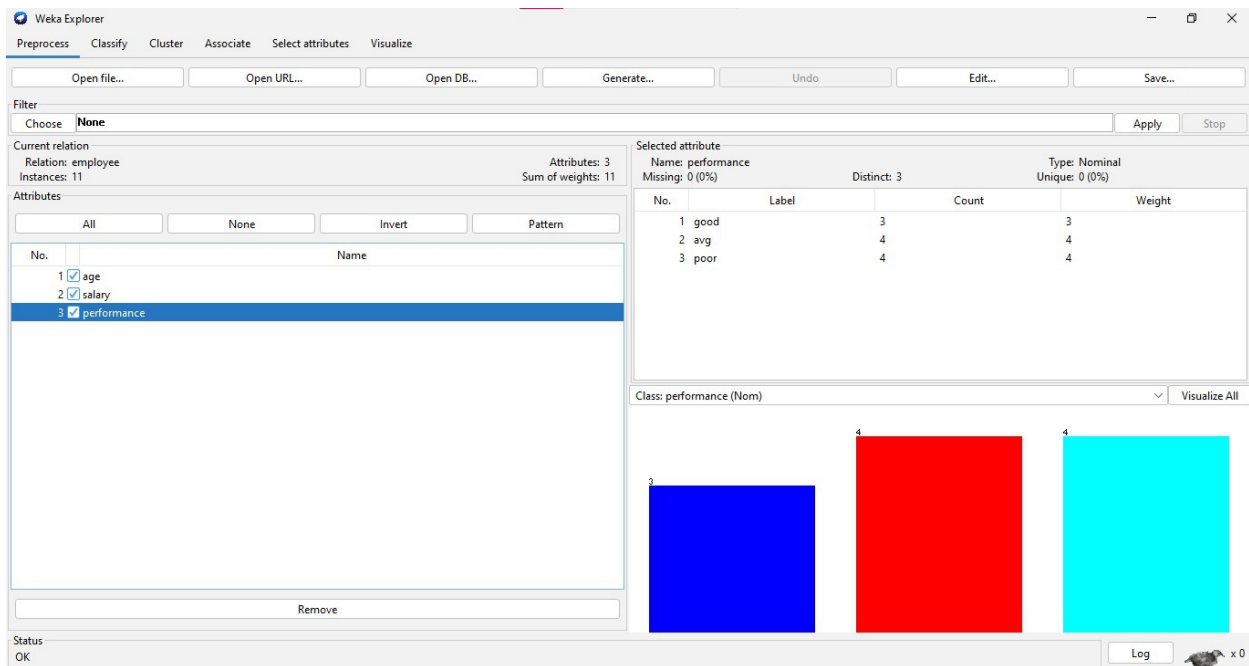# Data Mining file

**Pawan Bhati**

**23/SCA/BSC.IT/005**

1. Demonstration of preprocessing on dataset employee.arff



2.Demonstration of preprocessing on dataset employee.arff

Preprocess  Classify  Cluster  **Associate**  Select attributes  Visualize

**Associator**

Choose | **Apriori** -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c -1

Start | Stop

Associator output

Result list (right-click for ...)

09:10:20 - Apriori

```
=== Run information ===

Scheme:       weka.associations.Apriori -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c -1
Relation:     employee
Instances:    11
Attributes:   3
              age
              salary
              performance
=== Associator model (full training set) ===


Apriori
=======

Minimum support: 0.1 (1 instances)
Minimum metric <confidence>: 0.9
Number of cycles performed: 18

Generated sets of large itemsets:

Size of set of large itemsets L(1): 17

Size of set of large itemsets L(2): 25

Size of set of large itemsets L(3): 11

Best rules found:

 1. age=27 2 ==> performance=poor 2    <conf:(1)> lift:(2.75) lev:(0.12) [1] conv:(1.27)
 2. age=29 2 ==> performance=avg 2     <conf:(1)> lift:(2.75) lev:(0.12) [1] conv:(1.27)
 3. age=30 2 ==> performance=avg 2     <conf:(1)> lift:(2.75) lev:(0.12) [1] conv:(1.27)
```

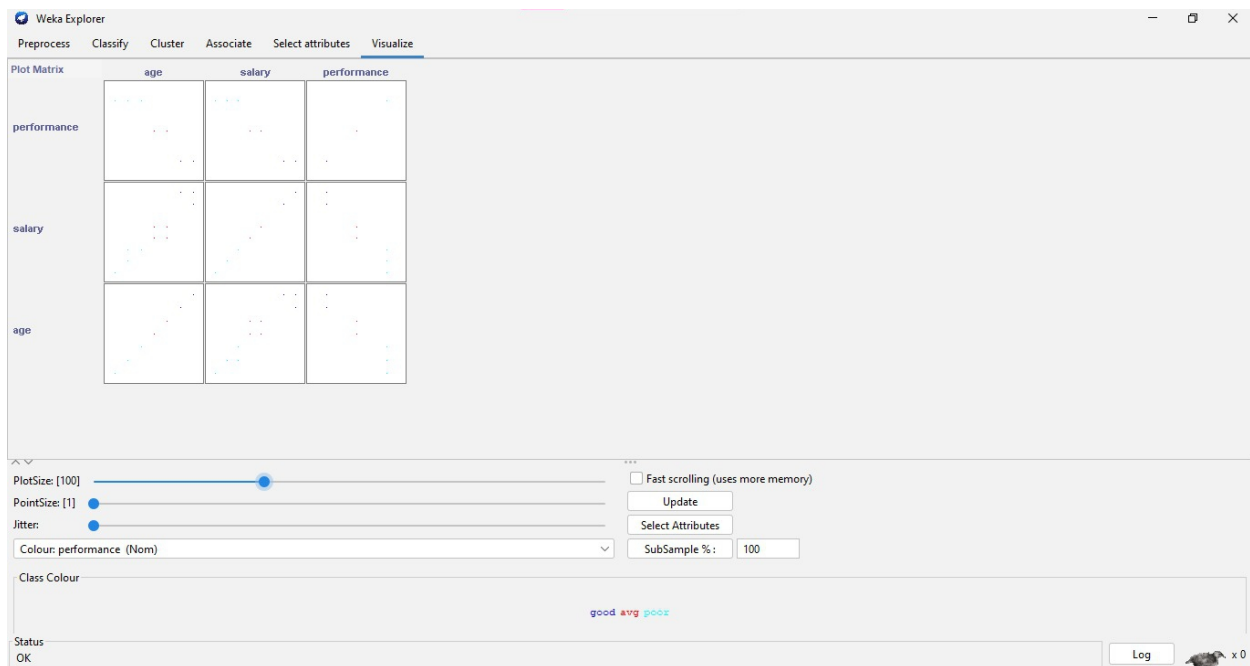Status
OK

Log | x 0

---

```
 4. age=48 2 ==> performance=good 2    <conf:(1)> lift:(3.67) lev:(0.13) [1] conv:(1.45)
 5. salary=17k 2 ==> performance=poor 2    <conf:(1)> lift:(2.75) lev:(0.12) [1] conv:(1.27)
 6. salary=20k 2 ==> performance=avg 2     <conf:(1)> lift:(2.75) lev:(0.12) [1] conv:(1.27)
 7. salary=25k 2 ==> performance=avg 2     <conf:(1)> lift:(2.75) lev:(0.12) [1] conv:(1.27)
 8. salary=32k 2 ==> performance=good 2    <conf:(1)> lift:(3.67) lev:(0.13) [1] conv:(1.45)
 9. salary=10k 1 ==> age=25 1    <conf:(1)> lift:(11) lev:(0.08) [0] conv:(0.91)
10. age=25 1 ==> salary=10k 1    <conf:(1)> lift:(11) lev:(0.08) [0] conv:(0.91)
```
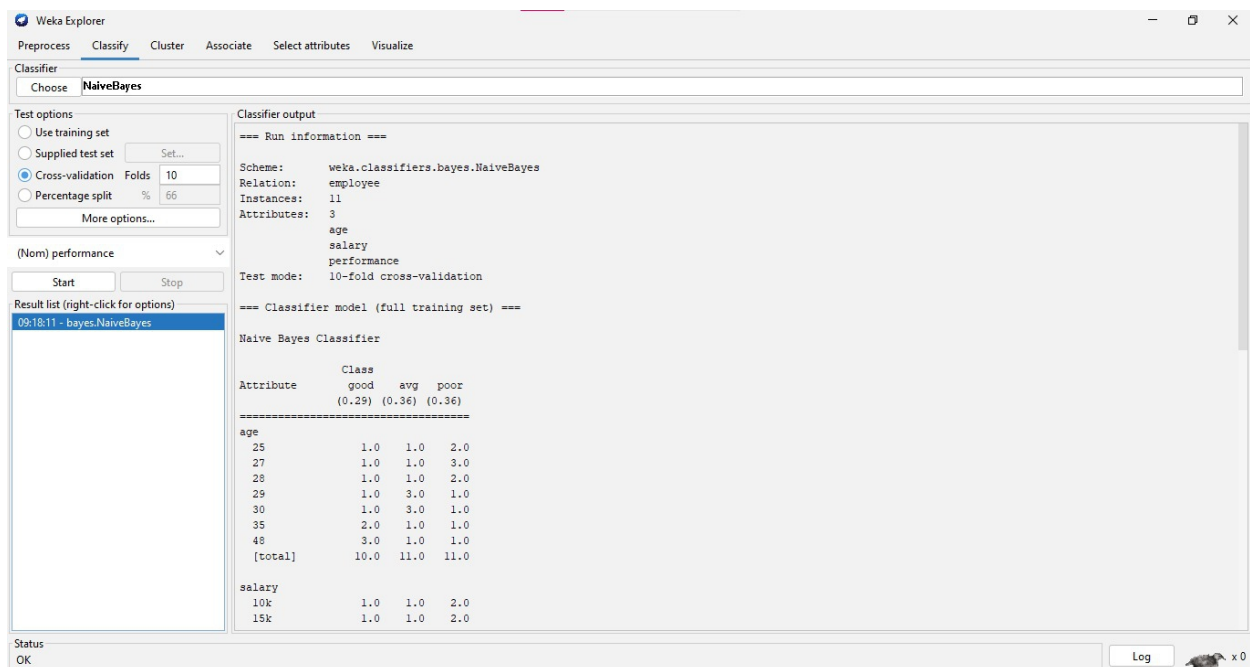
Status
OK

Log | x 0

3. Demonstration of classification rule
   process on dataset employee.arff using naïve Bayes algorithm

# Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

**Classifier**

Choose | NaiveBayes

**Test options**
- ○ Use training set
- ○ Supplied test set   Set...
- ● Cross-validation   Folds  10
- ○ Percentage split   %  66

More options...

(Nom) performance

Start | Stop

**Result list (right-click for options)**

09:18:11 - bayes.NaiveBayes

**Classifier output**

```
17k          1.0    1.0    3.0
20k          1.0    3.0    1.0
25k          1.0    3.0    1.0
30k          1.0    1.0    1.0
35k          2.0    1.0    1.0
32k          3.0    1.0    1.0
[total]     11.0   12.0   12.0


Time taken to build model: 0 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances          10               90.9091 %
Incorrectly Classified Instances         1                9.0909 %
Kappa statistic                          0.8625
Mean absolute error                      0.2899
Root mean squared error                  0.3171
Relative absolute error                 61.3111 %
Root relative squared error             63.0158 %
Total Number of Instances               11

=== Detailed Accuracy By Class ===

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
                 1.000    0.000    1.000      1.000   1.000      1.000  1.000     1.000     good
                 1.000    0.143    0.800      1.000   0.889      0.828  1.000     1.000     avg
                 0.750    0.000    1.000      0.750   0.857      0.810  1.000     1.000     poor
Weighted Avg.    0.909    0.052    0.927      0.909   0.908      0.868  1.000     1.000
```

=== Confusion Matrix ===

**Status**

OK | Log | x 0

---

```
=== Confusion Matrix ===

 a b c   <-- classified as
 3 0 0 | a = good
 0 4 0 | b = avg
 0 1 3 | c = poor
```

**Status**

OK | Log | x 0

---

# Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

**Plot Matrix**

|  | age | salary | performance |
|---|---|---|---|
| performance | | | |
| salary | | | |
| age | | | |

PlotSize: [100]
PointSize: [1]
Jitter:

☐ Fast scrolling (uses more memory)

Update

Select Attributes

Colour: performance (Nom) | SubSample % : 100

**Class Colour**

good avg poor

**Status**

OK | Log | x 0

## 4. Demonstration of clustering rule process on dataset Employee.arff using simple k- means



```
=== Run information ===

Scheme:        weka.clusterers.SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 2 -A "*
Relation:      employee
Instances:     11
Attributes:    3
               age
               salary
               performance
Test mode:     evaluate on training data


=== Clustering model (full training set) ===


kMeans
======

Number of iterations: 2
Within cluster sum of squared errors: 17.0

Initial starting points (random):

Cluster 0: 30,25k,avg
Cluster 1: 25,10k,poor

Missing values globally replaced with mean/mode

Final cluster centroids:
                               Cluster#
Attribute      Full Data          0           1
               (11.0)         (7.0)       (4.0)
```



```
Final cluster centroids:
                               Cluster#
Attribute      Full Data          0           1
               (11.0)         (7.0)       (4.0)
==================================================
age                27             29          27
salary            17k            20k         17k
performance       avg            avg        poor



Time taken to build model (full training data) : 0 seconds

=== Model and evaluation on training set ===

Clustered Instances

0       7 ( 64%)
1       4 ( 36%)
```