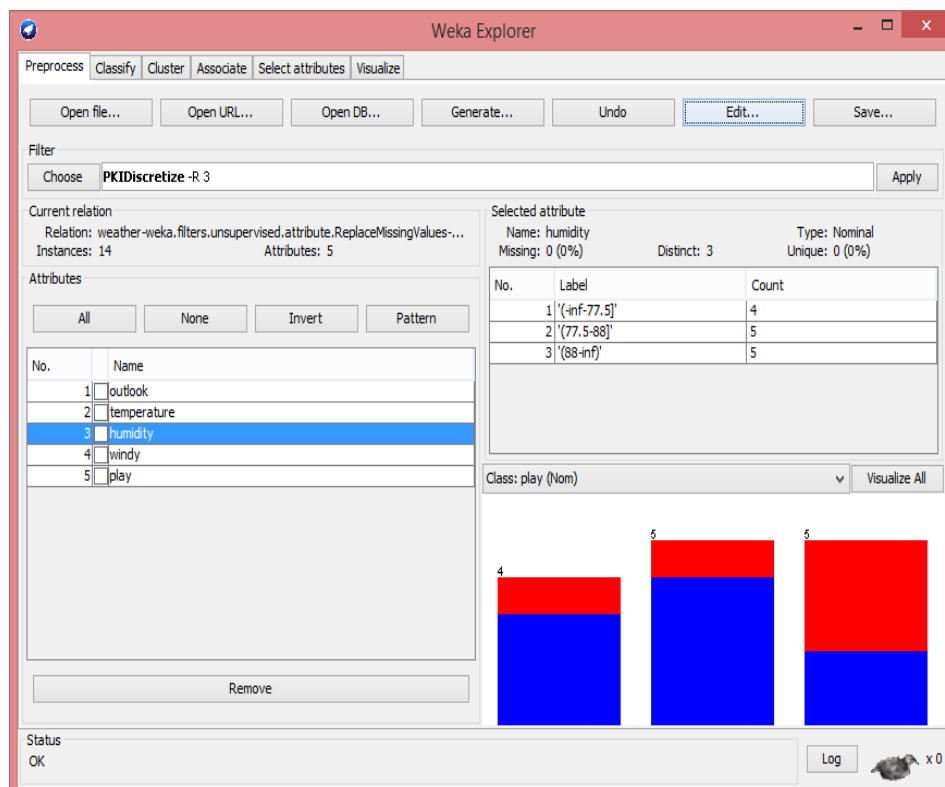
**Viewer**

Relation: weather-weka.filters.unsupervised.attribute.ReplaceMissing...

| No. | outlook Nominal | temperature Nominal | humidity Nominal | windy Nominal | play Nominal |
|---|---|---|---|---|---|
| 1 | sunny | '(78-inf)' | '(77.5-... | FALSE | no |
| 2 | sunny | '(78-inf)' | '(88-inf)' | TRUE | no |
| 3 | overcast | '(78-inf)' | '(77.5-... | FALSE | yes |
| 4 | rainy | '(71-78]' | '(88-inf)' | FALSE | yes |
| 5 | rainy | '(-inf-71]' | '(77.5-... | FALSE | yes |
| 6 | rainy | '(-inf-71]' | '(-inf-7... | TRUE | no |
| 7 | overcast | '(-inf-71]' | '(-inf-7... | TRUE | yes |
| 8 | sunny | '(71-78]' | '(88-inf)' | FALSE | no |
| 9 | sunny | '(-inf-71]' | '(77.5-... | FALSE | yes |
| 10 | rainy | '(71-78]' | '(77.5-... | FALSE | yes |
| 11 | sunny | '(71-78]' | '(-inf-7... | TRUE | yes |
| 12 | overcast | '(71-78]' | '(88-inf)' | TRUE | yes |
| 13 | overcast | '(78-inf)' | '(-inf-7... | FALSE | yes |
| 14 | rainy | '(-inf-71]' | '(88-inf)' | TRUE | no |

Undo    OK    Cancel