
Florida Atlantic University (FAU)

Assignment 2



Project objective: Modeling assignment: Classification using decision trees

Written By: Kevin Tudor

Table of Contents

Requirements	3
Part 1 - <i>Initial tree</i>	4
Part 2 - <i>Unpruned tree</i>	6
Part 3 - <i>Confidence Factor</i>	8
Part 4 - <i>Cost sensitivity</i>	10
Raw Weka Data	11
J48 trees	11
Best Cost Sensitive Classifier	23

Requirements

Part 1: *Initial tree*

This part of the project will allow you to predict a class (fp, nfp) using J48 (C4.5), a decision tree-based classification algorithm.

- Build a classification model using J48 (C4.5) using the fit data set and 10-fold cross validation.
- Determine the misclassification error rates (%) for both types of misclassifications **from the confusion matrix**.
 1. Type I: a nfp module is classified as fp
 2. Type II: a fp module is classified as nfp
- Record the number of leaves and nodes in the selected tree and represent the tree in the same way as in the textbook.
- Repeat the previous tasks using the test data set to evaluate the model.

Part 2: *Unpruned tree*

- Now in the J48 options, set the unpruned option to true. Rebuild the model in the same way as above, repeat all steps.
- Now that you have represented the unpruned tree, compare with the tree generated above, and determine the part that was pruned.

Part 3: *Confidence Factor*

- Now in the J48 options, set the confidence factor (C) to 0.01. Rebuild the model in the same way as for the initial tree (Part 1), repeat all the steps (of Part 1)
- How does the size of the new tree compare to one built in Part 1? Explain why. What part was pruned?

Part 4: *Cost sensitivity*

- Till now, we did not make any distinction between a Type I and a Type II error. However, in Software Quality Classification, **a Type II error is more serious than a Type I error**. Here, our objective is to **obtain a balanced misclassification rates with Type II as low as possible**.
- Use the cost sensitive classifier combined with J48 and determine the optimal cost ratio (set cost of a type I error to 1 and vary the cost of the Type II error), using 10-fold cross validation on the fit data set. Observe the trends in the misclassification rates. What happens when the cost of a Type II error decreases/increases?

Evaluate all the models on the test data set.

*** For tips on performing cost sensitive classification, [click here](#). ***

Part 1 - *Initial tree*

J48 pruned tree - Test mode: 10-fold cross-validation

Weka Tree

```

TOTOTORS <= 405
| NUMUORS <= 25: nfp (115.0/2.0)
| NUMUORS > 25
| | NLOGIC <= 7
| | | VG <= 8: fp (2.0)
| | | VG > 8: nfp (13.0/2.0)
| | NLOGIC > 7: fp (2.0)
TOTOTORS > 405
| NUMUORS <= 42
| | VG <= 48: fp (29.0/5.0)
| | VG > 48
| | | TOTOTORS <= 1522: nfp (7.0)
| | | TOTOTORS > 1522: fp (5.0/1.0)
| NUMUORS > 42: fp (19.0)

```

Confusion Matrix	fp	nfp
fp	40	15
nfp	10	127

Type I: $10/192 = 0.052083 = \mathbf{5.21\%}$

Type II: $15/192 = 0.078125 = \mathbf{7.81\%}$

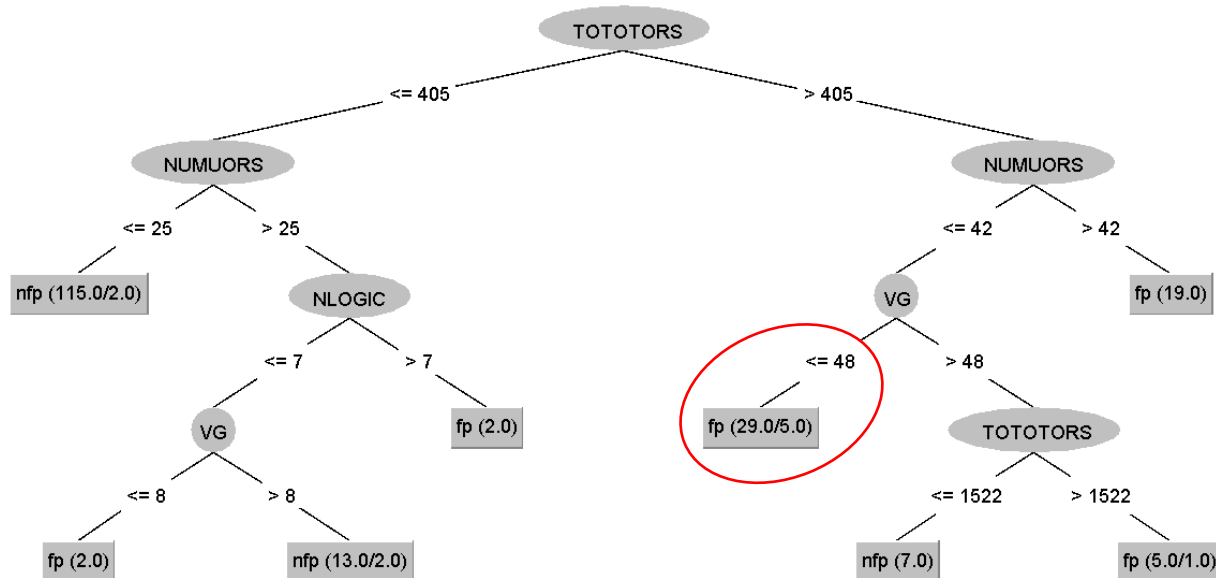
Number of Leaves: 8

Size of the tree: 15

Pruned branches and leaves from:

VG <= 48

Visualized tree



J48 pruned tree - Test mode: User Supplied**Weka Tree**

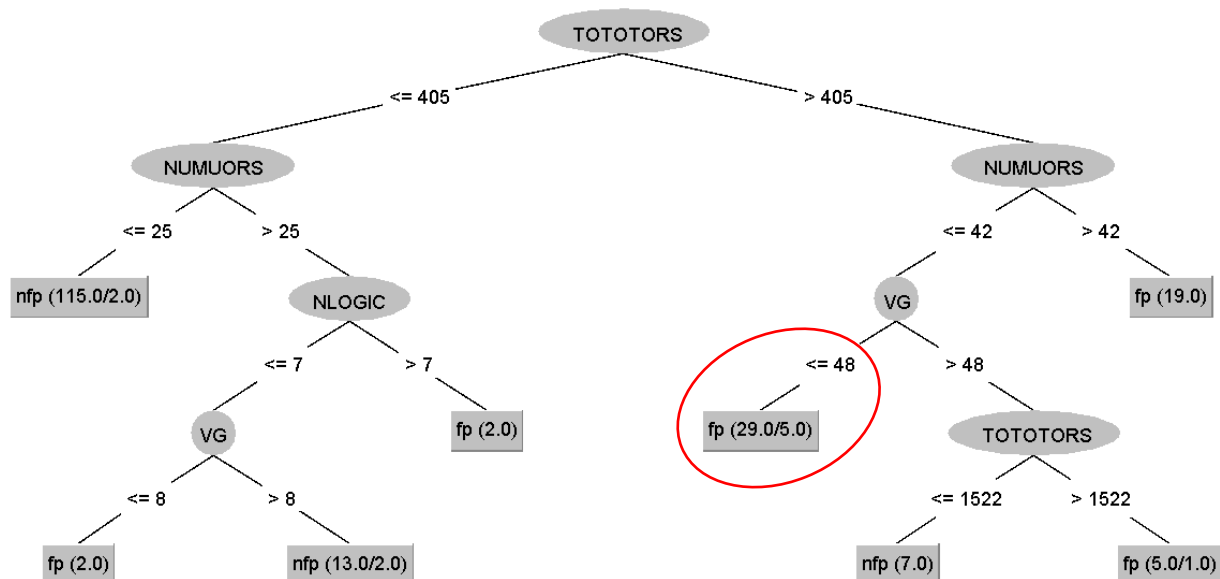
TOTOTORS <= 405
 | NUMUORS <= 25: nfp (115.0/2.0)
 | NUMUORS > 25
 | | NLOGIC <= 7
 | | | VG <= 8: fp (2.0)
 | | | VG > 8: nfp (13.0/2.0)
 | | NLOGIC > 7: fp (2.0)
 TOTOTORS > 405
 | NUMUORS <= 42
 | | **VG <= 48: fp (29.0/5.0)**
 | | VG > 48
 | | | TOTOTORS <= 1522: nfp (7.0)
 | | | TOTOTORS > 1522: fp (5.0/1.0)
 | NUMUORS > 42: fp (19.0)

Confusion Matrix	fp	nfp
fp	18	10
nfp	9	59

Type I: $9/192 = 0.046875 = \mathbf{4.69\%}$
 Type II: $10/192 = 0.052083 = \mathbf{5.21\%}$

Number of Leaves: 8
 Size of the tree: 15

Pruned branches and leaves from:
VG <= 48

Visualized tree

Part 2 - *Unpruned tree*

J48 unpruned tree - Test mode: 10-fold cross-validation

Weka Tree

TOTOTORS <= 405

| NUMUORS <= 25: nfp (115.0/2.0)

| NUMUORS > 25

| | NLOGIC <= 7

| | | VG <= 8: fp (2.0)

| | | VG > 8: nfp (13.0/2.0)

| | NLOGIC > 7: fp (2.0)

TOTOTORS > 405

| NUMUORS <= 42

| | **VG <= 48**

| | | NLOGIC <= 2: fp (19.0/1.0)

| | | NLOGIC > 2

| | | | NUMUORS <= 29: fp (3.0)

| | | | NUMUORS > 29

| | | | | VG <= 34: fp (4.0/1.0)

| | | | | VG > 34: nfp (3.0)

| | VG > 48

| | | TOTOTORS <= 1522: nfp (7.0)

| | | TOTOTORS > 1522: fp (5.0/1.0)

| NUMUORS > 42: fp (19.0)

Confusion Matrix	fp	nfp
fp	41	14
nfp	10	127

Type I: $10/192 = 0.052083 = \mathbf{5.21\%}$

Type II: $14/192 = 0.072916 = \mathbf{7.29\%}$

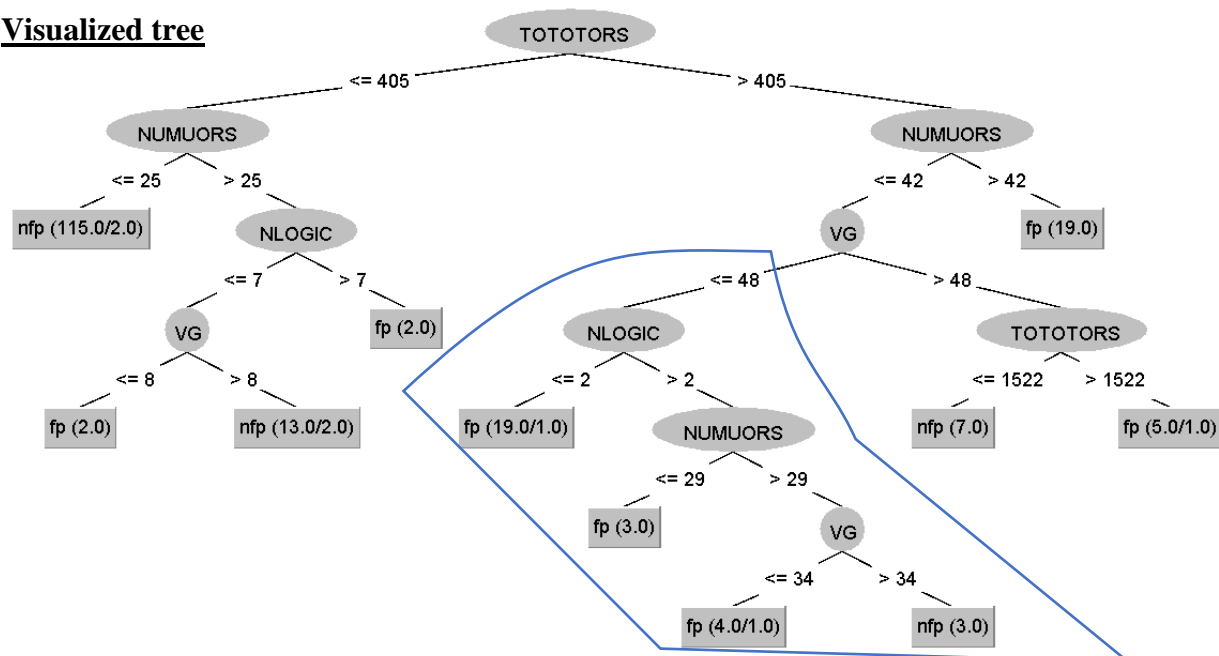
Number of Leaves: 11

Size of the tree: 21

Unpruned branches and leaves

from: **VG <= 48**

Visualized tree



J48 unpruned tree - Test mode: User Supplied**Weka Tree**

```

TOTOTORS <= 405
| NUMUORS <= 25: nfp (115.0/2.0)
| NUMUORS > 25
| | NLOGIC <= 7
| | | VG <= 8: fp (2.0)
| | | VG > 8: nfp (13.0/2.0)
| | NLOGIC > 7: fp (2.0)
TOTOTORS > 405
| NUMUORS <= 42
| | VG <= 48
| | | NLOGIC <= 2: fp (19.0/1.0)
| | | NLOGIC > 2
| | | | NUMUORS <= 29: fp (3.0)
| | | | NUMUORS > 29
| | | | | VG <= 34: fp (4.0/1.0)
| | | | | VG > 34: nfp (3.0)
| | | VG > 48
| | | | TOTOTORS <= 1522: nfp (7.0)
| | | | TOTOTORS > 1522: fp (5.0/1.0)
| NUMUORS > 42: fp (19.0)

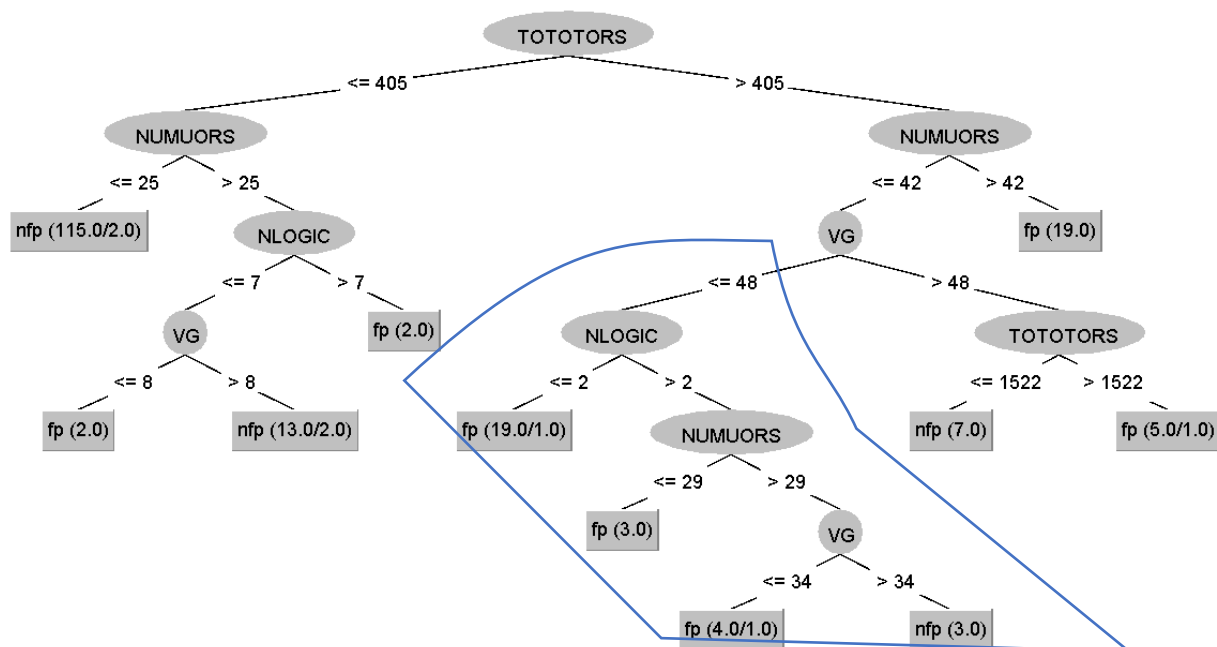
```

Confusion Matrix	fp	nfp
fp	18	10
nfp	9	59

Type I: $9/192 = 0.046875 = \mathbf{4.69\%}$
Type II: $10/192 = 0.052083 = \mathbf{5.21\%}$

Number of Leaves: 11
Size of the tree: 21

**Unpruned branches and leaves
from: $\mathbf{VG \leq 48}$**

Visualized tree

Part 3 - *Confidence Factor*

J48 pruned tree (C) = 0.01 Test mode: 10-fold cross-validation

Weka Tree

TOTOTORS ≤ 405 : nfp (132.0/8.0)

TOTOTORS > 405 : fp (60.0/13.0)

The Tree is much smaller with confidence 0.01. Both branches from TOTOTORS have been pruned and replaced with Leaves. This may be because the confidence of 0.01 causes the model to generalize the child Leaves during training resulting in a total of two leaf nodes.

Confusion Matrix	fp	nfp
fp	38	17
nfp	12	125

Type I: $12/192 = 0.0625 = 6.25\%$

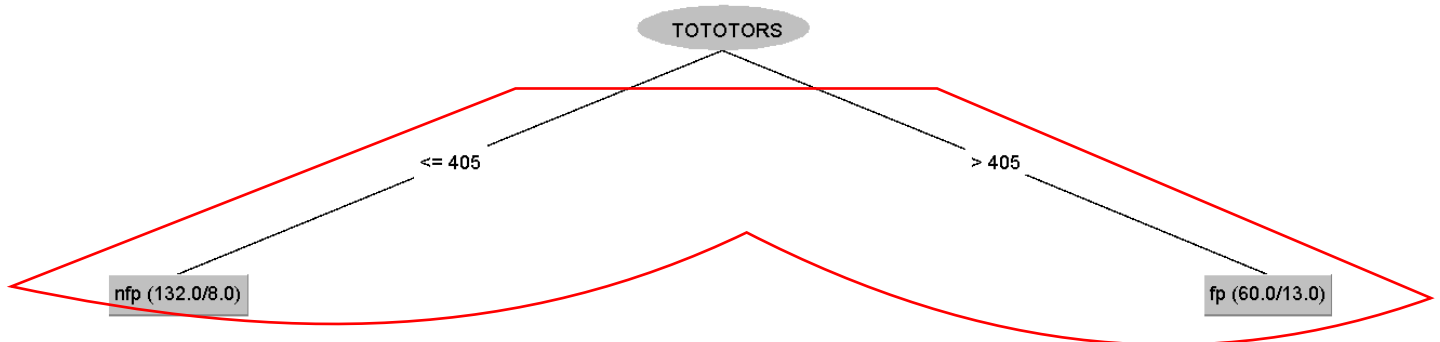
Type II: $17/192 = 0.08854 = 8.85\%$

Number of Leaves: 2

Size of the tree: 3

Pruned branches and leaves from:
TOTOTORS

Visualized tree



J48 pruned tree (C) = 0.01 Test mode: User Supplied**Weka Tree**

TOTOTORS ≤ 405 : nfp (132.0/8.0)

TOTOTORS > 405 : fp (60.0/13.0)

The Tree is much smaller with confidence 0.01. Both branches from TOTOTORS have been pruned and replaced with Leaves. This may be because the confidence of 0.01 causes the model to generalize the child Leaves during training resulting in a total of two leaf nodes.

Confusion Matrix	fp	nfp
fp	24	4
nfp	10	58

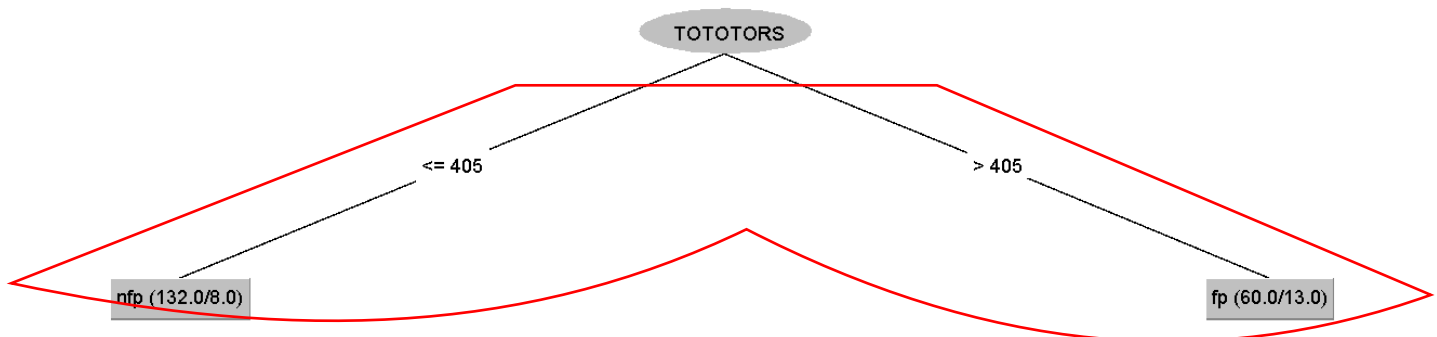
Type I: $10/192 = 0.052083 = \mathbf{5.21\%}$

Type II: $4/192 = 0.020833 = \mathbf{2.08\%}$

Number of Leaves: 2

Size of the tree: 3

Pruned branches and leaves from:
TOTOTORS

Visualized tree

Part 4 - *Cost sensitivity*

Sensitivity (∓ 0.1)	Type I %	Type II %
II:0.5	8/192	14/192
II:0.6	11/192	14/192
II:0.7	9/192	14/192
II:0.8	9/192	15/192
II:0.9	10/192	16/192
II:1.0	10/192	15/192
II:1.1	14/192	11/192
II:1.2	15/192	11/192
II:1.3	14/192	11/192
II:1.4	15/192	11/192
II:1.5	15/192	11/192

Sensitivity (∓ 0.05)	Type I %	Type II %
II:1.05	12/192	12/192
II:1.10	14/192	11/192
II:1.15	15/192	11/192

Sensitivity (BEST)	Type I %	Type II %
II:1.02	11/192	11/192

- Obtain a balanced misclassification rates with Type II as low as possible.
- Observe the trends in the misclassification rates. What happens when the cost of a Type II error decreases/increases?

By altering the cost of Type II, the optimal classification rates (balanced with Type II as low as possible) can be found. **When the cost of Type II is increased the Type I misclassifications go up while the Type II misclassifications go down. When the cost of Type II is decreased the Type I misclassifications go down while the Type II misclassifications go up.** Although increasing the cost of Type II lowers the Type II misclassifications the ratio is *unbalanced*. Around Type II cost of 1.1 *the lowest Type II cost is found at 11/192* whereas the Type I Misclassifications are at 14/192. To further optimize the model, Type II can be incremented or decremented by 0.05 from 1.10. Upon 0.05 modifications it is evident the most optimal ratio is with *Type II at 1.05 with a balanced 12/192 for both Types*. To optimize the model further, the cost can be incremented by 0.01 from 1.05. Ultimately **the best model was found with a Type II cost of 1.02 giving both Types a low and balanced value of 11/192 (correctly classified 88.5% overall best!).**

Raw Weka Data

J48 trees

=== Run information ===

Scheme: weka.classifiers.trees.J48 -C 0.25 -M 2

Relation: TudorFitClassifier

Instances: 192

Attributes: 9

NUMUORS
 NUMUANDS
 TOTOTORS
 TOTOPANDS
 VG
 NLOGIC
 LOC
 ELOC
 class

Test mode: 10-fold cross-validation

=== Classifier model (full training set) ===

J48 pruned tree

```

-----
TOTOTORS <= 405
| NUMUORS <= 25: nfp (115.0/2.0)
| NUMUORS > 25
| | NLOGIC <= 7
| | | VG <= 8: fp (2.0)
| | | VG > 8: nfp (13.0/2.0)
| | NLOGIC > 7: fp (2.0)
TOTOTORS > 405
| NUMUORS <= 42
| | VG <= 48: fp (29.0/5.0)
| | VG > 48
| | | TOTOTORS <= 1522: nfp (7.0)
| | | TOTOTORS > 1522: fp (5.0/1.0)
| NUMUORS > 42: fp (19.0)

```

Number of Leaves : 8

Size of the tree : 15

Time taken to build model: 0.01 seconds

=== Stratified cross-validation ===

=== Summary ===

Correctly Classified Instances	167	86.9792 %
Incorrectly Classified Instances	25	13.0208 %
Kappa statistic	0.6726	
Mean absolute error	0.1522	
Root mean squared error	0.3457	
Relative absolute error	37.1236 %	
Root relative squared error	76.4495 %	
Total Number of Instances	192	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area
Class								
	0.727	0.073	0.800	0.727	0.762	0.674	0.810	fp
	0.927	0.273	0.894	0.927	0.910	0.674	0.810	nfp
Weighted Avg.	0.870	0.216	0.867	0.870	0.868	0.674	0.810	0.813

=== Confusion Matrix ===

a b <-- classified as

40 15 | a = fp

10 127 | b = nfp

=== Run information ===

Scheme: weka.classifiers.trees.J48 -C 0.25 -M 2

Relation: TudorFitClassifier

Instances: 192

Attributes: 9

NUMUORS

NUMUANDS

TOTOTORS

TOTOPANDS

VG

NLOGIC

LOC

ELOC

class

Test mode: user supplied test set: size unknown (reading incrementally)

=== Classifier model (full training set) ===

J48 pruned tree

TOTOTORS <= 405

| NUMUORS <= 25: nfp (115.0/2.0)

| NUMUORS > 25

| | NLOGIC <= 7

| | | VG <= 8: fp (2.0)

| | | VG > 8: nfp (13.0/2.0)

| | NLOGIC > 7: fp (2.0)

TOTOTORS > 405

| NUMUORS <= 42

| | VG <= 48: fp (29.0/5.0)

| | VG > 48

| | | TOTOTORS <= 1522: nfp (7.0)

| | | TOTOTORS > 1522: fp (5.0/1.0)

| NUMUORS > 42: fp (19.0)

Number of Leaves : 8

Size of the tree : 15

Time taken to build model: 0 seconds

=== Evaluation on test set ===

Time taken to test model on supplied test set: 0 seconds

=== Summary ===

Correctly Classified Instances	77	80.2083 %
Incorrectly Classified Instances	19	19.7917 %
Kappa statistic	0.5159	
Mean absolute error	0.2103	
Root mean squared error	0.4136	
Relative absolute error	51.0445 %	
Root relative squared error	90.9892 %	
Total Number of Instances	96	

=== Detailed Accuracy By Class ===

Class	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area
	0.643	0.132	0.667	0.643	0.655	0.516	0.697	fp
	0.868	0.357	0.855	0.868	0.861	0.516	0.697	nfp
Weighted Avg.	0.802	0.292	0.800	0.802	0.801	0.516	0.697	0.774

=== Confusion Matrix ===

```

a b <-- classified as
18 10 | a = fp
9 59 | b = nfp

```

=== Run information ===

Scheme: weka.classifiers.trees.J48 -U -M 2

Relation: TudorFitClassifier

Instances: 192

Attributes: 9

NUMUORS

NUMUANDS

TOTOTORS

TOTOPANDS

VG

NLOGIC

LOC

ELOC

class

Test mode: 10-fold cross-validation

=== Classifier model (full training set) ===

J48 unpruned tree

TOTOTORS <= 405

| NUMUORS <= 25: nfp (115.0/2.0)

| NUMUORS > 25

| | NLOGIC <= 7

| | | VG <= 8: fp (2.0)

| | | VG > 8: nfp (13.0/2.0)

| | NLOGIC > 7: fp (2.0)

TOTOTORS > 405

| NUMUORS <= 42

| | VG <= 48

| | | NLOGIC <= 2: fp (19.0/1.0)

| | | NLOGIC > 2

| | | | NUMUORS <= 29: fp (3.0)

| | | | NUMUORS > 29

| | | | | VG <= 34: fp (4.0/1.0)

| | | | | VG > 34: nfp (3.0)

| | | VG > 48

| | | TOTOTORS <= 1522: nfp (7.0)

| | | TOTOTORS > 1522: fp (5.0/1.0)

| NUMUORS > 42: fp (19.0)

Number of Leaves : 11

Size of the tree : 21

Time taken to build model: 0 seconds

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

==== Summary ====

Correctly Classified Instances	168	87.5 %
Incorrectly Classified Instances	24	12.5 %
Kappa statistic	0.6874	
Mean absolute error	0.1426	
Root mean squared error	0.3388	
Relative absolute error	34.7902 %	
Root relative squared error	74.9185 %	
Total Number of Instances	192	

==== Detailed Accuracy By Class ====

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area
Class								
	0.745	0.073	0.804	0.745	0.774	0.688	0.817	fp
	0.927	0.255	0.901	0.927	0.914	0.688	0.817	nfp
Weighted Avg.	0.875	0.203	0.873	0.875	0.874	0.688	0.817	0.823

==== Confusion Matrix ====

```

a  b  <-- classified as
41 14 | a = fp
10 127 | b = nfp

```


=== Run information ===

Scheme: weka.classifiers.trees.J48 -U -M 2

Relation: TudorFitClassifier

Instances: 192

Attributes: 9

NUMUORS
 NUMUANDS
 TOTOTORS
 TOTOPANDS
 VG
 NLOGIC
 LOC
 ELOC
 class

Test mode: user supplied test set: size unknown (reading incrementally)

=== Classifier model (full training set) ===

J48 unpruned tree

```
TOTOTORS <= 405
| NUMUORS <= 25: nfp (115.0/2.0)
| NUMUORS > 25
| | NLOGIC <= 7
| | | VG <= 8: fp (2.0)
| | | VG > 8: nfp (13.0/2.0)
| | NLOGIC > 7: fp (2.0)
TOTOTORS > 405
| NUMUORS <= 42
| | VG <= 48
| | | NLOGIC <= 2: fp (19.0/1.0)
| | | NLOGIC > 2
| | | | NUMUORS <= 29: fp (3.0)
| | | | NUMUORS > 29
| | | | | VG <= 34: fp (4.0/1.0)
| | | | | VG > 34: nfp (3.0)
| | | VG > 48
| | | TOTOTORS <= 1522: nfp (7.0)
| | | TOTOTORS > 1522: fp (5.0/1.0)
| NUMUORS > 42: fp (19.0)
```

Number of Leaves : 11

Size of the tree : 21

Time taken to build model: 0 seconds

==== Evaluation on test set ====

Time taken to test model on supplied test set: 0 seconds

==== Summary ====

Correctly Classified Instances	77	80.2083 %
Incorrectly Classified Instances	19	19.7917 %
Kappa statistic	0.5159	
Mean absolute error	0.2128	
Root mean squared error	0.4292	
Relative absolute error	51.6503 %	
Root relative squared error	94.4258 %	
Total Number of Instances	96	

==== Detailed Accuracy By Class ====

Class	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area
	0.643	0.132	0.667	0.643	0.655	0.516	0.697	fp
	0.868	0.357	0.855	0.868	0.861	0.516	0.697	nfp
Weighted Avg.	0.802	0.292	0.800	0.802	0.801	0.516	0.697	0.776

==== Confusion Matrix ====

```

a b <-- classified as
18 10 | a = fp
9 59 | b = nfp

```

=== Run information ===

Scheme: weka.classifiers.trees.J48 -C 0.01 -M 2

Relation: TudorFitClassifier

Instances: 192

Attributes: 9

NUMUORS

NUMUANDS

TOTOTORS

TOTOPANDS

VG

NLOGIC

LOC

ELOC

class

Test mode: 10-fold cross-validation

=== Classifier model (full training set) ===

J48 pruned tree

TOTOTORS <= 405: nfp (132.0/8.0)

TOTOTORS > 405: fp (60.0/13.0)

Number of Leaves : 2

Size of the tree : 3

Time taken to build model: 0 seconds

=== Stratified cross-validation ===

=== Summary ===

Correctly Classified Instances	163	84.8958 %
Incorrectly Classified Instances	29	15.1042 %
Kappa statistic	0.6202	
Mean absolute error	0.1996	
Root mean squared error	0.3649	
Relative absolute error	48.6889 %	
Root relative squared error	80.6766 %	
Total Number of Instances	192	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area
Class								
	0.691	0.088	0.760	0.691	0.724	0.622	0.787	0.658
	0.912	0.309	0.880	0.912	0.896	0.622	0.787	0.850
Weighted Avg.	0.849	0.246	0.846	0.849	0.847	0.622	0.787	0.795

=== Confusion Matrix ===

```

a  b  <-- classified as
38 17 | a = fp
12 125 | b = nfp

```

=== Run information ===

Scheme: weka.classifiers.trees.J48 -C 0.01 -M 2

Relation: TudorFitClassifier

Instances: 192

Attributes: 9

NUMUORS

NUMUANDS

TOTOTORS

TOTOPANDS

VG

NLOGIC

LOC

ELOC

class

Test mode: user supplied test set: size unknown (reading incrementally)

=== Classifier model (full training set) ===

J48 pruned tree

TOTOTORS <= 405: nfp (132.0/8.0)

TOTOTORS > 405: fp (60.0/13.0)

Number of Leaves : 2

Size of the tree : 3

Time taken to build model: 0 seconds

=== Evaluation on test set ===

Time taken to test model on supplied test set: 0 seconds

=== Summary ===

Correctly Classified Instances	82	85.4167 %
Incorrectly Classified Instances	14	14.5833 %
Kappa statistic	0.668	
Mean absolute error	0.2115	
Root mean squared error	0.3386	
Relative absolute error	51.3474 %	
Root relative squared error	74.4905 %	
Total Number of Instances	96	

=== Detailed Accuracy By Class ===

Class	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area
	0.857	0.147	0.706	0.857	0.774	0.675	0.855	0.647
	0.853	0.143	0.935	0.853	0.892	0.675	0.855	0.902
Weighted Avg.	0.854	0.144	0.869	0.854	0.858	0.675	0.855	0.828

=== Confusion Matrix ===

```
a b <-- classified as
24 4 | a = fp
10 58 | b = nfp
```

Best Cost Sensitive Classifier

=== Run information ===

Scheme: weka.classifiers.meta.**CostSensitiveClassifier** -cost-matrix "[0.0 1.02; 1.0 0.0]" -S 1

-W weka.classifiers.trees.J48 -- -C 0.25 -M 2

Relation: TudorFitClassifier

Instances: 192

Attributes: 9

NUMUORS

NUMUANDS

TOTOTORS

TOTOPANDS

VG

NLOGIC

LOC

ELOC

class

Test mode: 10-fold cross-validation

=== Classifier model (full training set) ===

CostSensitiveClassifier using reweighted training instances

weka.classifiers.trees.J48 -C 0.25 -M 2

Classifier Model

J48 pruned tree

TOTOTORS <= 405

| NUMUORS <= 25: nfp (114.38/2.03)

| NUMUORS > 25

| | ELOC <= 49: nfp (5.97)

| | ELOC > 49

| | | VG <= 15: fp (4.06)

| | | VG > 15: nfp (7.0/2.03)

TOTOTORS > 405

| NUMUORS <= 42

| | VG <= 48: fp (29.31/4.97)

| | VG > 48

| | | TOTOTORS <= 1522: nfp (6.96)

| | | TOTOTORS > 1522: fp (5.05/0.99)

| NUMUORS > 42: fp (19.27)

Number of Leaves : 8

Size of the tree : 15

Cost Matrix

```
0  1.02
1  0
```

Time taken to build model: 0 seconds

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

Correctly Classified Instances	170	88.5417 %
Incorrectly Classified Instances	22	11.4583 %
Kappa statistic	0.7197	
Mean absolute error	0.1342	
Root mean squared error	0.3251	
Relative absolute error	32.736 %	
Root relative squared error	71.8888 %	
Total Number of Instances	192	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area
Class								
	0.800	0.080	0.800	0.800	0.800	0.720	0.847	fp
	0.920	0.200	0.920	0.920	0.920	0.720	0.846	nfp
Weighted Avg.	0.885	0.166	0.885	0.885	0.885	0.720	0.846	0.842

=== Confusion Matrix ===

```
a  b  <-- classified as
44 11 | a = fp
11 126 | b = nfp
```