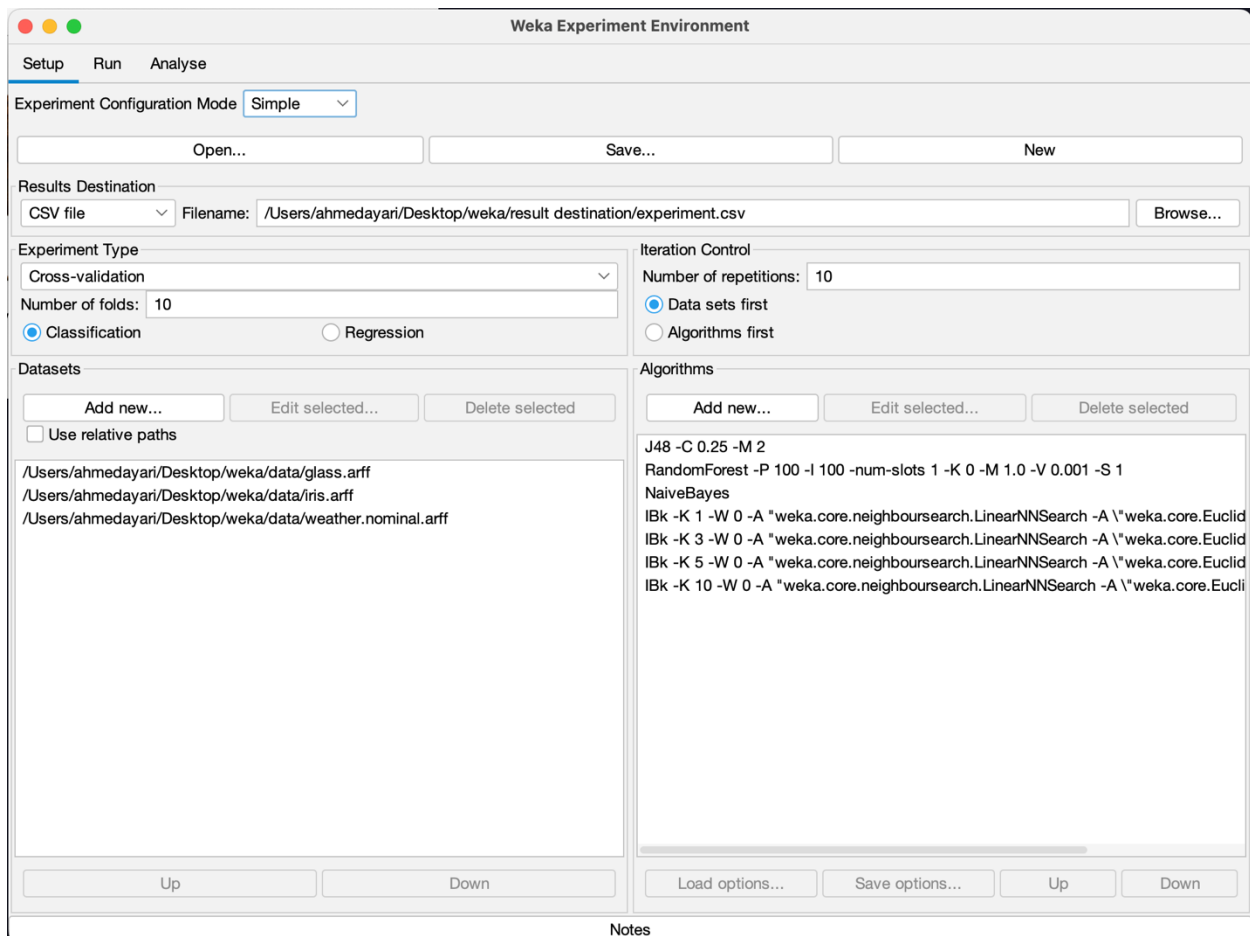


TP2 – Data Mining: Introduction to Weka Experimenter

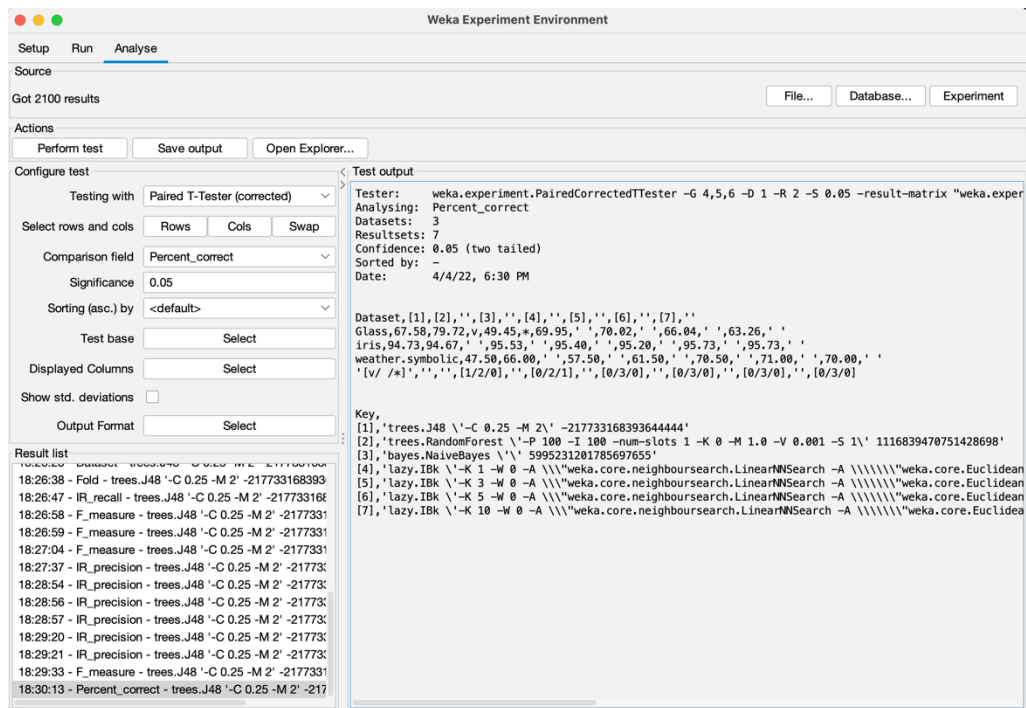
Realized by:
Ahmed Ayari,
Oussema Zouaghi,
Moahmed Karim Mallouli

Experience description:

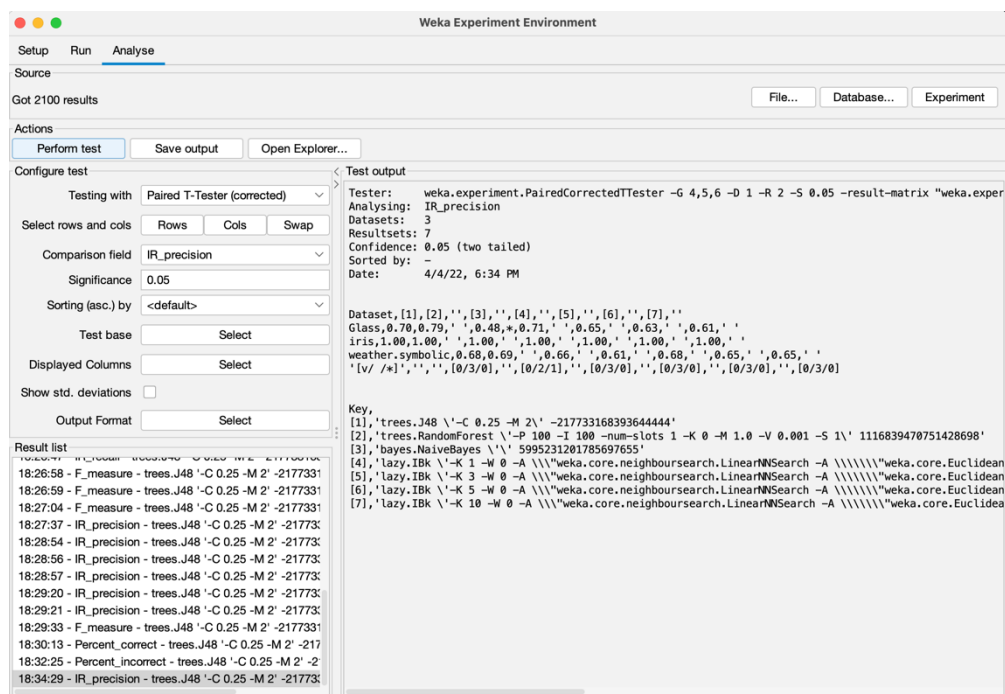
- Models: J48, RandomForest, Naive Bayes, IBK k=1, IBK k=3, IBK k=5, IBK k=10
- Datasets:
 - o Iris: A dataset containing observation of a lot of iris flowers and their races
 - o Glass: dataset containing observations about glass types and their chemical composition
 - o Weather: dataset containing observations of the weather status each corresponding to the fact that tennis (or other sport) went outside to practice or not
- Comparison attributes: Accuracy, Recall, Precision and f-value



- Percent correct (Accuracy):



- Precision:



- Recall

The screenshot shows the Weka Experiment Environment interface. The 'Analyse' tab is selected. The 'Source' section indicates 'Got 2100 results'. The 'Actions' section includes buttons for 'Perform test', 'Save output', and 'Open Explorer...'. The 'Configure test' section is set to 'Paired T-Tester (corrected)' with 'IR_recall' as the comparison field, a significance level of 0.05, and 'trees.J48' as the test base. The 'Test output' section displays the following information:

```
Tester: weka.experiment.PairedCorrectedTTester -G 4,5,6 -D 1 -R 2 -S 0.05 -result-matrix "weka.exper
Analysing: IR_recall
Datasets: 3
Resultsets: 7
Confidence: 0.05 (two tailed)
Sorted by: -
Date: 4/4/22, 6:35 PM
```

The 'Dataset' section shows the following data:

```
Dataset, [1], [2], [3], [4], [5], [6], [7], ''
Glass, 0.71, 0.86, v, 0.74, ' ', 0.75, ' ', 0.83, ' ', 0.83, ' ', 0.86, v
iris, 0.98, 1.00, ' ', 1.00, ' ', 1.00, ' ', 1.00, ' ', 1.00, ' ', 1.00, ' '
weather, symbolic, 0.47, 0.78, ' ', 0.71, ' ', 0.70, ' ', 0.91, v, 1.00, v, 1.00, v
[v/ /*], ' ', ' ', [1/2/0], ' ', [0/3/0], ' ', [0/3/0], ' ', [1/2/0], ' ', [1/2/0], ' ', [2/1/0]
```

The 'Key' section shows the following information:

```
Key,
[1], 'trees.J48 \'-C 0.25 -M 2\' -21773316839364444'
[2], 'trees.RandomForest \'-P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1\' 1116839470751428698'
[3], 'bayes.NaiveBayes \\' 5995231201785697655'
[4], 'lazy.IBk \'-K 1 -W 0 -A \\'weka.core.neighboursearch.LinearNNSearch -A \\'weka.core.Euclidean
[5], 'lazy.IBk \'-K 3 -W 0 -A \\'weka.core.neighboursearch.LinearNNSearch -A \\'weka.core.Euclidean
[6], 'lazy.IBk \'-K 5 -W 0 -A \\'weka.core.neighboursearch.LinearNNSearch -A \\'weka.core.Euclidean
[7], 'lazy.IBk \'-K 10 -W 0 -A \\'weka.core.neighboursearch.LinearNNSearch -A \\'weka.core.Euclidean
```

The 'Result list' section shows the following results:

```
18:26:59 - F_measure - trees.J48 \'-C 0.25 -M 2\' -2177331
18:27:04 - F_measure - trees.J48 \'-C 0.25 -M 2\' -2177331
18:27:37 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:28:54 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:28:56 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:28:57 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:29:20 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:29:21 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:29:33 - F_measure - trees.J48 \'-C 0.25 -M 2\' -2177331
18:30:13 - Percent_correct - trees.J48 \'-C 0.25 -M 2\' -217
18:32:25 - Percent_incorrect - trees.J48 \'-C 0.25 -M 2\' -2
18:34:29 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:35:05 - IR_recall - trees.J48 \'-C 0.25 -M 2\' -217733168
```

- F-value (calculated from the precision and recall of the test)

The screenshot shows the Weka Experiment Environment interface. The 'Analyse' tab is selected. The 'Source' section indicates 'Got 2100 results'. The 'Actions' section includes buttons for 'Perform test', 'Save output', and 'Open Explorer...'. The 'Configure test' section is set to 'Paired T-Tester (corrected)' with 'F_measure' as the comparison field, a significance level of 0.05, and 'trees.J48' as the test base. The 'Test output' section displays the following information:

```
Tester: weka.experiment.PairedCorrectedTTester -G 4,5,6 -D 1 -R 2 -S 0.05 -result-matrix "weka.exper
Analysing: F_measure
Datasets: 3
Resultsets: 7
Confidence: 0.05 (two tailed)
Sorted by: -
Date: 4/4/22, 6:37 PM
```

The 'Dataset' section shows the following data:

```
Dataset, [1], [2], [3], [4], [5], [6], [7], ''
Glass, 0.70, 0.82, v, 0.58, *, 0.72, ' ', 0.72, ' ', 0.71, ' ', 0.71, ' '
iris, 0.99, 1.00, ' ', 1.00, ' ', 1.00, ' ', 1.00, ' ', 1.00, ' ', 1.00, ' '
weather, symbolic, 0.80, 0.83, ' ', 0.80, ' ', 0.76, ' ', 0.84, ' ', 0.82, ' ', 0.82, ' '
[v/ /*], ' ', ' ', [1/2/0], ' ', [0/2/1], ' ', [0/3/0], ' ', [0/3/0], ' ', [0/3/0], ' ', [0/3/0]
```

The 'Key' section shows the following information:

```
Key,
[1], 'trees.J48 \'-C 0.25 -M 2\' -21773316839364444'
[2], 'trees.RandomForest \'-P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1\' 1116839470751428698'
[3], 'bayes.NaiveBayes \\' 5995231201785697655'
[4], 'lazy.IBk \'-K 1 -W 0 -A \\'weka.core.neighboursearch.LinearNNSearch -A \\'weka.core.Euclidean
[5], 'lazy.IBk \'-K 3 -W 0 -A \\'weka.core.neighboursearch.LinearNNSearch -A \\'weka.core.Euclidean
[6], 'lazy.IBk \'-K 5 -W 0 -A \\'weka.core.neighboursearch.LinearNNSearch -A \\'weka.core.Euclidean
[7], 'lazy.IBk \'-K 10 -W 0 -A \\'weka.core.neighboursearch.LinearNNSearch -A \\'weka.core.Euclidean
```

The 'Result list' section shows the following results:

```
18:27:04 - F_measure - trees.J48 \'-C 0.25 -M 2\' -2177331
18:27:37 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:28:54 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:28:56 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:28:57 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:29:20 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:29:21 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:29:33 - F_measure - trees.J48 \'-C 0.25 -M 2\' -2177331
18:30:13 - Percent_correct - trees.J48 \'-C 0.25 -M 2\' -217
18:32:25 - Percent_incorrect - trees.J48 \'-C 0.25 -M 2\' -2
18:34:29 - IR_precision - trees.J48 \'-C 0.25 -M 2\' -2177331
18:35:05 - IR_recall - trees.J48 \'-C 0.25 -M 2\' -217733168
18:37:20 - F_measure - trees.J48 \'-C 0.25 -M 2\' -2177331
```

Results Analysis:

- For accuracy attribute:
 - o Glass dataset: best result with Random Forest Model
 - o Iris dataset: best result with IBK model with K=5 and K=10
 - o Weather dataset: best result with IBK model with K=5
- For precision attribute:
 - o Glass dataset: best result with J48 Model
 - o Iris dataset: all results=1
 - o Weather dataset: best result with Random Forest model
- For recall attribute:
 - o Glass dataset: best result with Random Forest and IBK with k=10 Models
 - o Iris dataset: all results=1
 - o Weather dataset: best result with IBK model with K=5 and K=10
- For f-value attribute:
 - o Glass dataset: best result with Random Forest Model
 - o Iris dataset: all results=1
 - o Weather dataset: best result with IBK model with K=3

Conclusions:

- Glass dataset: the model should predict the type of glass created depending on the matter used to create it.
The best model here is the one with higher precision rate (less false positive prediction), having higher precision will result in lesser costs.
⇒ Best model: Random Forest
- Iris dataset: based on the context, there are no losses caused by false negative or false positive (not a metric to choose the best model).
Accuracy would be the right attribute to choose the corresponding model, and based on the results above:
⇒ Best Model: IBK with K=5 and K=10
- Weather dataset: model is used to determine whether a player will go to practice or not based on the weather conditions.
If this model is used in training spaces (stadium preparation), false negative would have more impact than false positive, a model with less false negative is more suitable: lesser false negative implies higher recall.
⇒ Best models: IBK with K=5 and K=10