# Experiment 4

## Association Rule Mining on Employee Dataset

Aim: To create an employee.arff dataset and demonstrate Association rule process on it using apriori algorithm

Tasks:

1. Create employee.arff dataset and load it into Weka.
2. Apply Apriori algorithm with default parameters.
3. Change the parameters and observe the results

Task 1: Create employee.arff dataset and load it.

Create employee.arff with following categorical attributes and load it in to Weka.

| Attribute | States |
|---|---|
| Designation | i.   Manager<br>ii.  Developer<br>iii. Tester |
| Beneficiary | i.  Yes<br>ii. No |
| GPF | i.  Yes<br>ii. No |
| Salary | i.   Low<br>ii.  Medium<br>iii. High |
| CreditRating | i.   Poor<br>ii.  Fair<br>iii. Excellent |
| BankLoan | i.  Yes<br>ii. No |

Task 2: Apply Apriori algorithm with default parameters

Association Rule Mining is a process that finds features which occur together or features that are correlated. Popular applications are Market Basket Analysis and Cross Marketing.

Association rules are mined out after frequent itemsets in a big dataset which can be found using algorithms such as Apriori and FP Growth.

Frequent Itemset mining mines data using support and confidence measures.

$$support\ (A \Rightarrow B) = \frac{number\ of\ instances\ containing\ both\ A\ and\ B}{total\ number\ of\ instances}$$

$$confidence\ (A \Rightarrow B) = p(B/A) = \frac{number\ of\ instances\ containing\ both\ A\ and\ B}{number\ of\ instances\ containing\ A}$$

Apriori Rule Learner in Weka implements Apriori algorithm. It iteratively reduces the minimum support from its upperBound until (i) it finds the required number of rules or the minimum support reaches lowerBound.

Default values of the some important parameters:

lowerBoundMinSupport = 0.1 (10%)
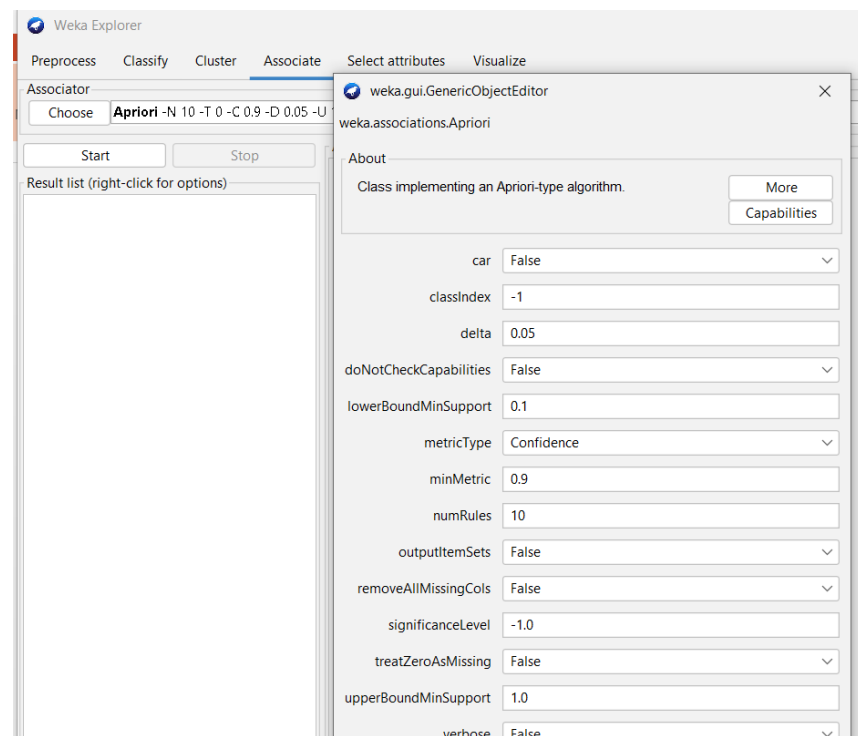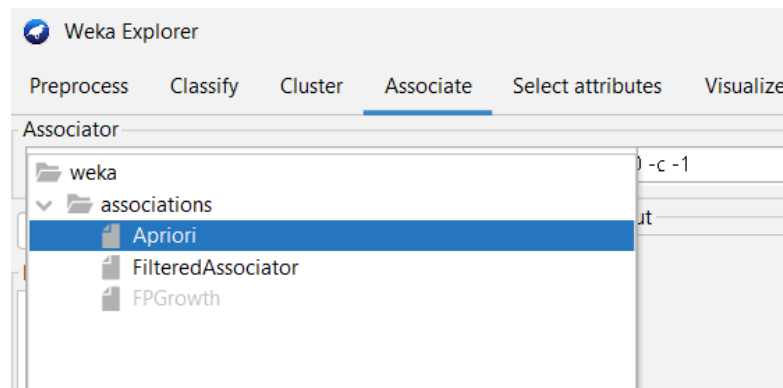
upperBoundMinSupport = 1.0 (100%)

metricType = Confidence

minMetric = 0.9 (90%)

numRules = 10

Steps:

    I.  Select Associate → Choose → associations → Apriori.

   II.  Observe the default parameter values.

  III.  Click on Start.

Observations:

| Default Parameters | Observations |
|---|---|
| lowerBoundMinSupport = | Minimum support = |
| upperBoundMinSupport = | Minimum Metric <Confidence> = |
| metricType = | Number of cycles performed = |
| minMetric = | Best rules found with confidence: |
| numRules = | 1. |
| | 2. |
| | 3. |
| | 4. |
| | 5. |
| | 6. |
| | 7. |
| | 8. |
| | 9. |
| | 10 |
| | . |

Task 3: Apply Apriori algorithm with required parameters

Observations: Change the default parameter values and perform the experiment

| Parameters | Observations |
|---|---|
| lowerBoundMinSupport = <br><br> upperBoundMinSupport = <br><br> metricType = <br><br> minMetric = <br><br> numRules = | Minimum support = <br><br> Minimum Metric <Confidence> = <br><br> Number of cycles performed = <br><br> Best rules found with confidence: |

Conclusion:

# Experiment 5

## Classification Using J48

Aim: To demonstrate Classification process on iris.arff dataset using j48 algorithm with percentage split.
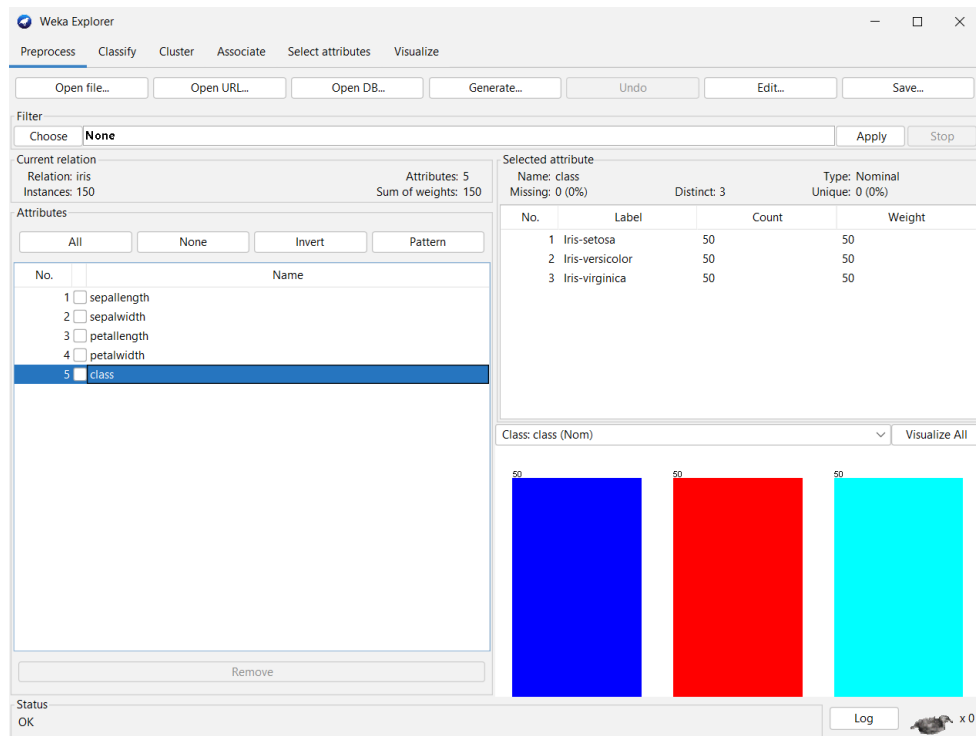
Tasks:

1. Load iris.arff dataset and explore it.
2. Build a classification model using J48 algorithm with percentage split.
3. Make predictions on new data.

Task 1: Load iris.arff dataset and explore it.

Load iris.arff from the Weka's data folder.

Observations:

| Attribute | Type | Range / States |
|-----------|------|----------------|
|           |      |                |
|           |      |                |
|           |      |                |
|           |      |                |
|           |      |                |
| Number of Instances: | | |

Task 2: Build a classification model using J48 algorithm with percentage split.

Classification is a process of determining the class (state) of the given instance. Examples:

Determining Play or Not play based on weather conditions.

Determining the digit (0 – 9) given the image pixel data.

Determining the Spam or Not-spam based on mail text.

J48 is Weka's Java implementation of the C4.5 algorithm. It can generate pruned or unpruned tress with both nominal and numerical attributes for classification.
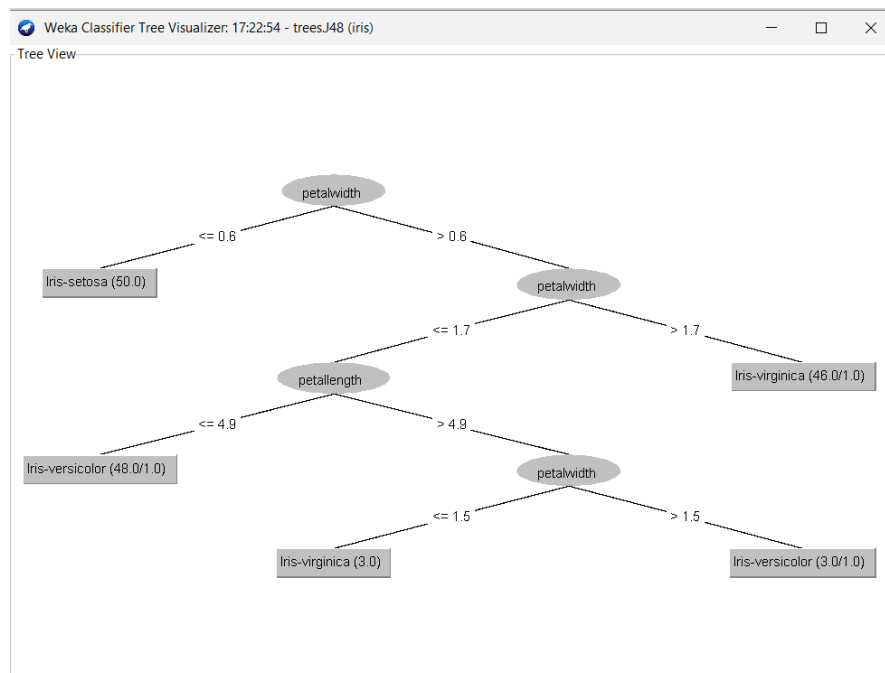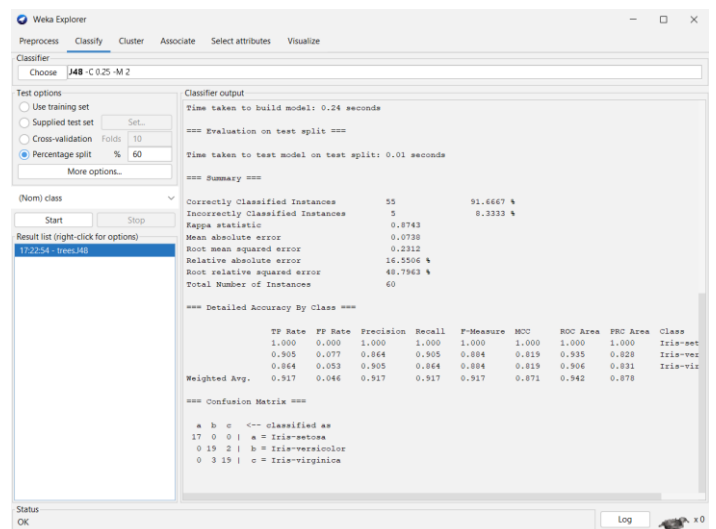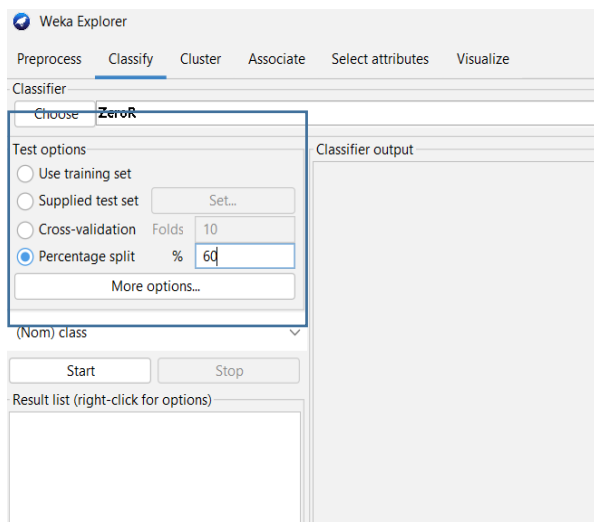
Percentage Split:

Splits the data into training and test subsets.

Training set is used to build the model.

Test set is used to evaluate the model performance.

Steps to build the model:

1. Click on Classify and select Percentage split with required training percentage under Test options group.

2. Select Choose → classifiers → trees → J48.

3. Clock on Start.

4. Right click on the model and click on Visualize tree

Observations:

- Total number of instances:
- Correctly classified instances:
- Incorrectly classified instances:
- Accuracy:
- Calculation of Accuracy from Confusion Matrix:

Task 3: Make predictions on new data

Steps:

i. Create an ARFF file with unlabeled (use ? in the place of class label) instances.

ii. On the "Classify" tab, select the "Supplied test set" option in the "Test options" pane.

iii. Click the "Set" button, click the "Open file" button on the options window and select the new dataset.

iv. Click the "More options…" button and for the "Output predictions" option click the "Choose" button and select "PlainText".

v. Right click on the model in the "Results list" pane and Select "Re-evaluate model on current test set".

Observations:

| Instance No. | Predicted class |
|---|---|
|  |  |

Conclusion:

# Experiment 6

## Classification Using J48

**Aim:** To demonstrate Classification process on StudentResult.arff dataset using j48 algorithm with cross-validation

**Tasks:**

1. Create StudentResult.arff dataset and load it into Weka.
2. Build a classification model using J48 algorithm with k-fold cross validation.
3. Make predictions on new data.

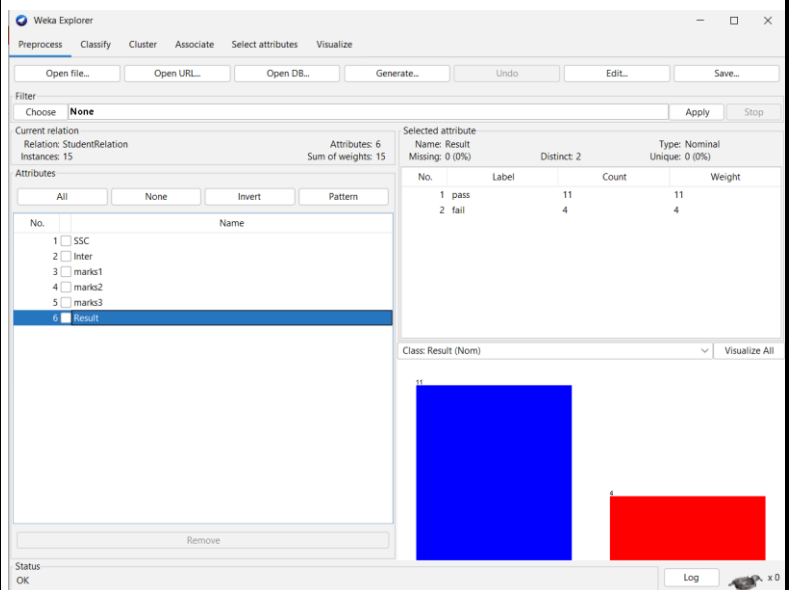Task 1: Create StudentResult.arff dataset and load it into Weka.

Create StudentResults.arff with following attributes and load it into Weka

| Attribute | Type |
|-----------|------|
| SSC | Nominal with states First, Second, Third. |
| Inter | Nominal with states First, Second, Third. |
| Marks1 | Numeric |
| Marks2 | Numeric |
| Marks3 | Numeric |
| Result | Nominal with states Pass & Fail. |

```
@relation StudentRelation
@attribute SSC {first,second,third}
@attribute Inter {first,second,third}
@attribute marks1 numeric
@attribute marks2 numeric
@attribute marks3 numeric
@attribute Result {pass, fail}
@data
first,second,55,57,62,pass
first,first,63,63,55,pass
first,first,65,67,66,pass
first,first,76,77,82,pass
second,third,32,43,23,fail
first,first,67,76,57,pass
second,second,56,54,45,pass
second,second,56,65,57,pass
third,third,34,23,12,fail
second,third,23,34,23,fail
second,first,65,64,56,pass
first,first,65,66,67,pass
third,third,45,32,23,fail
first,first,56,63,73,pass
second,first,65,56,57,pass
```

Task 2: Build a classification model using J48 algorithm with k-fold cross validation

Classification is a process of determining the class (state) of the given instance. Examples:

Determining Play or Not play based on weather conditions.

Determining the digit $(0 - 9)$ given the image pixel data.

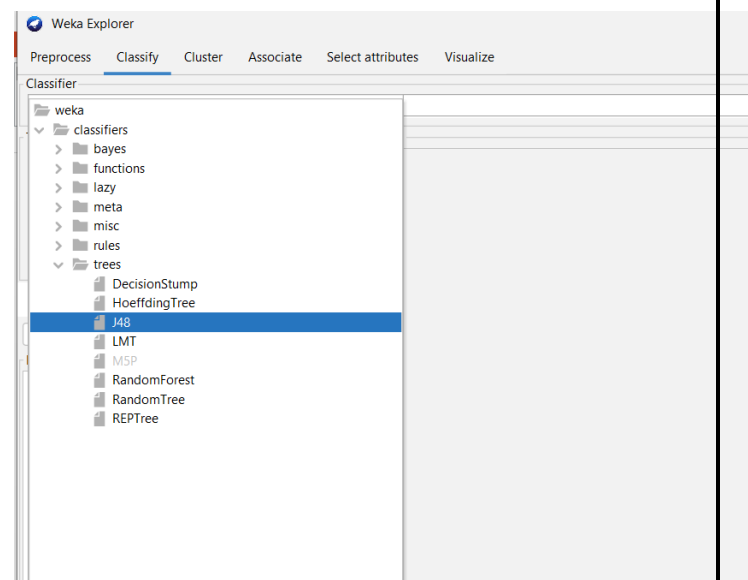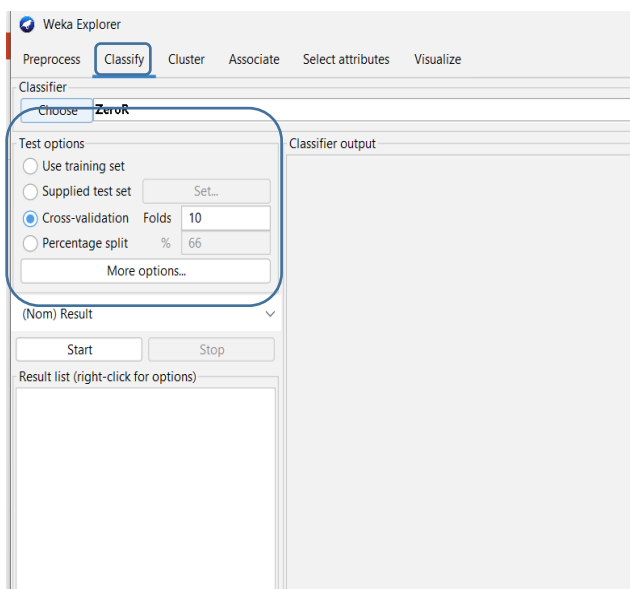Determining the Spam or Not-spam based on mail text.

J48 is Weka's Java implementation of the C4.5 algorithm. It can generate pruned or unpruned tress with both nominal and numerical attributes for classification.
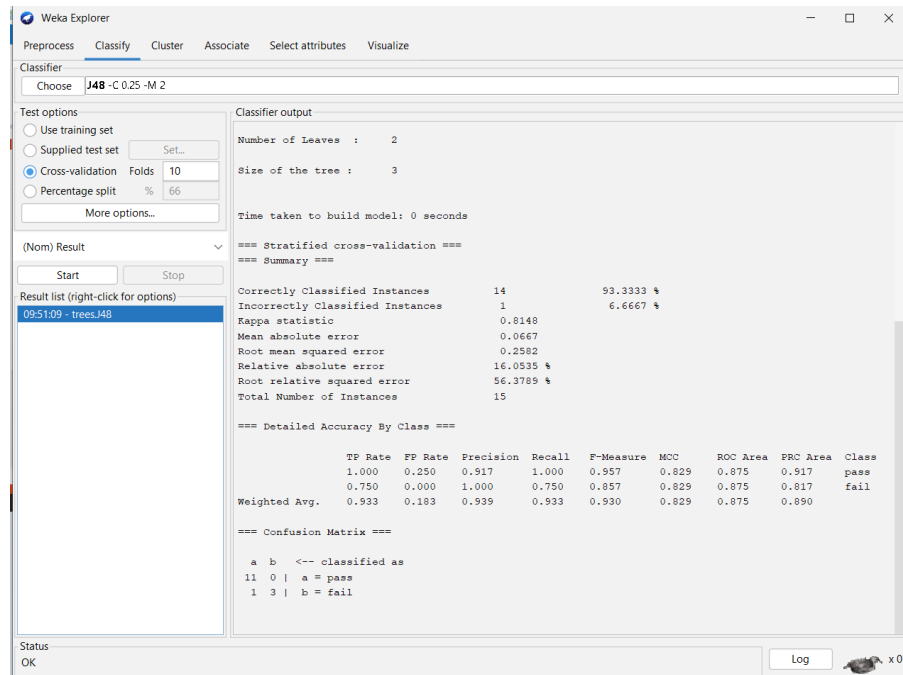
K-fold Cross Validation is a resampling procedure used to evaluate data mining models on a limited data set. It's process is

I.   Split the input dataset into K groups

II.  For i from 1 to k

- Take $i^{th}$ group as test dataset.

- Use remaining K-1 groups as training dataset.

- Fit the model using training set and evaluate its performance on test set.

Steps to build the model:

1. Click on Classify and select Cross-validation with default 10 folds under Test options group.

2. Select Choose → classifiers → trees → J48.

3. Clock on Start.

4. Right click on the model and click on Visualize tree

```
Weka Explorer                                                                    —   □   ×

Preprocess   Classify   Cluster   Associate   Select attributes   Visualize

Classifier
  Choose    J48 -C 0.25 -M 2

Test options                        Classifier output
○ Use training set
○ Supplied test set      Set...      Number of Leaves  :     2
● Cross-validation  Folds  10
○ Percentage split    %   66         Size of the tree :      3
      More options...
                                     Time taken to build model: 0 seconds
(Nom) Result
    Start         Stop               === Stratified cross-validation ===
                                     === Summary ===
Result list (right-click for options)
09:51:09 - trees.J48                 Correctly Classified Instances        14                93.3333 %
                                     Incorrectly Classified Instances       1                 6.6667 %
                                     Kappa statistic                        0.8148
                                     Mean absolute error                    0.0667
                                     Root mean squared error                0.2582
                                     Relative absolute error               16.0535 %
                                     Root relative squared error           56.3789 %
                                     Total Number of Instances             15

                                     === Detailed Accuracy By Class ===

                                                  TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
                                                   1.000    0.250    0.917     1.000    0.957     0.829   0.875     0.917    pass
                                                   0.750    0.000    1.000     0.750    0.857     0.829   0.875     0.817    fail
                                     Weighted Avg.  0.933    0.183    0.939     0.933    0.930     0.829   0.875     0.890

                                     === Confusion Matrix ===

                                      a  b   <-- classified as
                                     11  0 |  a = pass
                                      1  3 |  b = fail

Status
OK                                                                              Log
```

Tree generated by J48 algorithm for the given dataset is:

Observations:

- Total number of instances:

- Correctly classified instances:

- Incorrectly classified instances:

- Accuracy:

- Calculation of Accuracy from Confusion Matrix:
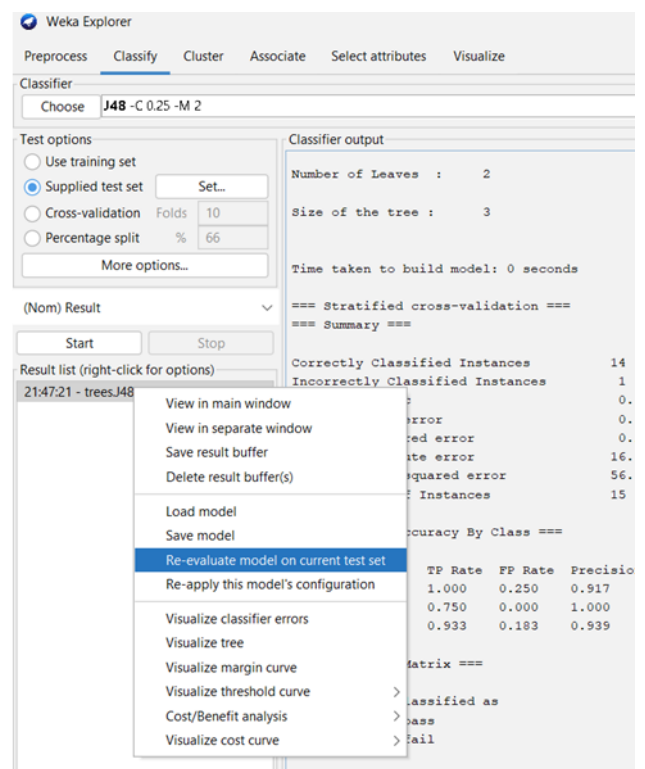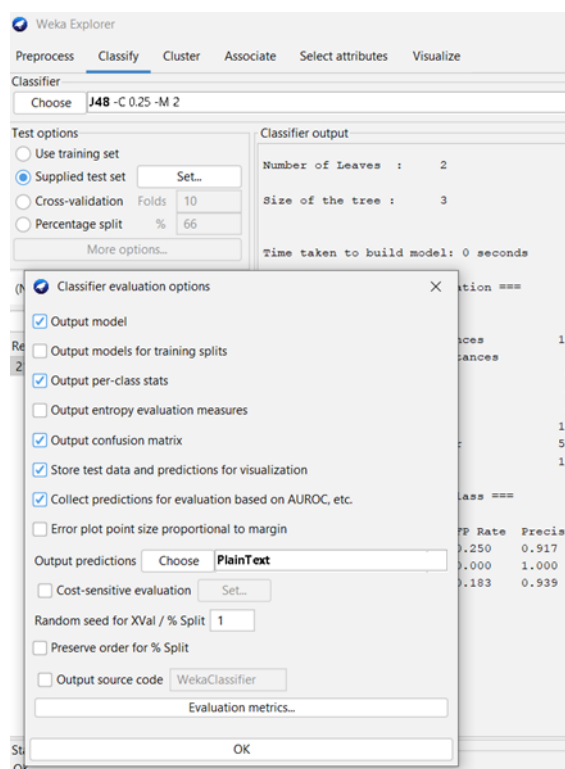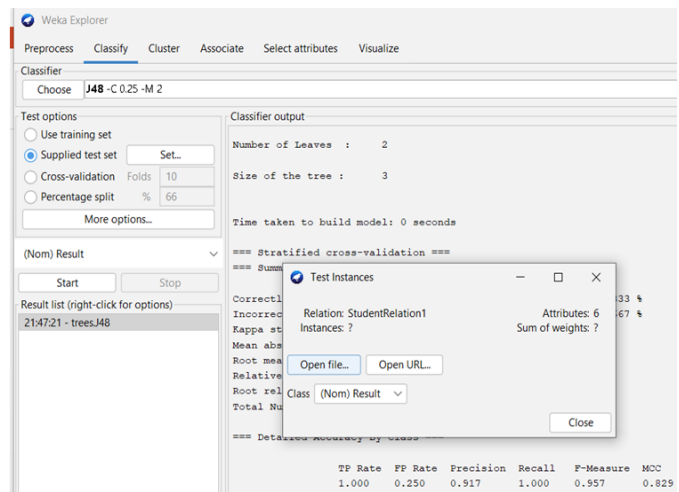
Task 3: Make predictions on new data

Steps:

i.    Create an ARFF file with unlabeled (use ? in the place of class label) instances.

ii.    On the "Classify" tab, select the "Supplied test set" option in the "Test options" pane.

iii.    Click the "Set" button, click the "Open file" button on the options window and select the new dataset.

iv.    Click the "More options…" button and for the "Output predictions" option click the "Choose" button and select "PlainText".

v.    Right click on the model in the "Results list" pane and Select "Re-evaluate model on current test set".

Weka Explorer

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

Classifier

Choose | J48 -C 0.25 -M 2

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

More options...

(Nom) Result

Start | Stop

Result list (right-click for options)

21:47:21 - trees.J48

Classifier output

```
=== Confusion Matrix ===

  a  b   <-- classified as
 11  0 |  a = pass
  1  3 |  b = fail


=== Re-evaluation on test set ===

User supplied test set
Relation:     StudentRelation1
Instances:    unknown (yet). Reading incrementally
Attributes:   6

=== Predictions on user test set ===

   inst#    actual  predicted error prediction
      1       1:?      1:pass       1
      2       1:?      2:fail       1
      3       1:?      1:pass       1


=== Summary ===

Total Number of Instances              0
Ignored Class Unknown Instances              3

=== Detailed Accuracy By Class ===

                 TP Rate FP Rate Precision Recall F-Measure MCC  ROC Area PRC Area Class
                 ?       ?       ?         ?      ?         ?    ?        ?        pass
                 ?       ?       ?         ?      ?         ?    ?        ?        fail
Weighted Avg.    ?       ?       ?         ?      ?         ?    ?        ?

=== Confusion Matrix ===
```

Status
OK

Log    x 0

Observations:

| Instance No. | Predicted class |
| --- | --- |
|  |  |

Conclusion:

# Experiment 7

## Classification Using ID3

**Aim:** To demonstrate Classification process on contact-lenses.arff dataset using ID3 algorithm with cross-validation.

**Tasks:**

1. Install required package for ID3
2. Load contact-lenses.arff dataset
3. Build a classification model using ID3 algorithm with k-fold cross-validation.
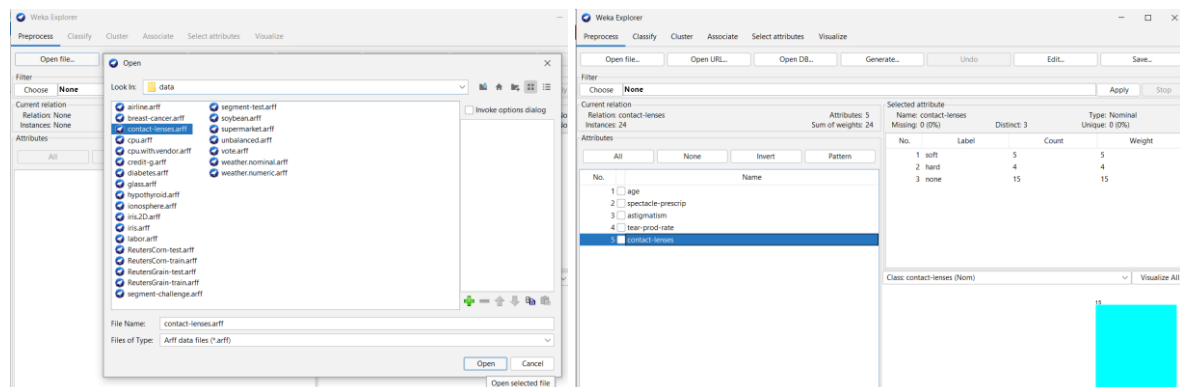4. Make predictions on new data.

**Task 1:** Install required package for ID3

Select GUI Chooser → Tools → Package Manager

Search, select & install *simpleEducationalLearningSchemes* package to get ID3 classifier under trees group.



**Task 2:** Load contact-lenses.arff dataset

Load contact-lenses.arff from the Weka's data folder.

Task 3: Build a classification model using ID3 algorithm with k-fold cross validation.

Classification is a process of determining the class (state) of the given instance.

ID3 stands for Iterative Dichotomiser 3 and is named such because the algorithm iteratively (repeatedly) dichotomizes(divides) features into two or more groups at each step. ID3 uses a top-down greedy approach to build a decision tree. ID3 is only used for classification problems with nominal features only.

K-fold Cross Validation is a resampling procedure used to evaluate data mining models on a limited data set. It's process is

I.   Split the input dataset into K groups

II.  For i from 1 to k

   • Take i<sup>th</sup> group as test dataset.

   • Use remaining K-1 groups as training dataset.

   • Fit the model using training set and evaluate its performance on test set.

Steps to build the model:

1. Click on Classify and select Cross-validation with some number folds under Test options group.

2. Select Choose → classifiers → trees → ID3.

3. Clock on Start.

   Note that we can't visualize ID3 tree.

Observations:

- Total number of instances:
- Correctly classified instances:
- Incorrectly classified instances:
- Accuracy:
- Calculation of Accuracy from Confusion Matrix:

Task 4: Make predictions on new data

Steps:

i. Create an ARFF file with unlabeled (use ? in the place of class label) instances.

ii. On the "Classify" tab, select the "Supplied test set" option in the "Test options" pane.

iii. Click the "Set" button, click the "Open file" button on the options window and select the new dataset.

iv. Click the "More options…" button and for the "Output predictions" option click the "Choose" button and select "PlainText".

v. Right click on the model in the "Results list" pane and Select "Re-evaluate model on current test set".

```
@relation contact-lenses1

@attribute age {young, pre-presbyopic, presbyopic}
@attribute spectacle-prescrip {myope, hypermetrope}
@attribute astigmatism {no, yes}
@attribute tear-prod-rate {reduced, normal}
@attribute contact-lenses {soft, hard, none}

@data
young,hypermetrope,no,normal,?
presbyopic,hypermetrope,no,normal,?
presbyopic,myope,no,normal,?
```

Observations:

| Instance No. | Predicted class |
| --- | --- |
|  |  |

Conclusion: