

Data Mining

Assignment 3

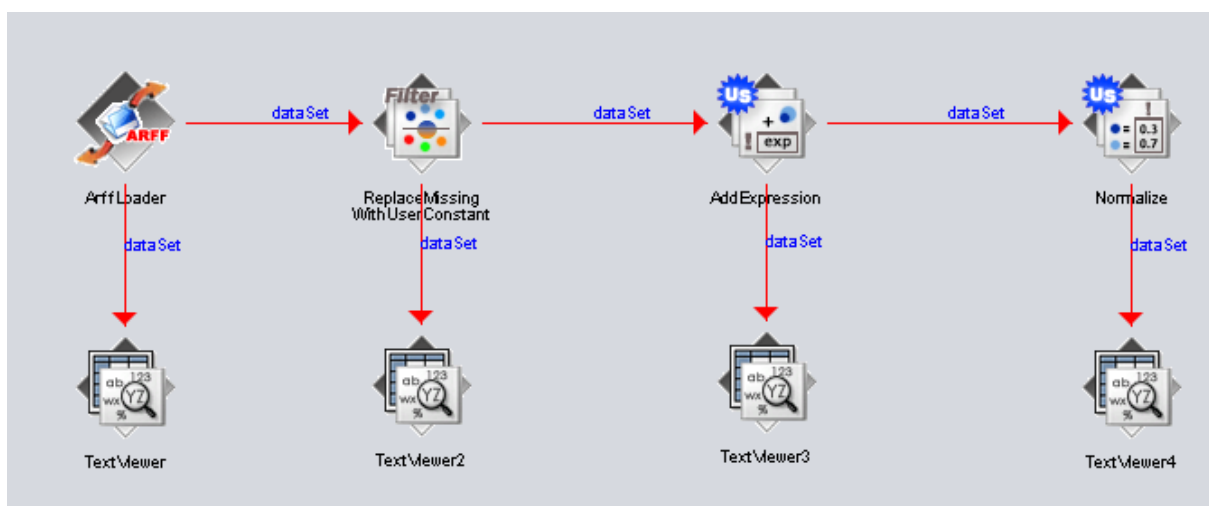
Name: Ayush Kumar

Reg. no.: 20214284


Group: CSE-6B1

Q1. Use of filters in Weka Knowledge Flow: Use Glass dataset and do the following operations using appropriate filters.

```
glass_arff.arff
File Edit View
@relation Glass
@attribute 'RI' real
@attribute 'Na' real
@attribute 'Mg' real
@attribute 'Al' real
@attribute 'Si' real
@attribute 'K' real
@attribute 'Ca' real
@attribute 'Ba' real
@attribute 'Fe' real
@attribute 'Type' { 'build wind float', 'build wind non-float', 'vehic wind float', 'vehic wind non-float', containers, tableware, headlamps}
@data
1.51793,12.79,3.5,1.12,73.03,0.64,8.77,0,0,'build wind float'
1.51643,12.16,3.52,1.35,72.89,0.57,8.53,0,0,'vehic wind float'
1.51793,13.21,3.48,1.41,72.64,0.59,8.43,0,0,'build wind float'
1.51299,14.4,1.74,1.54,74.55,0,7.59,0,0,'tableware'
1.53393,12.3,0,1,70.16,0.12,16.19,0,0.24,'build wind non-float'
1.51655,12.75,2.85,1.44,73.27,0.57,8.79,0,11,0.22,'build wind non-float'
1.51779,13.64,3.65,0.65,73.06,8.93,0,0,'vehic wind float'
1.51837,13.14,2.84,1.28,72.85,0.55,9.07,0,0,'build wind float'
1.51545,14.14,0,2.68,73.39,0.08,9.07,0.61,0.05,'headlamps'
1.51780,13.10,2.0,1.3,72.33,0.55,8.44,0,0.28,'build wind non-float'
Ln 107, Col 1 | 17,823 characters | 100% | Unix (LF) | UTF-8
```



a. Mark missing values with name NaN.



ReplaceMissingWithUserConstant options

About

Replaces all missing values for nominal, string, numeric and date attributes in the dataset with user-supplied constant values.

More

Capabilities

attributes

10

dateFormat

yyyy-MM-dd'THH:mm:ss

dateReplacementValue

debug

False

doNotCheckCapabilities

False

ignoreClass

False

nominalStringReplacementValue

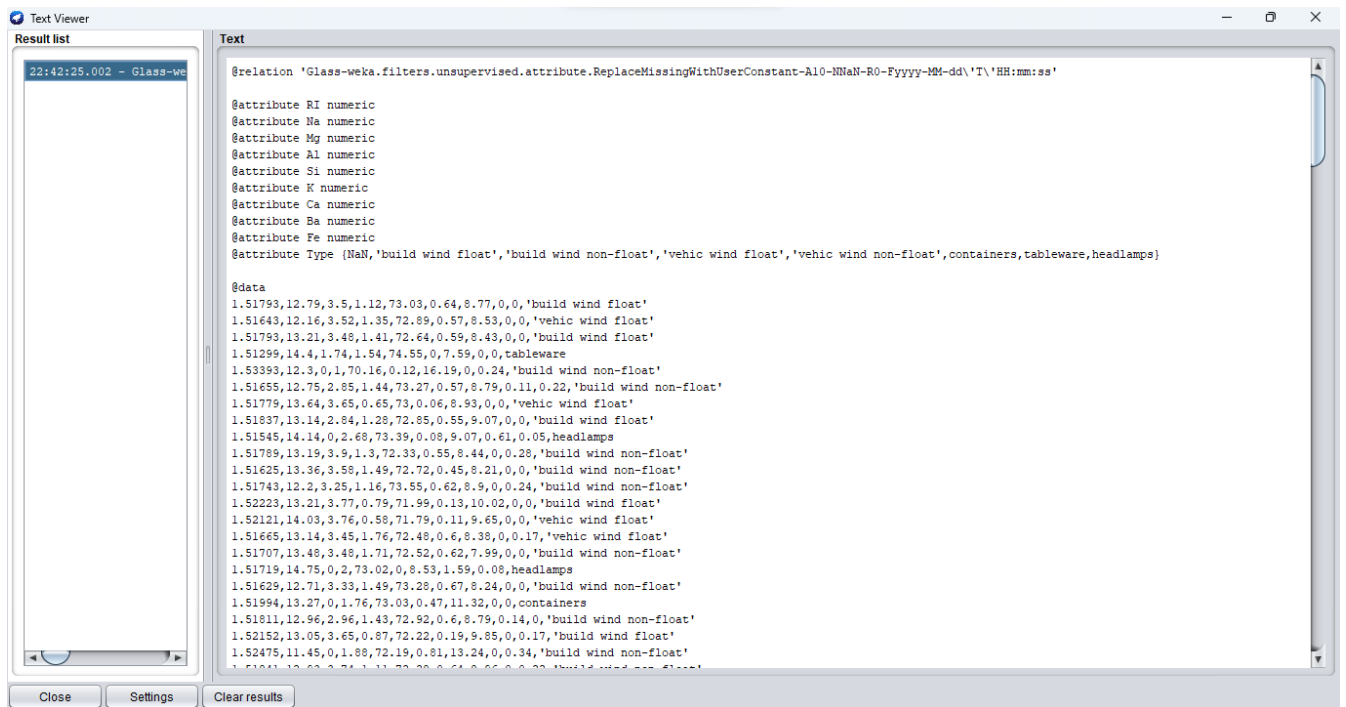
NaN

numericReplacementValue

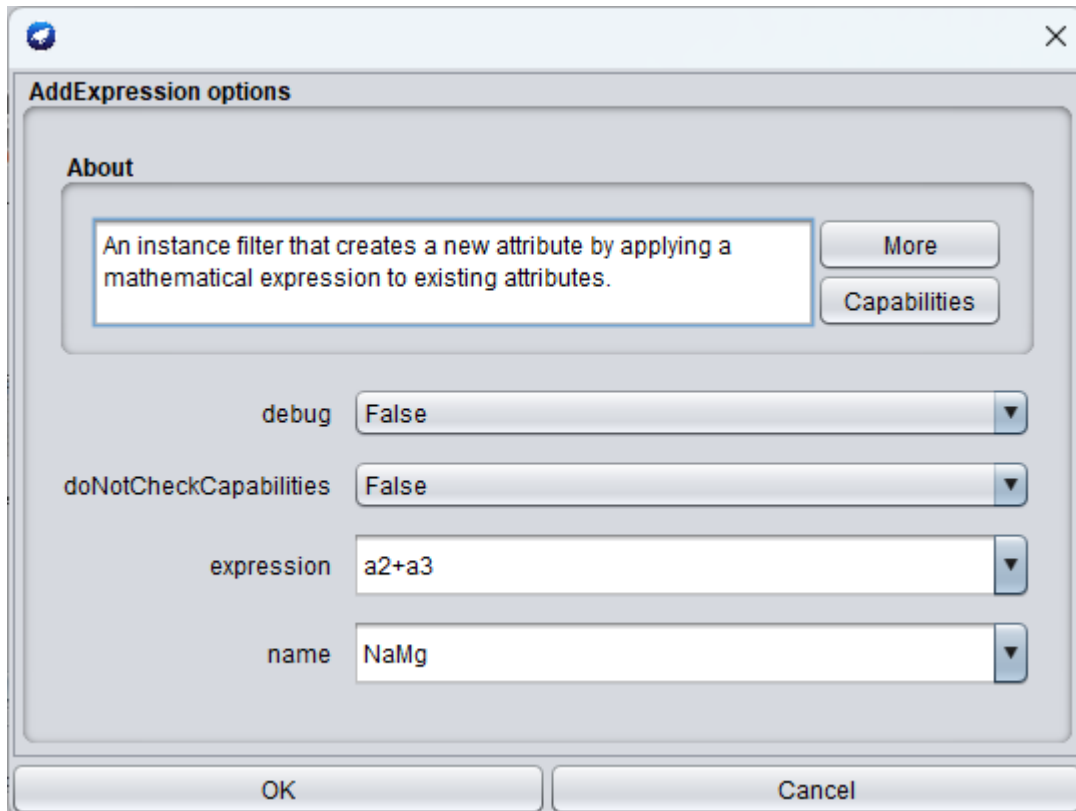
0

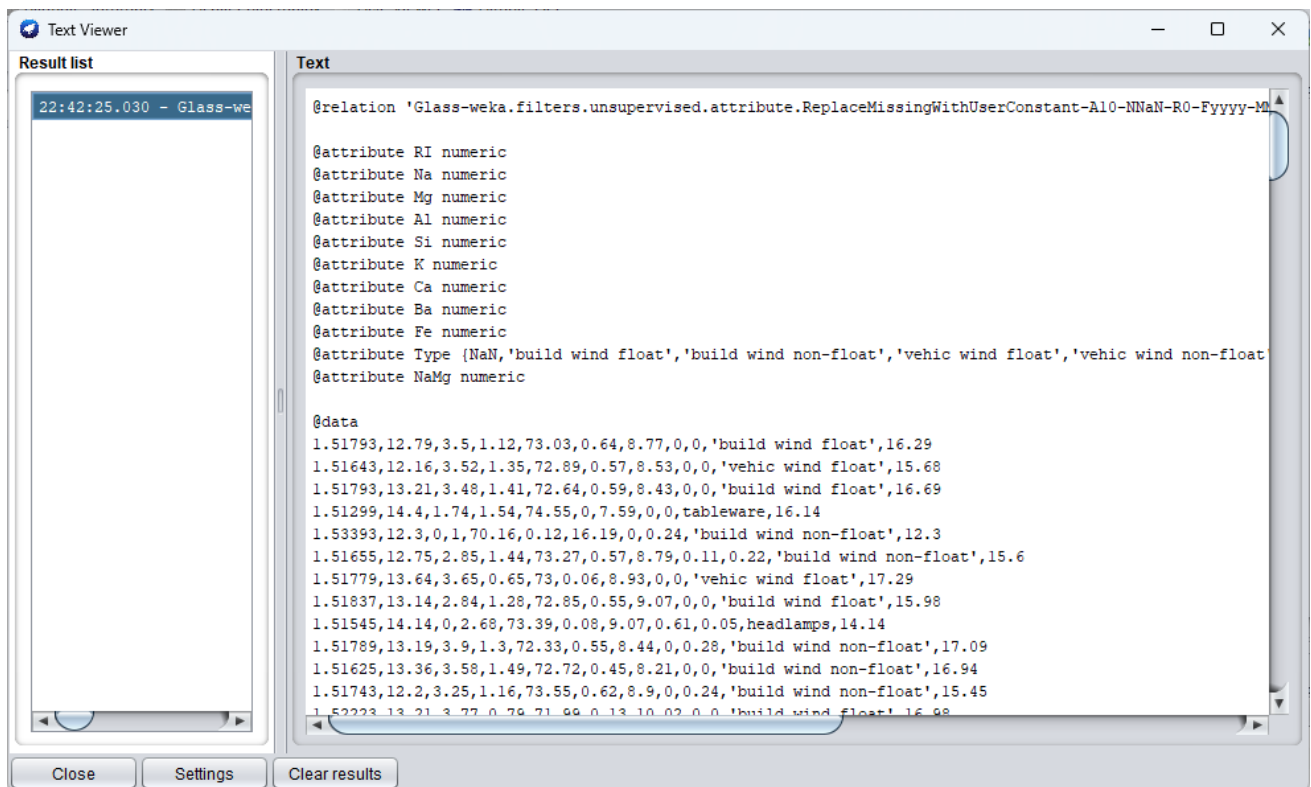
OK

Cancel

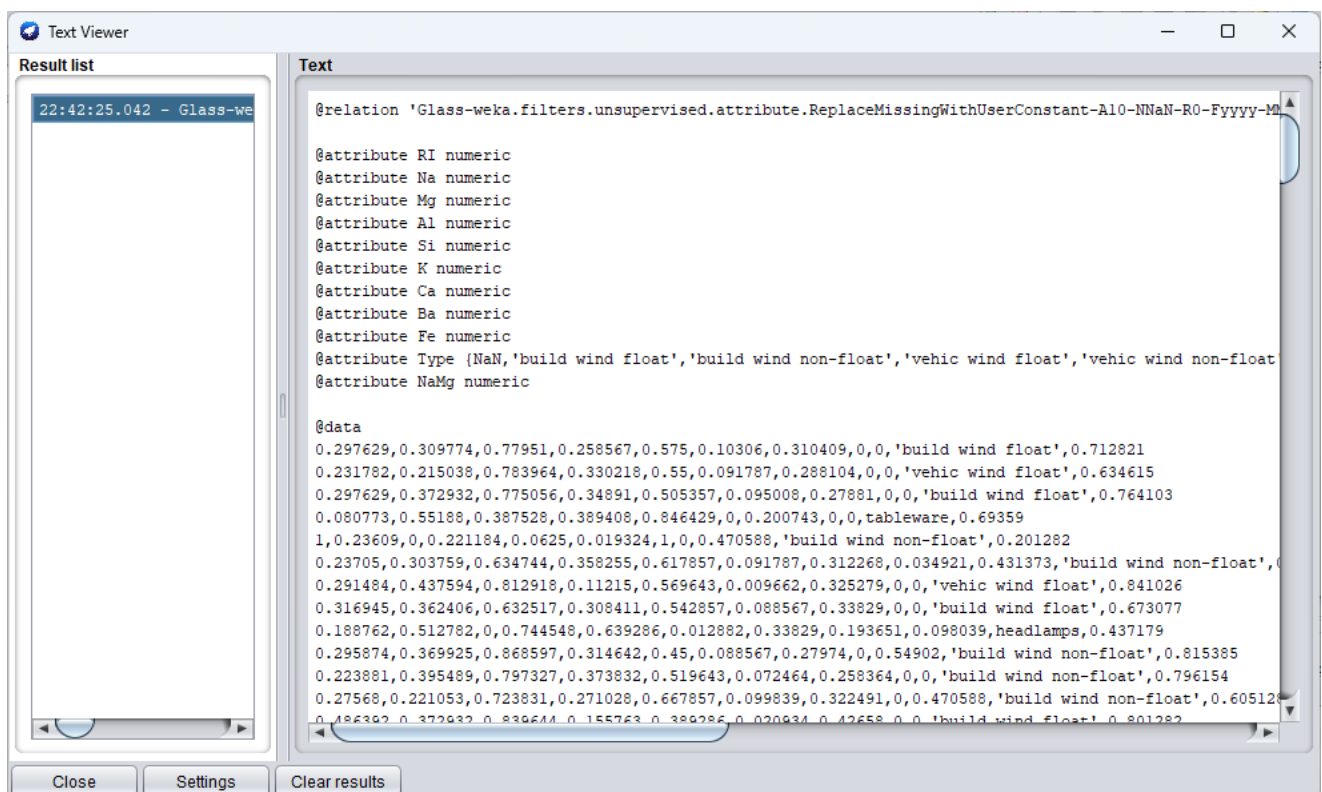


b. Merge any two similar numeric types attribute into new attributes.

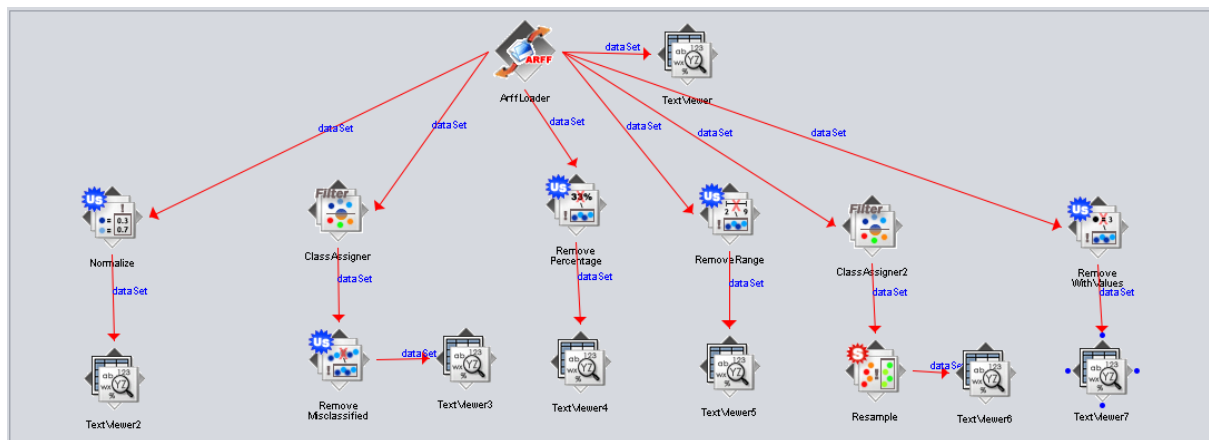
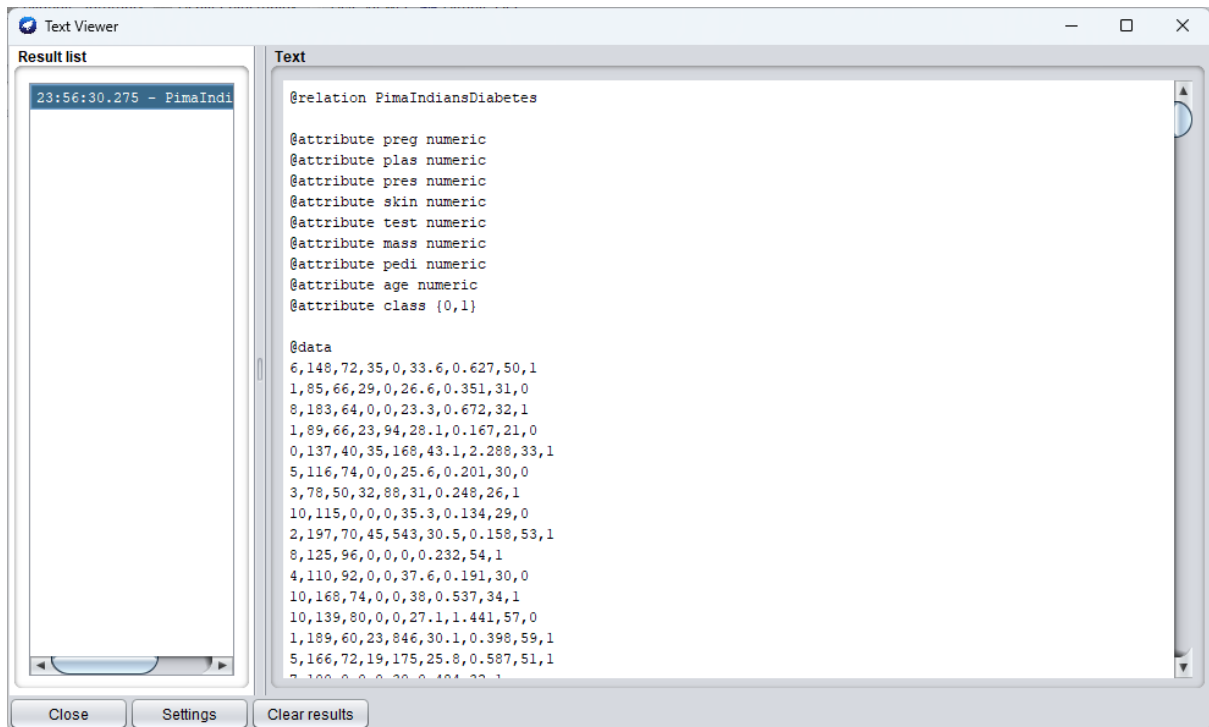




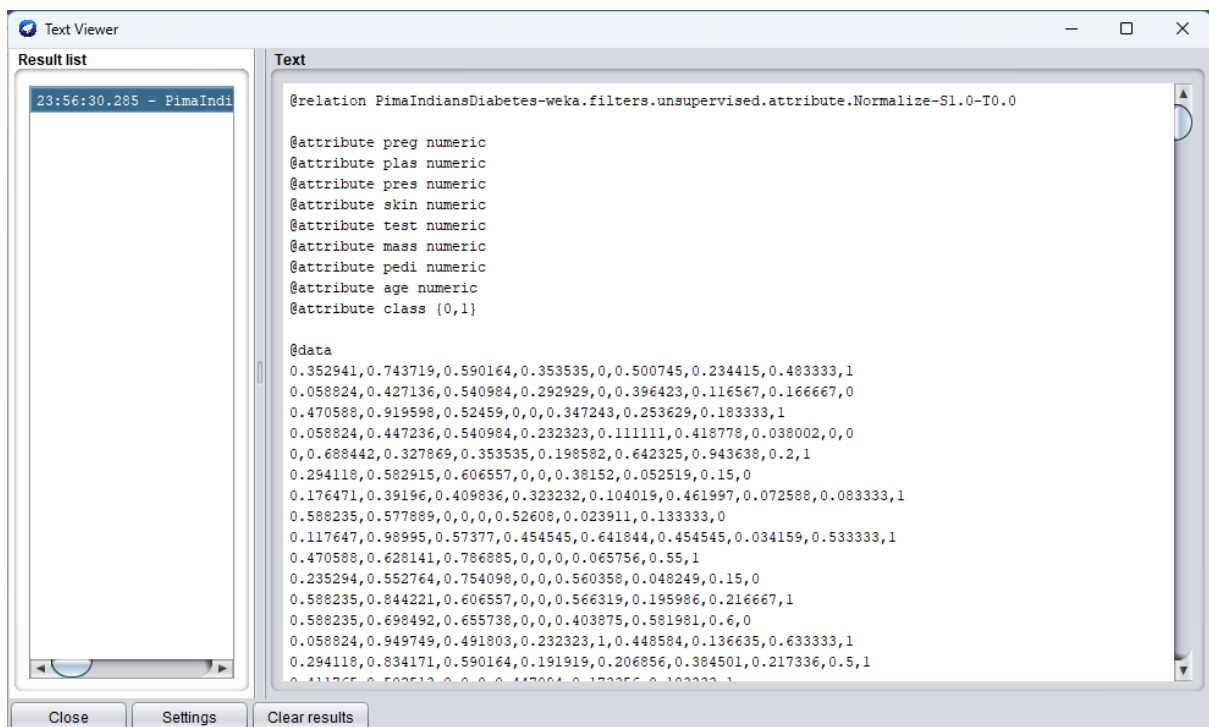
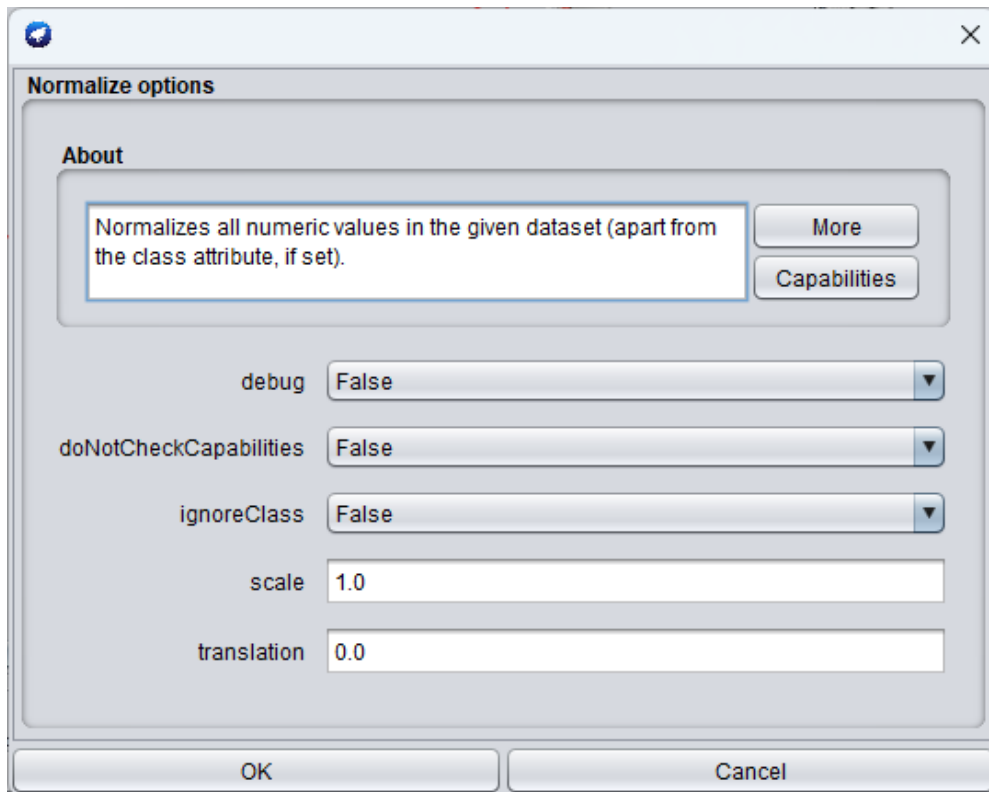
c. Select single /multi numeric attributes and convert it into numerical range between (Min-Max) ranges by formula given below: Note: normalization on other such as z-score, linear etc.



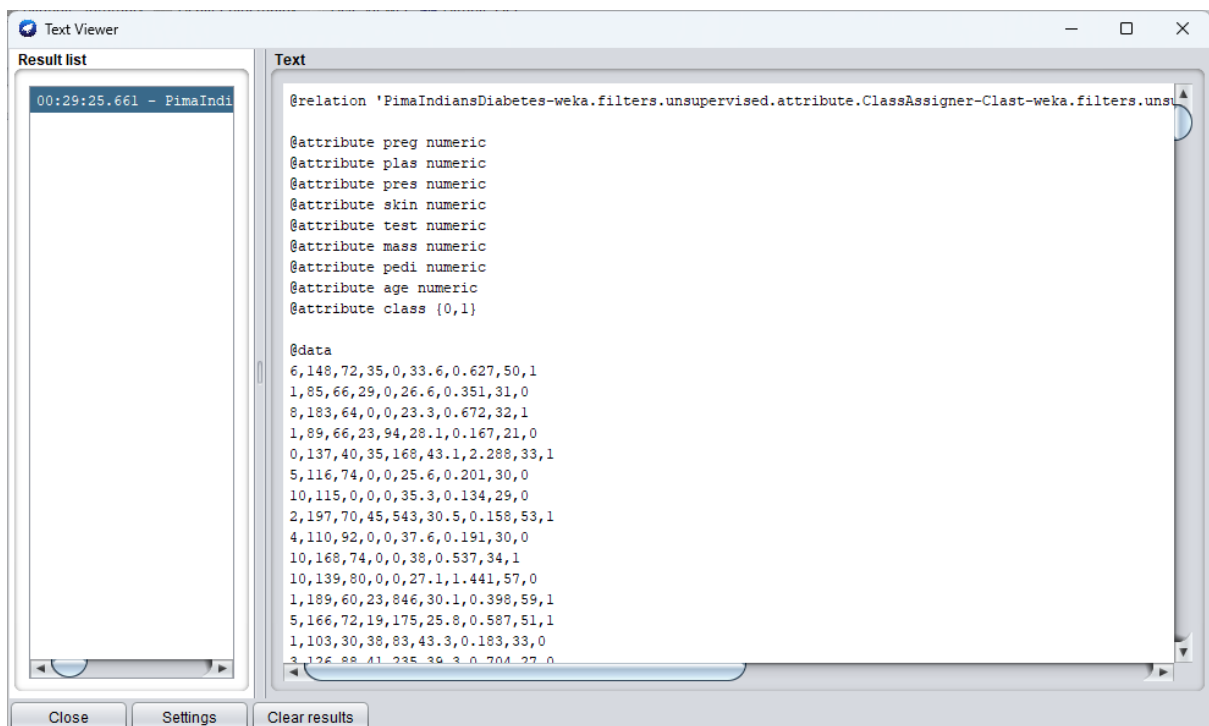
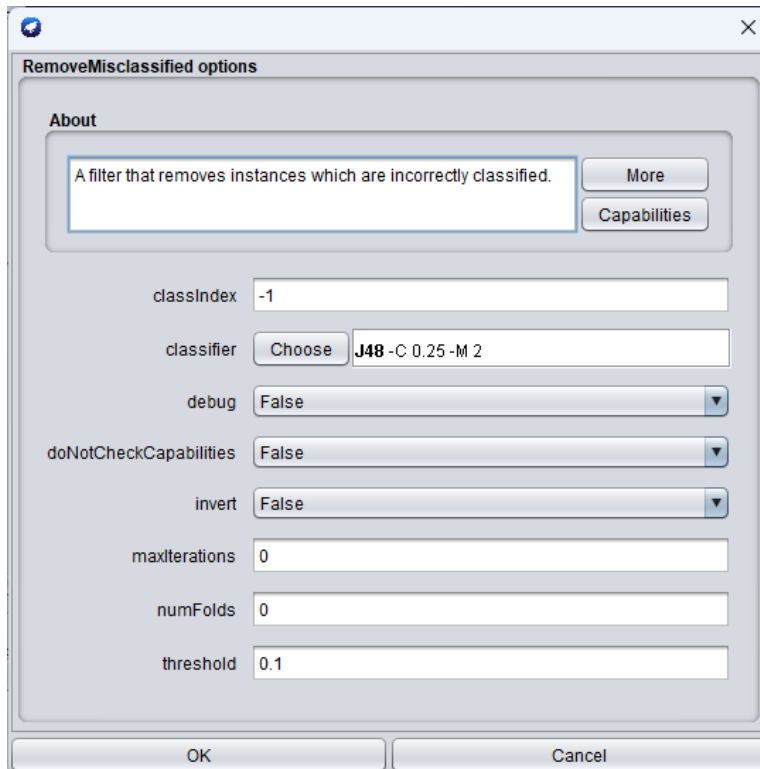
Q2. Download (<https://www.kaggle.com/datasets/kumargh/pimaIndiansdiabetescsv> Pima Indian Diabetes Dataset (PIDD) and do following operations and save all reports in separate file using data sink .csv format.



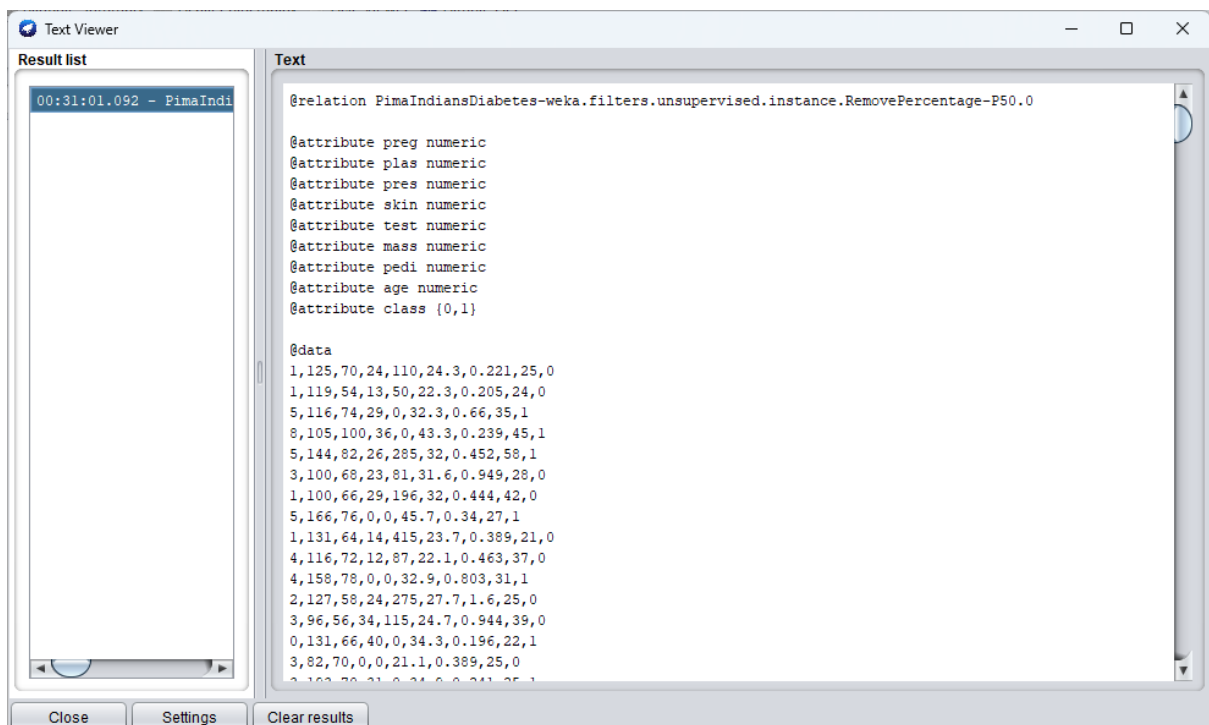
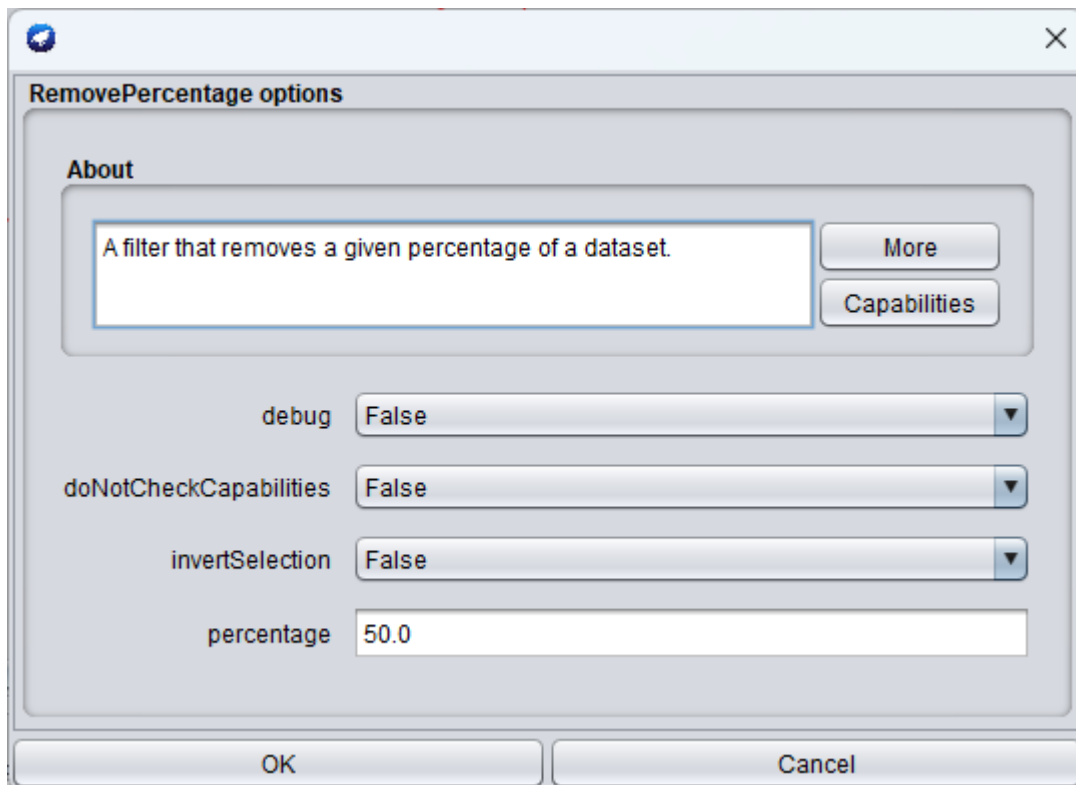
a. Normalize all numeric attributes



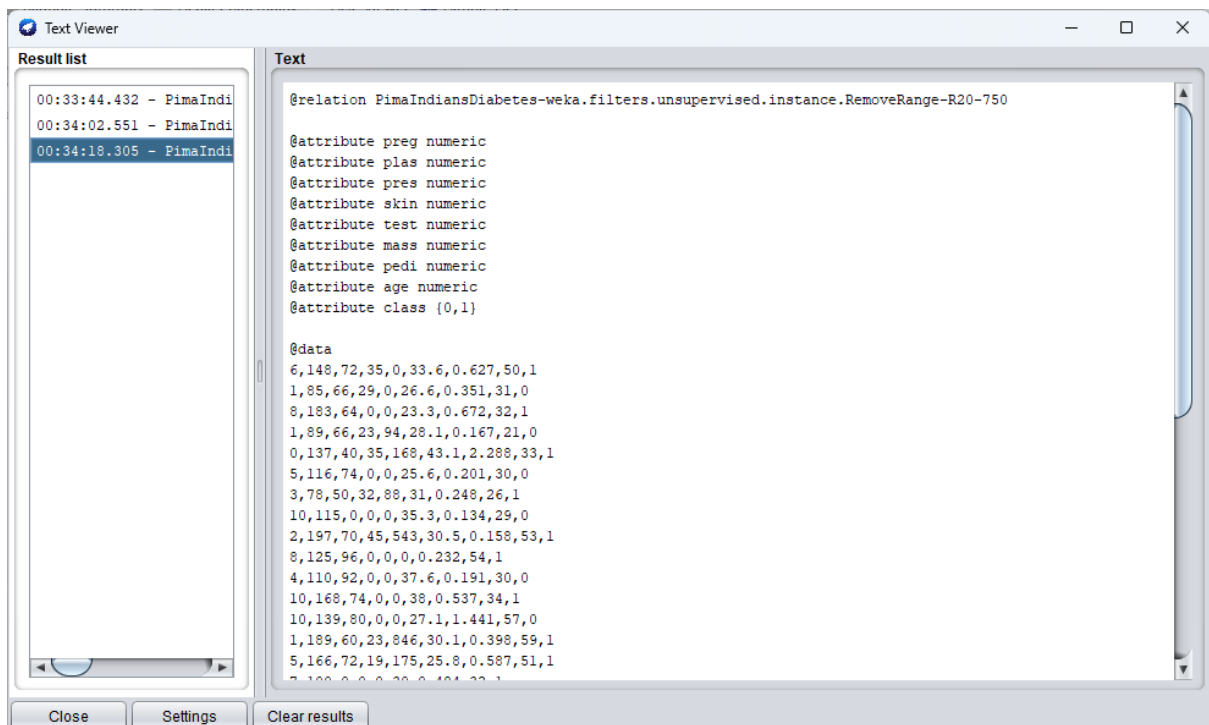
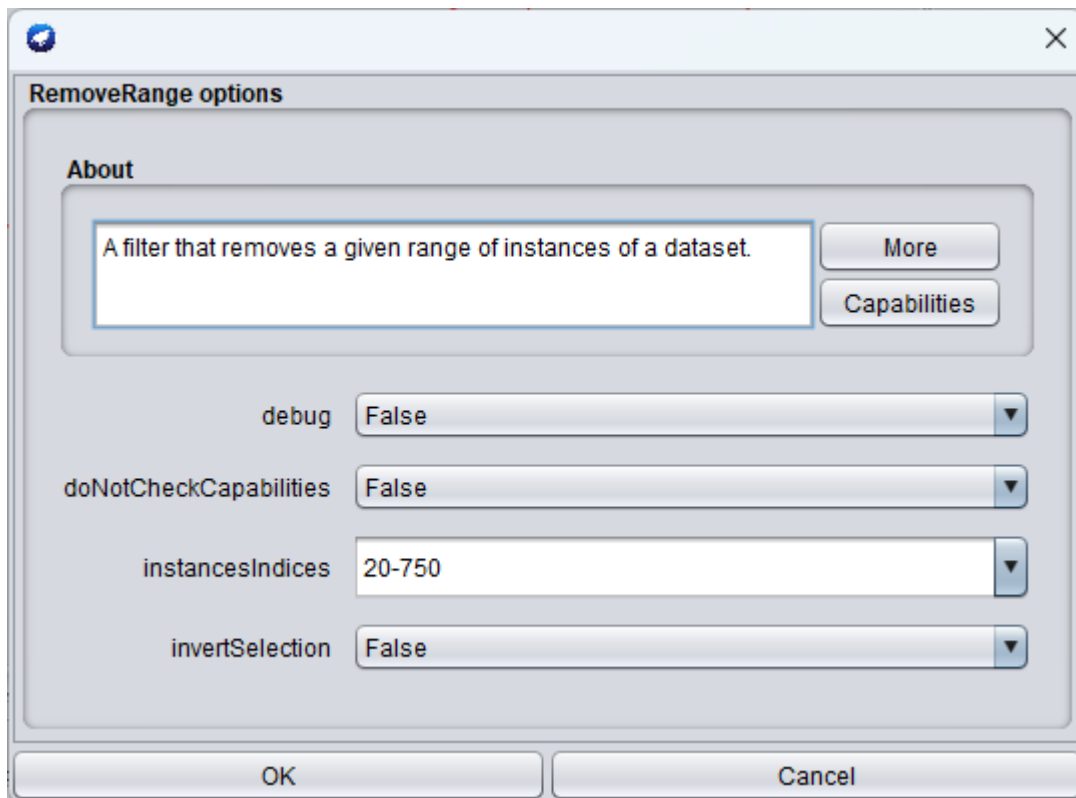
b. Remove all misclassified value by J48 class filter



c. Remove Percent



d. Remove Range



Resample options

About

Produces a random subsample of a dataset using either sampling with replacement or without replacement.

[More](#) [Capabilities](#)

biasToUniformClass 0.0

debug False

doNotCheckCapabilities False

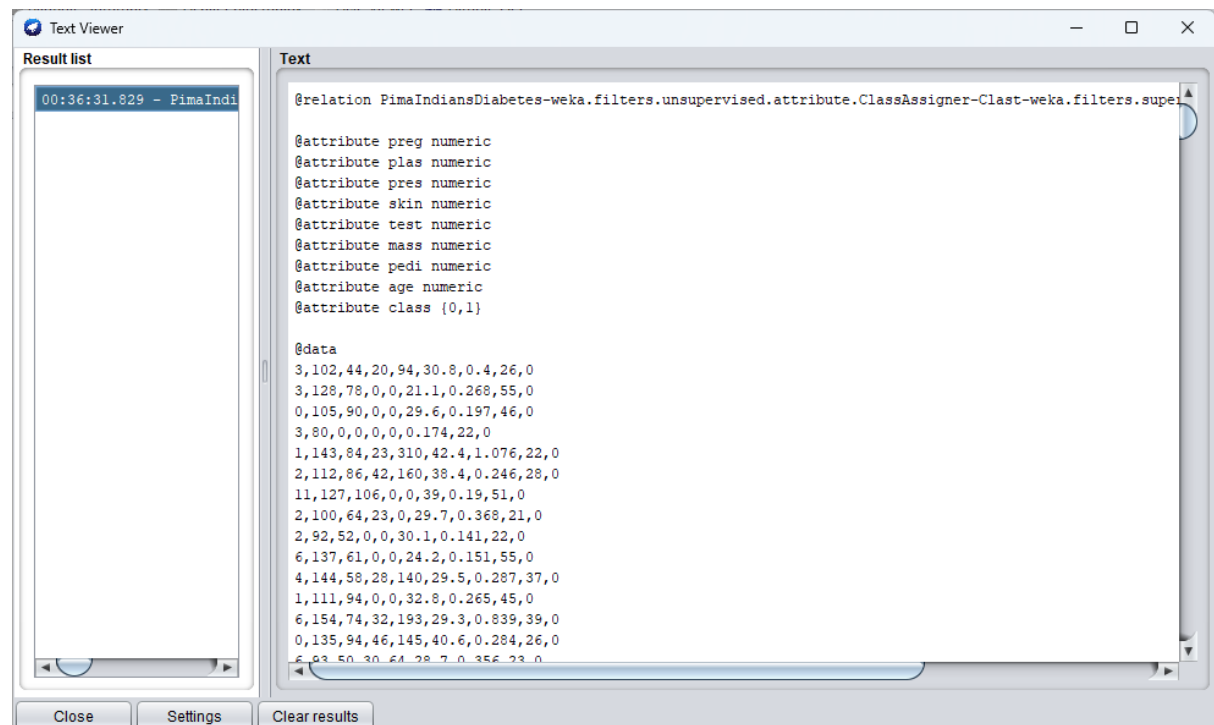
invertSelection False

noReplacement False

randomSeed 1

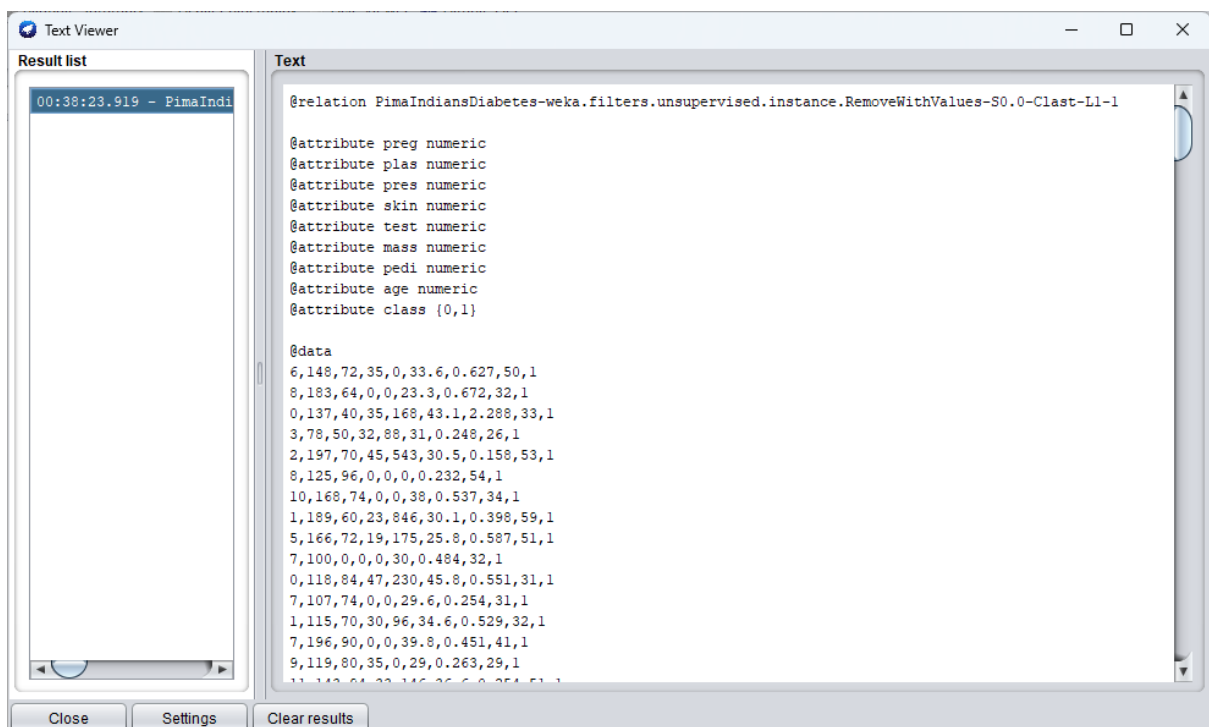
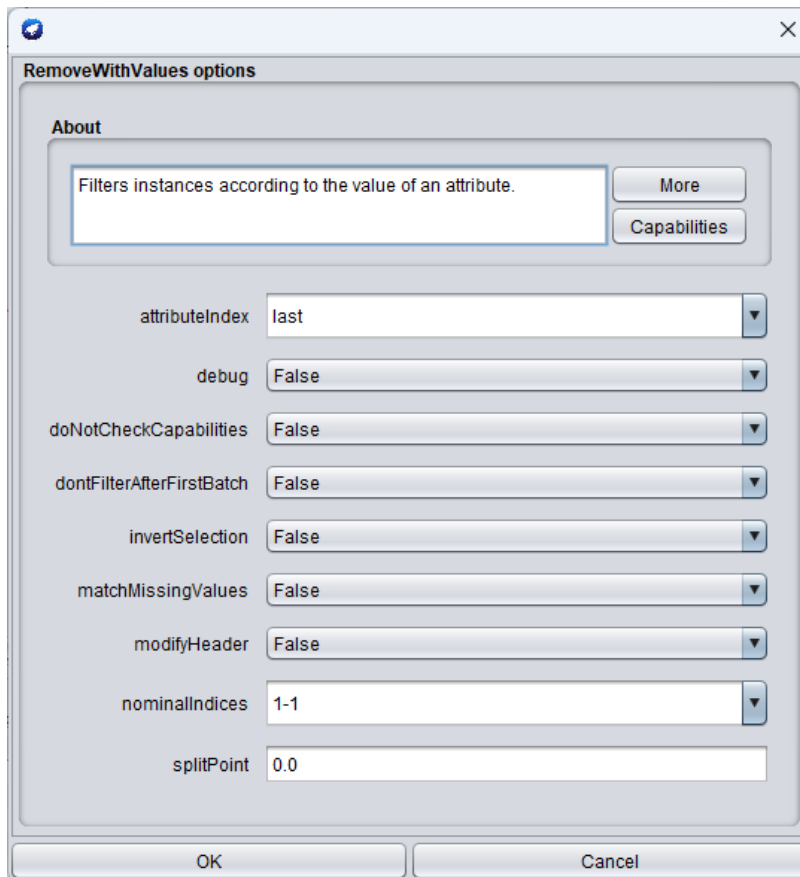
sampleSizePercent 100.0

OK Cancel



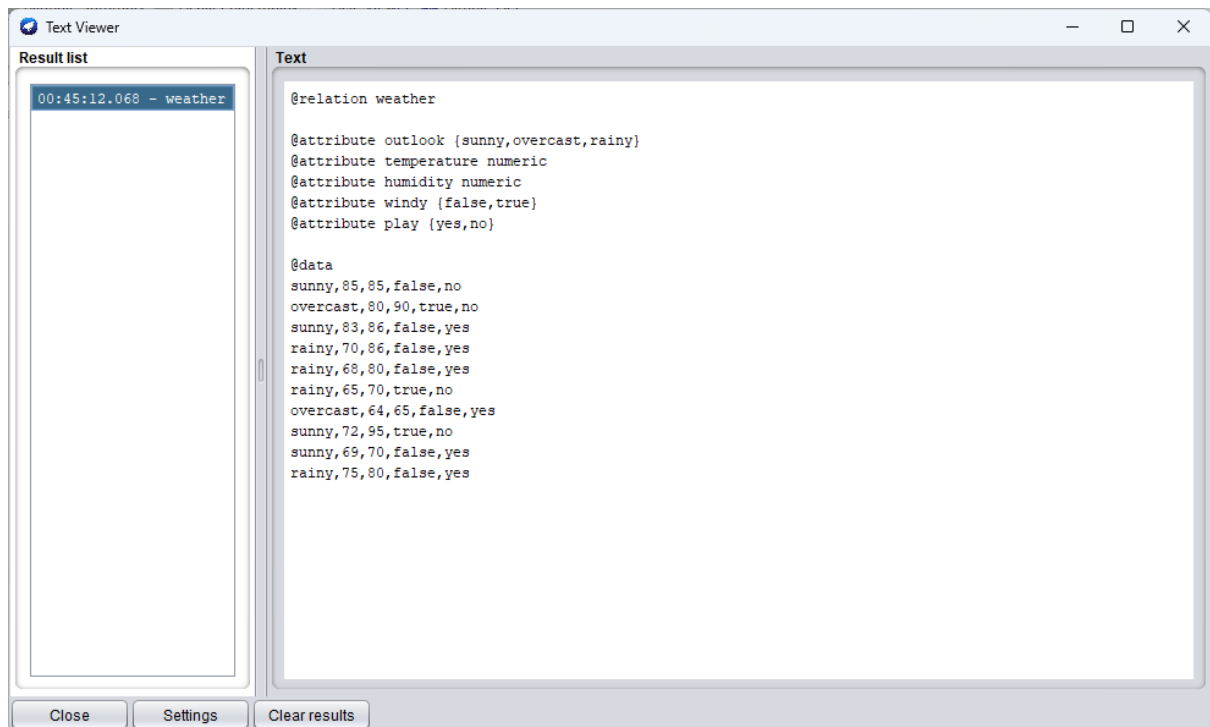
f. Remove With Values

instances with class=0 removed

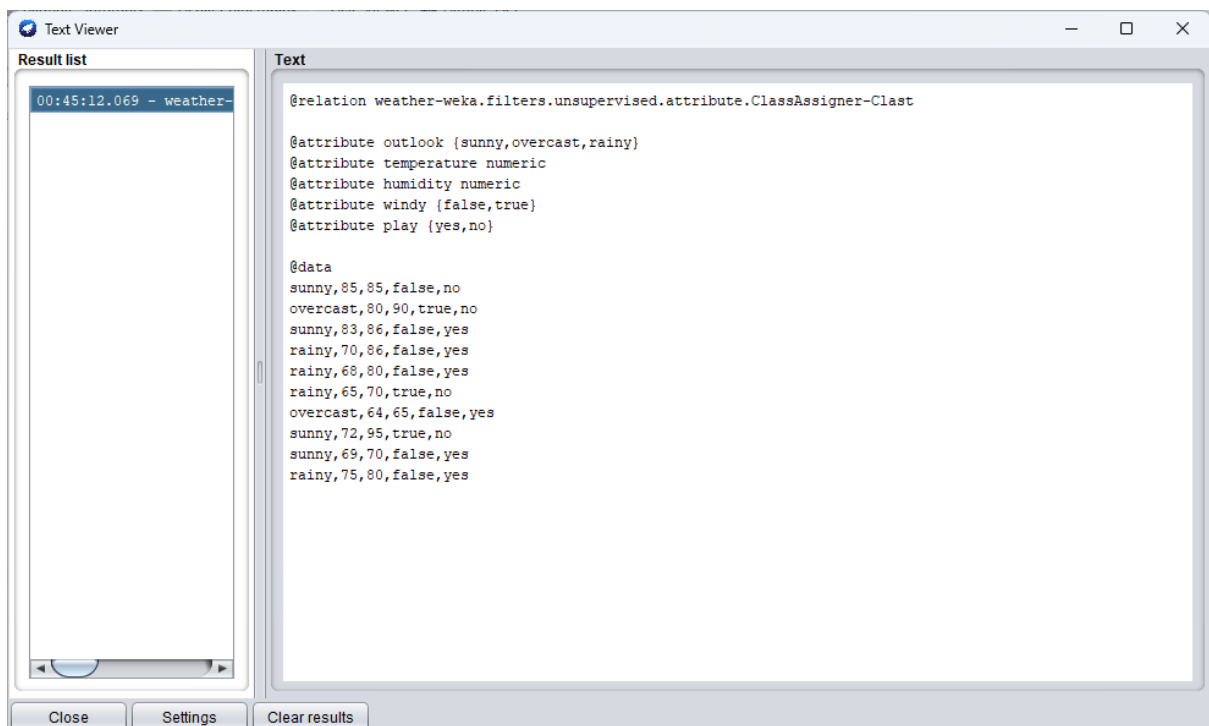
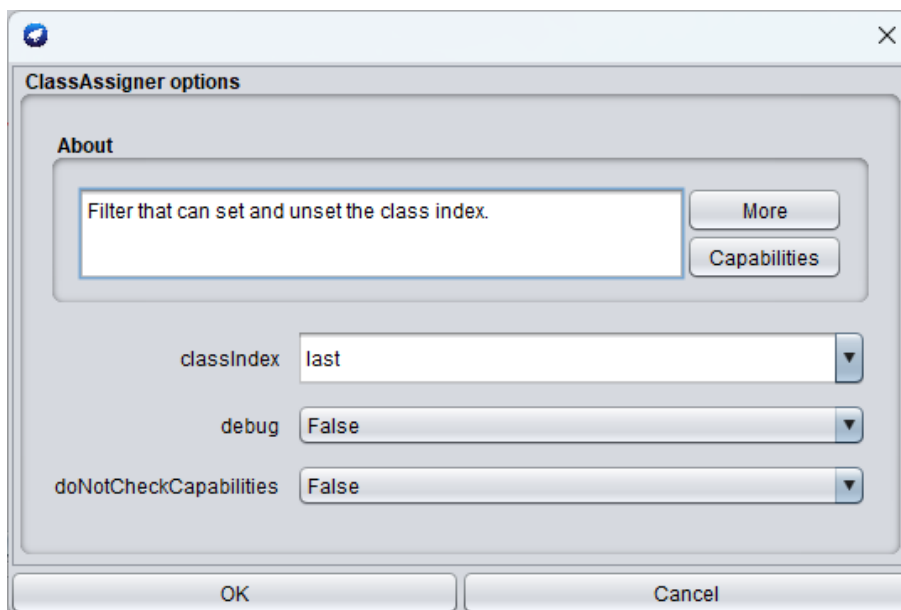
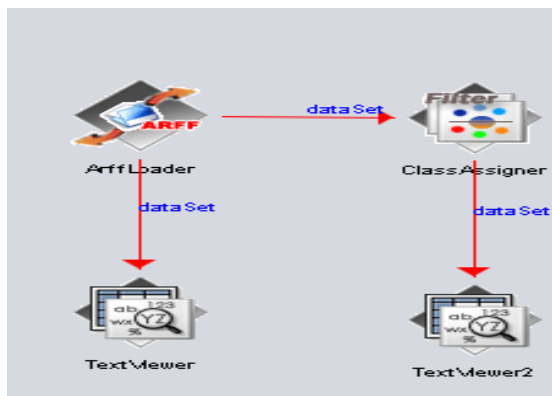


Q3.Training -test layout maker using Weka Knowledge Flow:

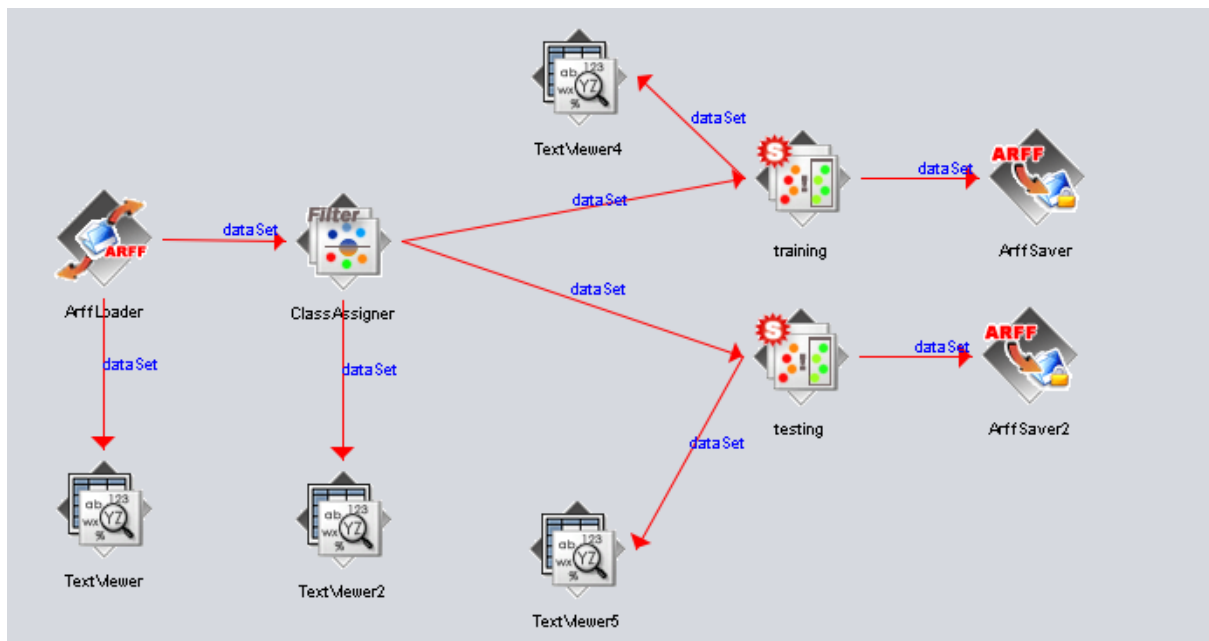
Load weather.arff dataset and do the following operation using:




a. Make play attribute as Class label.



b. Split dataset into training (80%) and testing(20%) part and save both splits in different folds in directory location



X

Resample options

About

Produces a random subsample of a dataset using either sampling with replacement or without replacement.

More

Capabilities

biasToUniformClass

0.0

debug

False

▼

doNotCheckCapabilities

False

▼

invertSelection

False

▼

noReplacement

False

▼

randomSeed

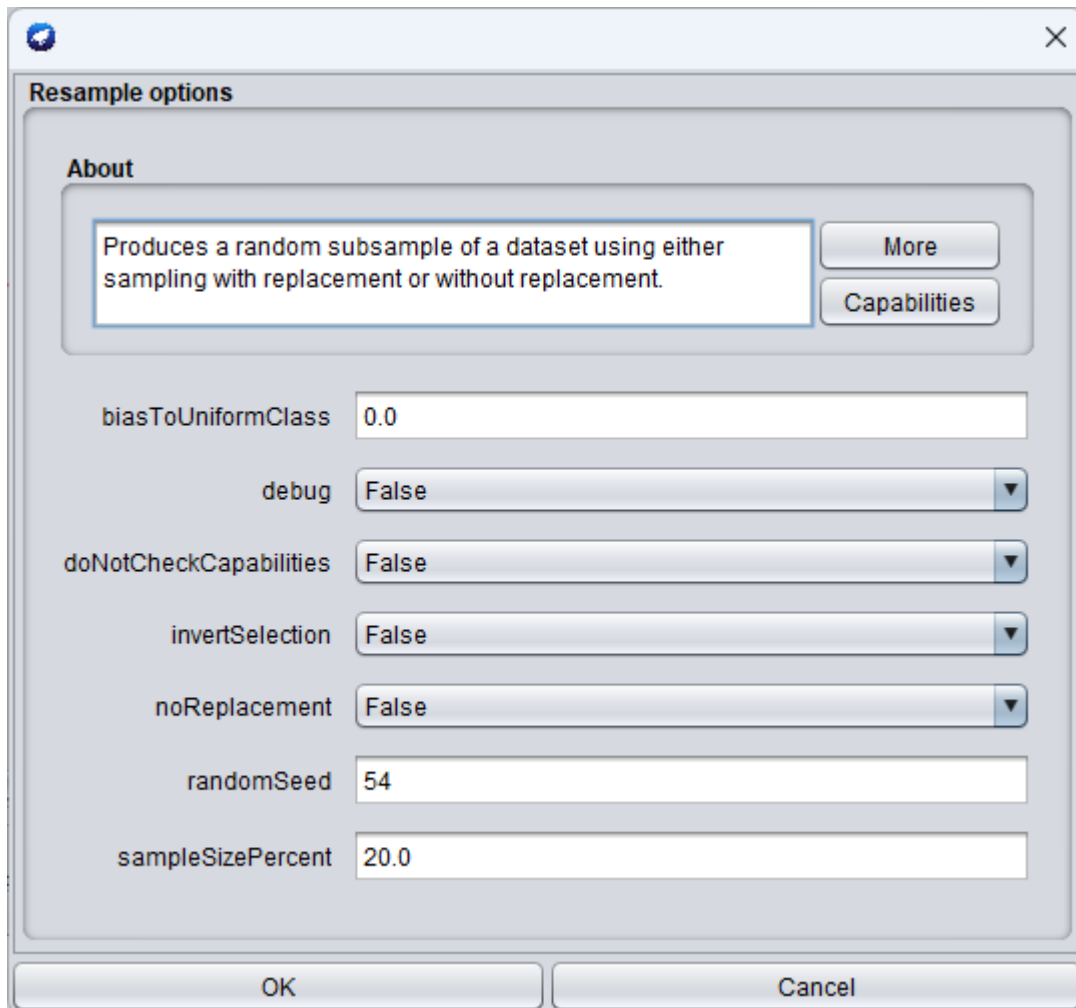
187

sampleSizePercent

80.0

OK

Cancel



The image shows a 'Resample options' dialog box with a title bar containing a blue icon and a close button. The dialog is divided into an 'About' section and a configuration section. The 'About' section contains a text box with the description 'Produces a random subsample of a dataset using either sampling with replacement or without replacement.' and two buttons: 'More' and 'Capabilities'. The configuration section contains several labeled input fields: 'biasToUniformClass' (text box with '0.0'), 'debug' (dropdown menu with 'False'), 'doNotCheckCapabilities' (dropdown menu with 'False'), 'invertSelection' (dropdown menu with 'False'), 'noReplacement' (dropdown menu with 'False'), 'randomSeed' (text box with '54'), and 'sampleSizePercent' (text box with '20.0'). At the bottom of the dialog are 'OK' and 'Cancel' buttons.

Resample options

About

Produces a random subsample of a dataset using either sampling with replacement or without replacement.

More

Capabilities

biasToUniformClass 0.0

debug False

doNotCheckCapabilities False

invertSelection False

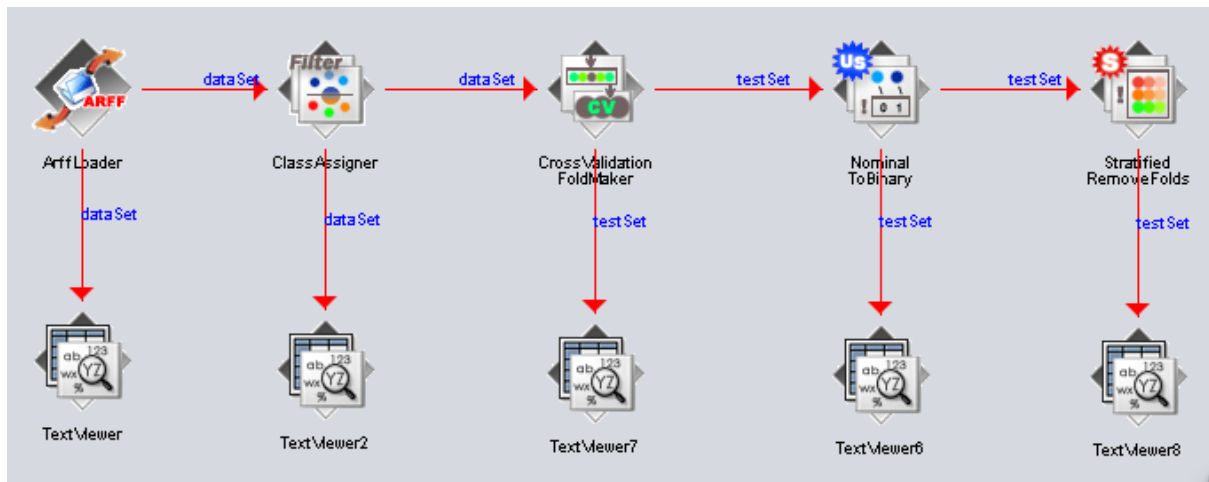
noReplacement False

randomSeed 54

sampleSizePercent 20.0

OK **Cancel**

c. Prepare 5 fold cross validation testing data and filter all such folds which having all given class label nominal values.



CrossValidationFoldMaker options

About

A Step that creates stratified cross-validation folds from incoming data [More](#)

Number of folds: 5

Preserve instances order: False

Random seed: 59

OK Cancel

StratifiedRemoveFolds options

About

This filter takes a dataset and outputs a specified fold for cross validation.

[More](#)

[Capabilities](#)

debug

doNotCheckCapabilities

fold

invertSelection

numFolds

seed

Text Viewer

Result list

- 01:59:21.532 - testSet:
- 01:59:21.532 - testSet:
- 01:59:21.533 - testSet:
- 01:59:21.533 - testSet:
- 01:59:21.533 - testSet:
- 01:59:21.533 - testSet:

Text

```
@relation weather-weka.filters.unsupervised.attribute.ClassAssigner-Clast-weka.filters.unsupervised.attrik
@attribute outlook=sunny numeric
@attribute outlook=overcast numeric
@attribute outlook=rainy numeric
@attribute temperature numeric
@attribute humidity numeric
@attribute windy=FALSE numeric
@attribute play {yes,no}

@data
0,1,0,80,75,1,yes
0,0,1,71,91,0,no
```