



RAPIDMINER 5.0

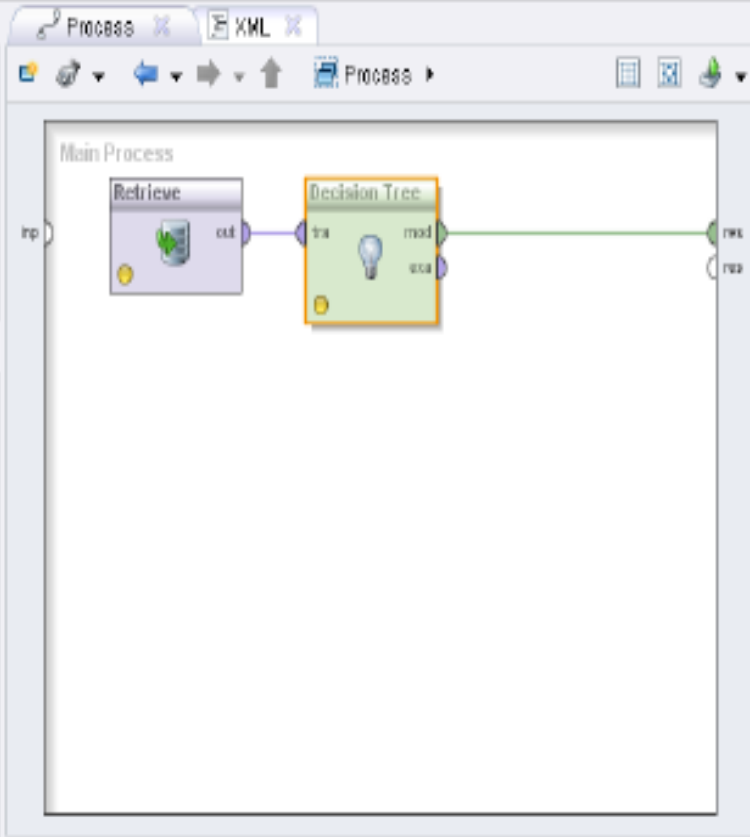
Ricardo Aler Mur

ÍNDICE

1. Árbol de decisión
2. Reglas
3. Árbol de decisión con validación cruzada
4. Árbol de decisión vs. knn con validación cruzada
5. Árbol de decisión vs. knn con validación cruzada evaluando el coste (con matriz de costes)
6. Árbol de decisión vs. Metacost con árboles de decisión y matriz de costes
7. Comparación de curvas ROC

1. ÁRBOL DE DECISIÓN

- Ejemplo 1 del tutorial:
- http://www.dataprix.com/files/RapidMiner_Tutorial_online_Operadores.pdf



- Operators Repositories
- [Filter]
- Process Control (5)
 - Utility (1)
 - Repository Access (2)
 - Retrieve
 - Store
 - Import (1)
 - Export (1)
 - Data Transformation
 - Modeling
 - Classification and Regression
 - Lazy Modeling (2)
 - Bayesian Modeling (2)
 - Tree Induction (8)
 - Decision Tree**
 - Decision Tree (Multiway)
 - Decision Tree (Weight-Based)
 - ID3
 - CHAID
 - Decision Stump
 - Random Tree
 - Random Forest
 - Rule Induction (5)
 - Neural Net Training (2)
 - Execution Utilities (6)

Parameters

Decision Tree

criterion:
 minimal size for split:
 minimal leaf size:
 minimal gain:
 maximal depth:
 confidence:

⚠ 3 hidden expert parameters

Problems Log

No Errors

Message	Files	Location

Help Comment

Synopsis

Learns a pruned decision tree which can handle both numerical and nominal attributes.

Description

This operator learns decision trees from both nominal and numerical data. Decision trees are powerful classification methods which often can also easily be understood. This decision tree learner works similar to Quinlan's C4.5 or CART. The actual type of the tree is determined by the criterion, e.g. using gain_ratio or gini for CART / C4.5.

Input

- training set: expects ExampleSet

Output

- model
- exampleSet

ÁRBOL DE DECISIÓN

The image shows a screenshot of the RapidMiner software interface. The main window displays a decision tree visualization. The root node is "Outlook", which branches into three categories: "overcast", "rain", and "sunny". The "overcast" branch leads to a leaf node labeled "yes" (red). The "rain" branch leads to a "Wind" node, which branches into "false" (leading to "yes", red) and "true" (leading to "no", blue). The "sunny" branch leads to a "Humidity" node, which branches into " > 77 " (leading to "no", blue) and " ≤ 77.500 " (leading to "yes", red).

On the right side, there is a "Repositories" panel showing a tree structure of data sources:

- NewLocalRepository (an anonymous)
- sample (an anonymous)
- samples (an anonymous)
- data (an anonymous)
- Golf-Testset1 (1 - anonymous - 429 bytes)
- Golf (1 - anonymous - 429 bytes)
- Iris (1 - anonymous - 7109 bytes)
- Labor-Negotiations (1 - anonymous - 3110 bytes)
- Market-Data (1 - anonymous - 133 bytes)
- Polynomial (1 - anonymous - 9780 bytes)
- Ripley-Set (1 - anonymous - 4377 bytes)
- Sonar (1 - anonymous - 102276 bytes)
- Transactions (1 - anonymous - 390 bytes)

At the bottom, there is a "Log" panel with the following entries:

```
Dec 5, 2009 9:36:17 PM INFO: Process file version is 50
Dec 5, 2009 10:02:21 PM INFO: No filename given for result file, using stdout for logging results!
Dec 5, 2009 10:02:21 PM INFO: Process starts
Dec 5, 2009 10:02:21 PM INFO: Process finished successfully after 0 s
Dec 5, 2009 10:07:02 PM INFO: No filename given for result file, using stdout for logging results!
Dec 5, 2009 10:07:02 PM INFO: Process starts
Dec 5, 2009 10:07:04 PM INFO: Process finished successfully after 1 s
```

On the far right, there is a "stem Monitor" panel showing a grid and a progress bar at 0%.

TECLAS IMPORTANTES

- F8 o View -> Perspectives -> Design
 - Vista de procesos
- F9 o View -> Perspectives -> Results
 - Vista de resultados

2. REGLA DE DECISIÓN

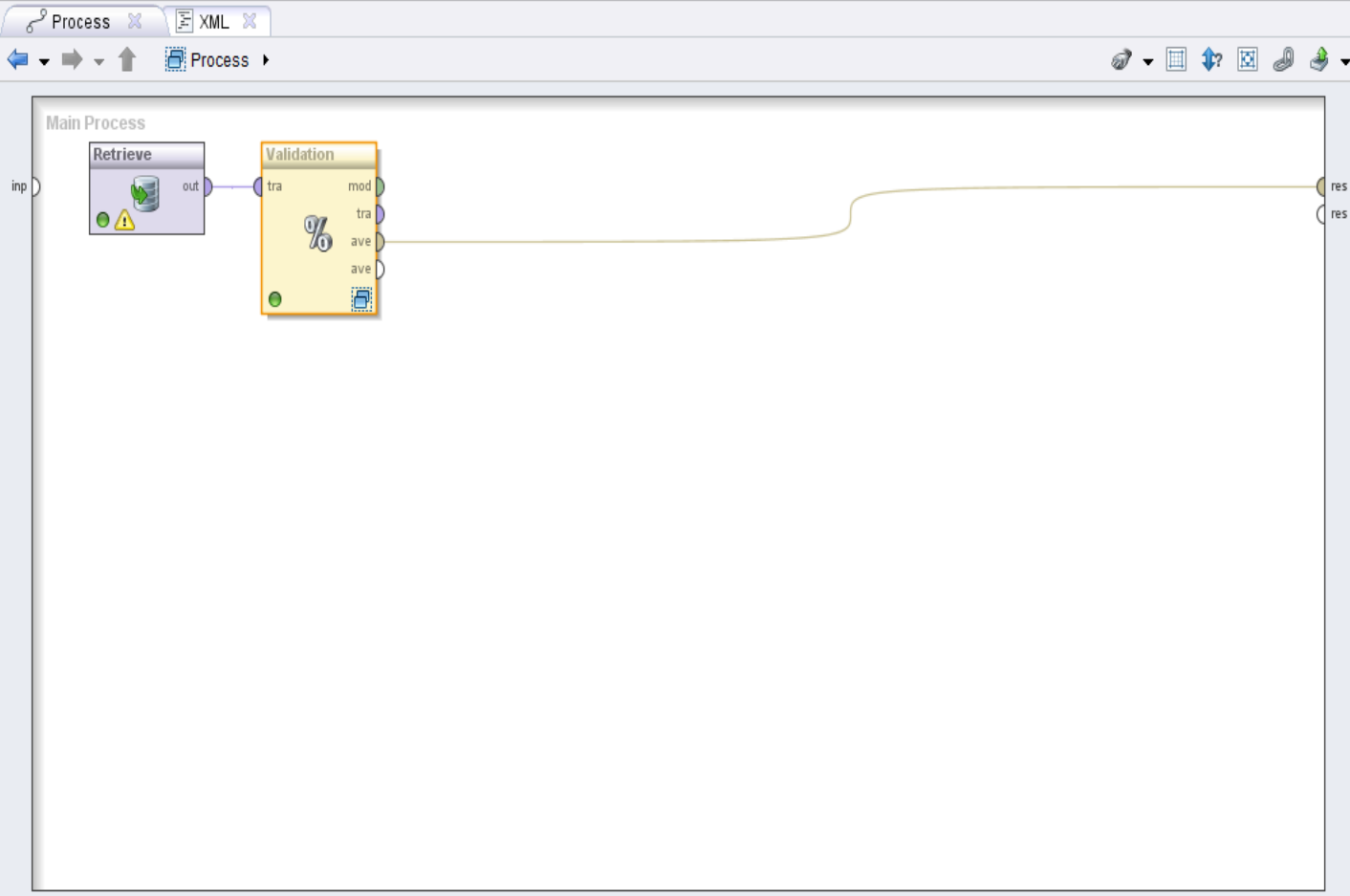
- Pinchar con el botón derecho del ratón sobre la caja del árbol de decisión y reemplazar el modelo para usar reglas (Rule induction)

8. Reemplazar el aprendiz por otro esquema de aprendizaje para tareas de clasificación. Hacer clic derecho sobre el operador **Decision Tree** y seleccionar **Replace Operator** → **Modeling** → **Classification and Regression** → **Rule Induction** → **Rule Induction**. Después de ejecutar el proceso cambiado con este ejemplo, se presenta el Nuevo modelo:

```
IF Cielo = Cubierto THEN Sí  
IF Temperatura ≤ 77.500 AND Ventoso = Falso AND Cielo = Lluvioso THEN  
Sí  
IF Cielo = Lluvioso THEN No  
IF Humedad > 77.500 THEN No ELSE Sí
```

3. CLASIFICACIÓN CON VALIDACIÓN CRUZADA

- Ejemplo 9 del tutorial, pero para clasificación



Parameters Context

% Validation (X-Validation)

- average performances only
- leave one out

number of validations

sampling type

- use local random seed
- parallelize training
- parallelize testing

Compatibility level

Problems Log

One potential problem

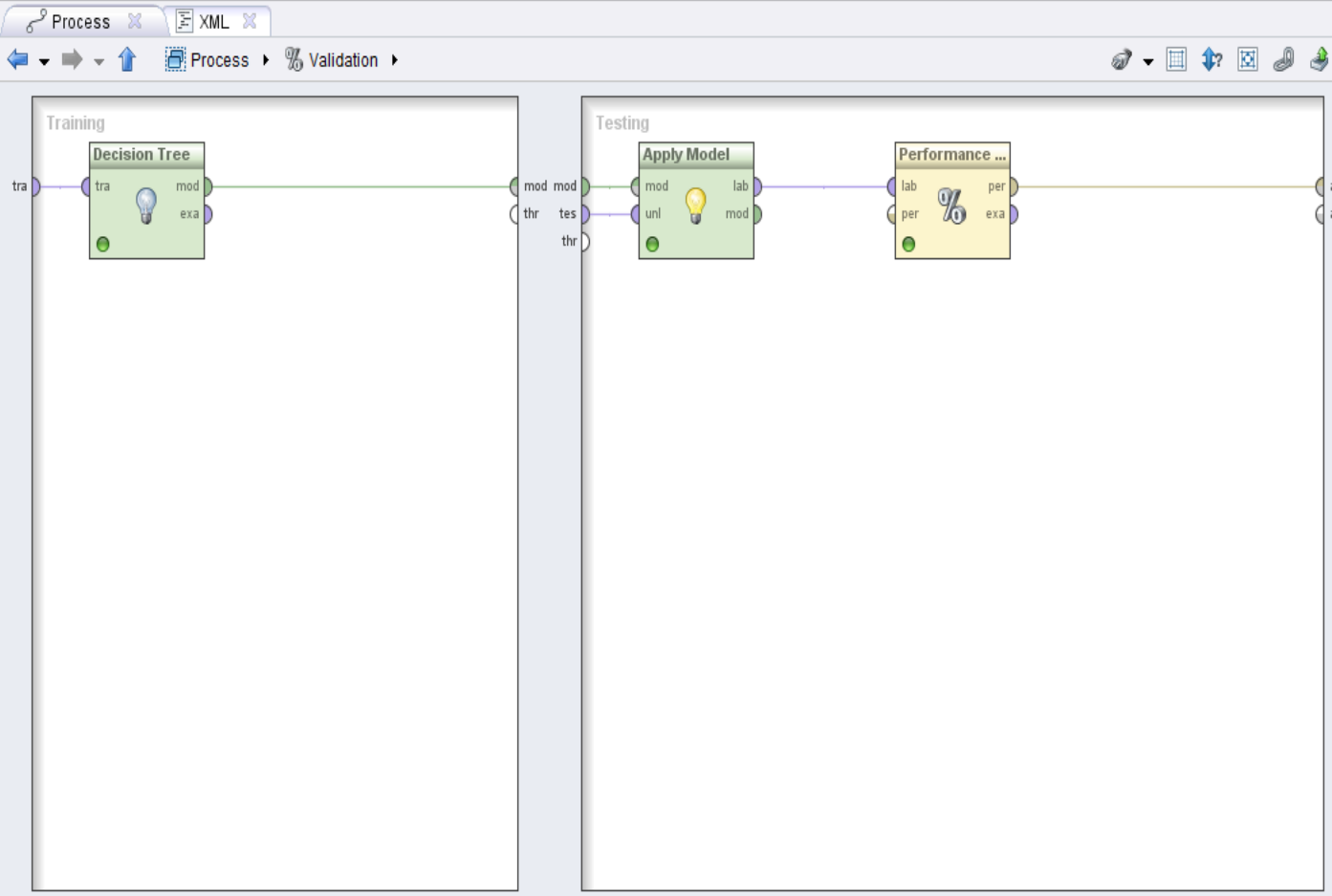
Message	Fixes	Location
Parameter 'repository entry' accesses a repository by name (//Local Repository/data/oil... No quick fix available		Retrieve

Help Comment

% X-Validation (RapidMiner Core)

Synopsis

This operator performs a cross-validation in order to estimate the statistical performance of a learning operator (usually on unseen data sets). It is mainly used to estimate how accurately a model (learnt by a particular learning operator) will perform in practice.



Parameters Context

% Validation (X-Validation)

- average performances only
- leave one out
- number of validations:
- sampling type:
- use local random seed
- parallelize training
- parallelize testing
- Compatibility level:

Problems Log

One potential problem

Message	Fixes	Location
Parameter 'repository entry' accesses a repository by name (//Local Repository/data/oil... ? No quick fix available	Retrieve	

Help Comment

% X-Validation (RapidMiner Core)

Synopsis

This operator performs a cross-validation in order to estimate the statistical performance of a learning operator (usually on unseen data sets). It is mainly used to estimate how accurately a model (learnt by a particular learning operator) will perform in practice.

4. ÁRBOL DE DECISIÓN Y KNN CON VALIDACIÓN CRUZADA

- Usar File -> Open **Template** -> Compare Learning Algorithms by Significance Test
- Reemplazar el primer clasificador por un árbol de decisión
- Reemplazar el segundo clasificador por KNN
- Reemplazar ambos “regression performance measurement” por “binomial performance measurement”. Hacer que se impriman, aparte de “accuracy”, también “AUC”



Table / Plot View Text View Annotations

- Criterion Selector
- classification_error
- AUC
- false_positive
- false_negative
- true_positive
- true_negative

Multiclass Classification Performance Annotations

Table View Plot View

classification_error: 7.15% +/- 0.82% (mikro: 7.15%)

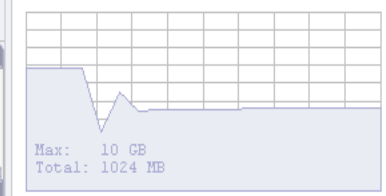
	true 2.000	true 1.000	class precision
pred. 2.000	6	32	15.79%
pred. 1.000	35	864	96.11%
class recall	14.63%	96.43%	

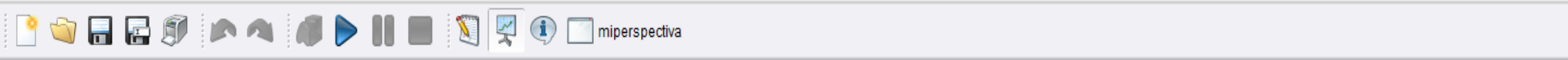
- Samples (none)
- DB
- Local Repository (Aler)

```

Nov 27, 2013 8:24:15 PM INFO: Executing process concurrently: Training
Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Testing
Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Training
Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Testing
Nov 27, 2013 8:24:19 PM INFO: Saving results.
Nov 27, 2013 8:24:19 PM INFO: Process //Local Repository/processes/PRUEBAS finished successfully after 2 s

```





- Criterion Selector
- classification_error
- AUC
- false_positive
- false_negative
- true_positive
- true_negative

Multiclass Classification Performance
 Annotations
 Table View
 Plot View

classification_error: 4.69% +/- 1.97% (mikro: 4.70%)

	true 2.000	true 1.000	class precision
pred. 2.000	12	15	44.44%
pred. 1.000	29	881	96.81%
class recall	29.27%	98.33%	

- Samples (none)
- DB
- Local Repository (Aler)

Nov 27, 2013 8:24:15 PM INFO: Executing process concurrently: Training
 Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Testing
 Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Training
 Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Testing
 Nov 27, 2013 8:24:19 PM INFO: Saving results.
 Nov 27, 2013 8:24:19 PM INFO: Process //Local Repository/processes/PRUEBAS finished successfully after 2 s

Max: 10 GB
 Total: 1024 MB



T-Test Significance

	0.047 +/- 0.020	0.071 +/- 0.008
0.047 +/- 0.020		0.002
0.071 +/- 0.008		

Probabilities for random values with the same result.
Bold values are smaller than alpha=0.050 which indicates a probably significant difference between the actual mean values!

- Samples (none)
- DB
- Local Repository (Aler)

```
Nov 27, 2013 8:24:15 PM INFO: Executing process concurrently: Training
Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Testing
Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Training
Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Testing
Nov 27, 2013 8:24:19 PM INFO: Saving results.
Nov 27, 2013 8:24:19 PM INFO: Process //Local Repository/processes/PRUEBAS finished successfully after 2 s
```

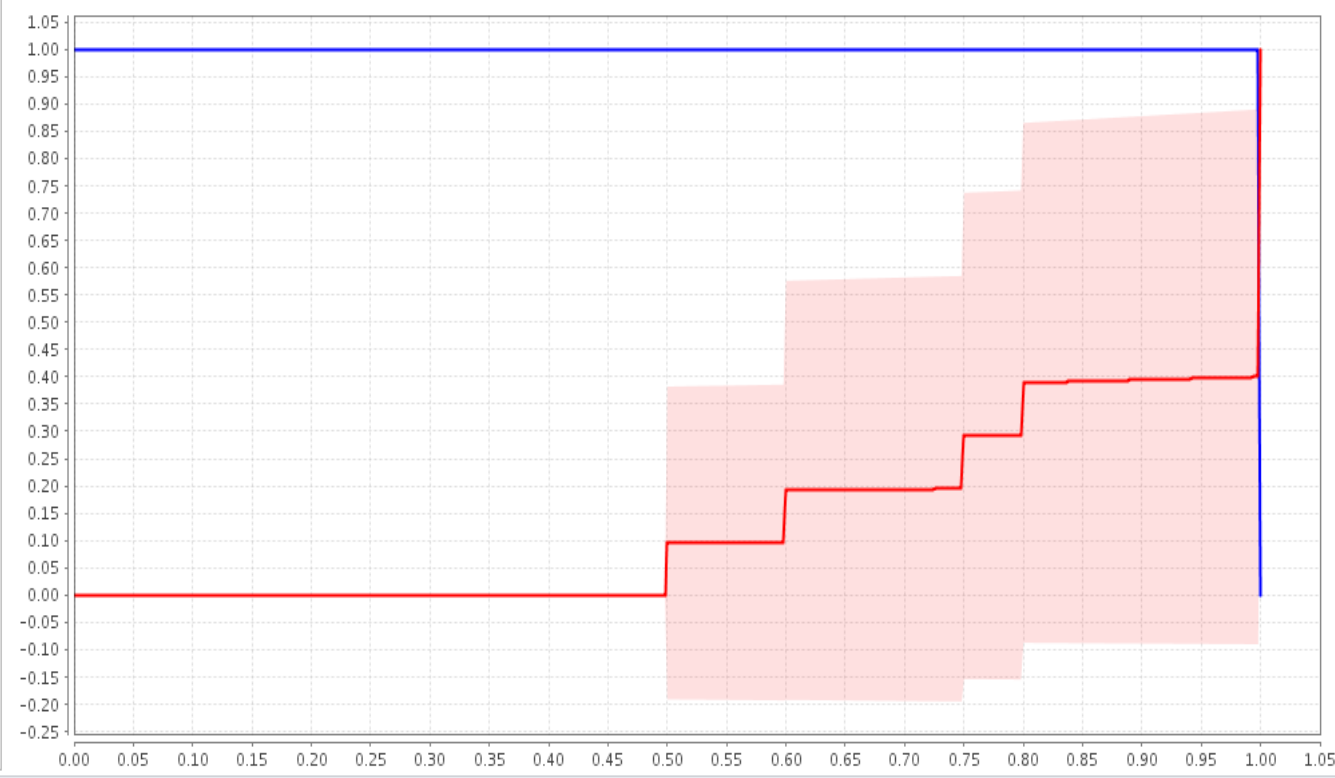

Max: 10 GB
Total: 1024 MB



- Criterion Selector
- classification_error
- AUC**
- false_positive
- false_negative
- true_positive
- true_negative

AUC: 0.500 +/- 0.000 (mikro: 0.500) (positive class: 1.000)

ROC ROC (Thresholds)



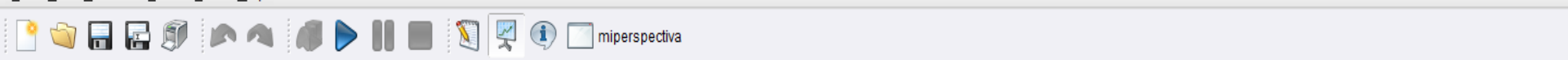
- Samples (none)
- DB
- Local Repository (Alert)

```

Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Training
Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Testing
Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Training
Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Testing
Nov 27, 2013 8:24:19 PM INFO: Saving results.
Nov 27, 2013 8:24:19 PM INFO: Process //Local Repository/processes/PRUEBAS finished successfully after 2 s

```


Max: 10 GB
Total: 1024 MB



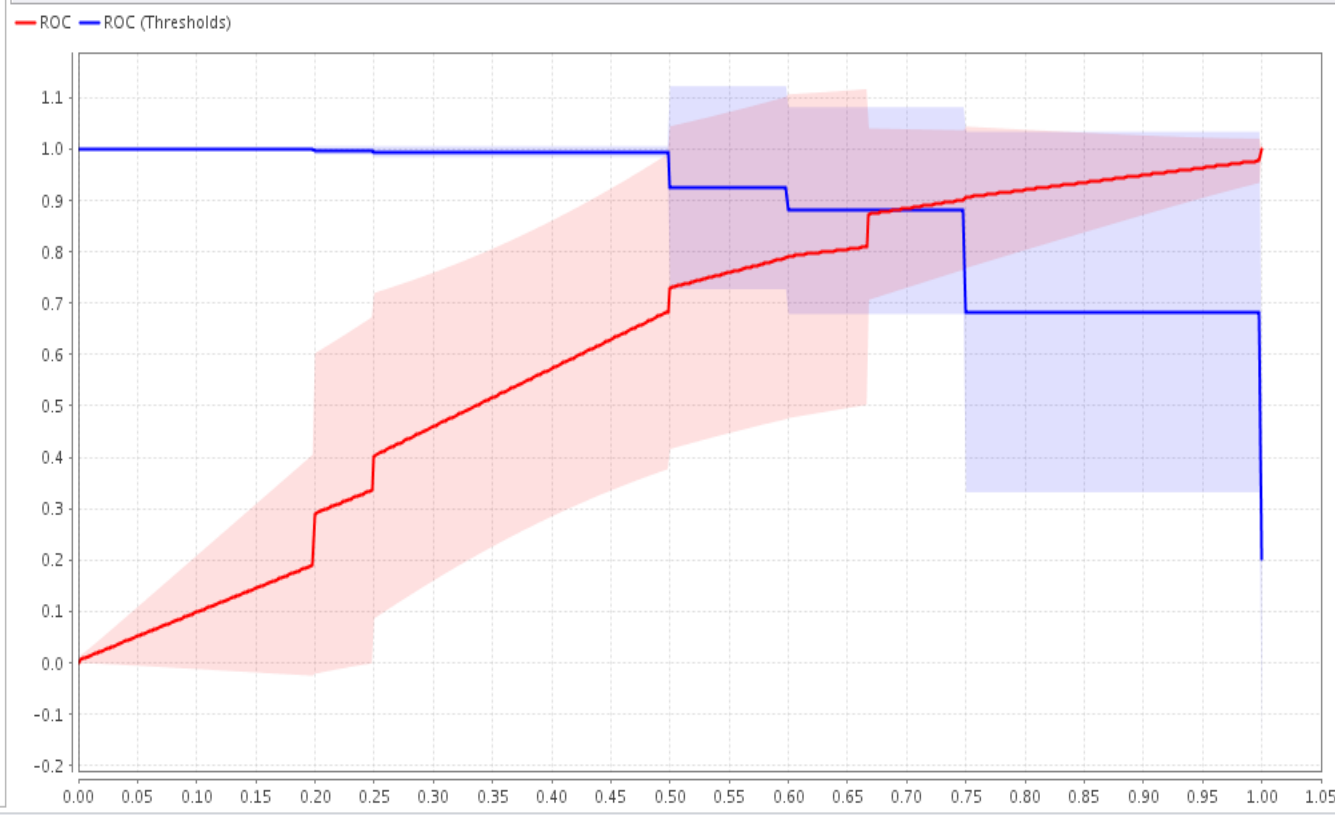
Result Overview x PerformanceVector (Performance (2)) x PerformanceVector (Performance) x Pairwise t-Test (T-Test) x Anova Test (Anova) x

Table / Plot View o Text View o Annotations

- Criterion Selector
- classification_error
- AUC
- false_positive
- false_negative
- true_positive
- true_negative

o Area Under Curve o Text View o Annotations

AUC: 0.610 +/- 0.174 (mikro: 0.610) (positive class: 1.000)



Repositories

- Samples (none)
- DB
- Local Repository (Aler)

Log x x x x

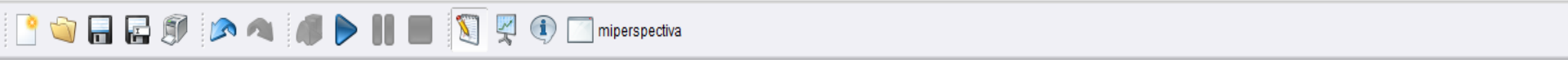
System Monitor

Max: 10 GB
Total: 1024 MB

Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Training
Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Testing
Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Training
Nov 27, 2013 8:24:19 PM INFO: Executing process concurrently: Testing
Nov 27, 2013 8:24:19 PM INFO: Saving results.
Nov 27, 2013 8:24:19 PM INFO: Process //Local Repository/processes/PRUEBAS finished successfully after 2 s

5. ÁRBOL DECISIÓN + CON MATRIZ DE COSTES)

- Reemplazar ambos “binomial performance measurement” por “performance costs”.
- Introducir en ambos una matriz de costes, de tal manera que el coste de confundir la clase 1 por la dos sea de



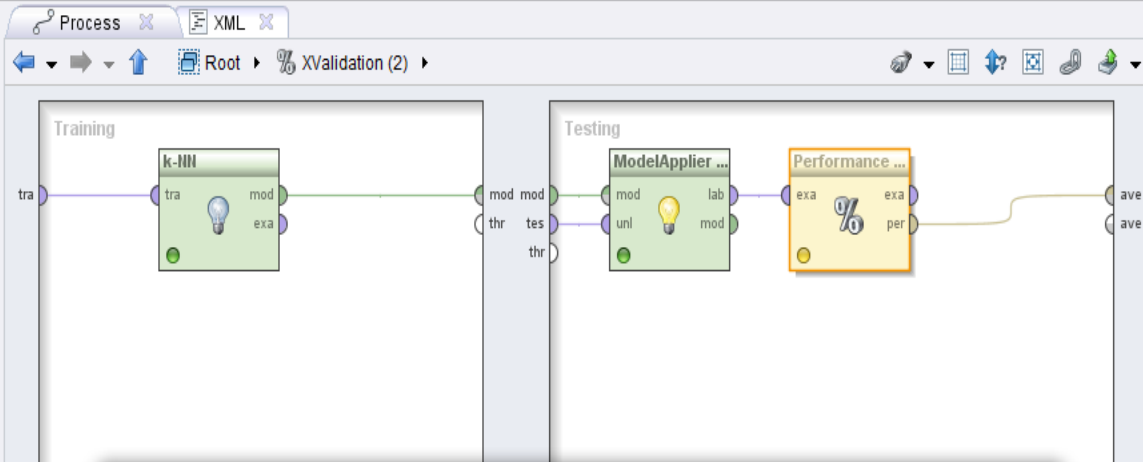
Operators

Search

- Process Control (39)
- Utility (52)
- Repository Access (6)
- Import (28)
- Export (18)
- Data Transformation (114)
- Modeling (262)
- Evaluation (32)
 - X-Validation
 - Split Validation
 - Bootstrapping Validation
 - Batch-X-Validation
 - Wrapper Split Validation
 - Wrapper-X-Validation
 - X-Validation (Parallel)
- Performance Measurement (20)
- Significance (2)
- Visual Evaluation (3)
- ISPR (28)

Repositories

- Samples (none)
- DB
- Local Repository (Aler)
 - data (Aler)
 - oil (Aler - v1, 11/27/13 5:39 PM - 389 kB)
 - oilr (Aler - v1, 11/27/13 6:56 PM - 321 kB)
 - processes (Aler)



Edit Parameter Matrix: cost matrix

Edit Parameter Matrix: **cost matrix**
The matrix of misclassification costs. Columns and Rows in order of internal mapping.

Cost Matrix	True Class 1	True Class 2
Predicted Class 1	0.0	1000.0
Predicted Class 2	25.0	0.0

Buttons: Increase Size, Decrease Size, OK, Cancel

Parameters Context

Performance (2) (Performance (Costs))

cost matrix

class order definition

Performance (Costs)
(RapidMiner Core)

Synopsis

This operator provides the ability to evaluate misclassification costs for performance evaluation of classification tasks.

6. ÁRBOL DE DECISIÓN VS. METACOST CON ÁRBOLES DE DECISIÓN Y MATRIZ DE COSTES

- Vamos a sustituir knn por metacost, que recordemos era capaz de estimar probabilidades y dada una matriz de costes, era capaz de devolver un clasificador óptimo para ella
- Metacost es un clasificador “meta” y trabaja con cualquier clasificador. En este caso utilizaremos como clasificador para metacost a los árboles de decisión

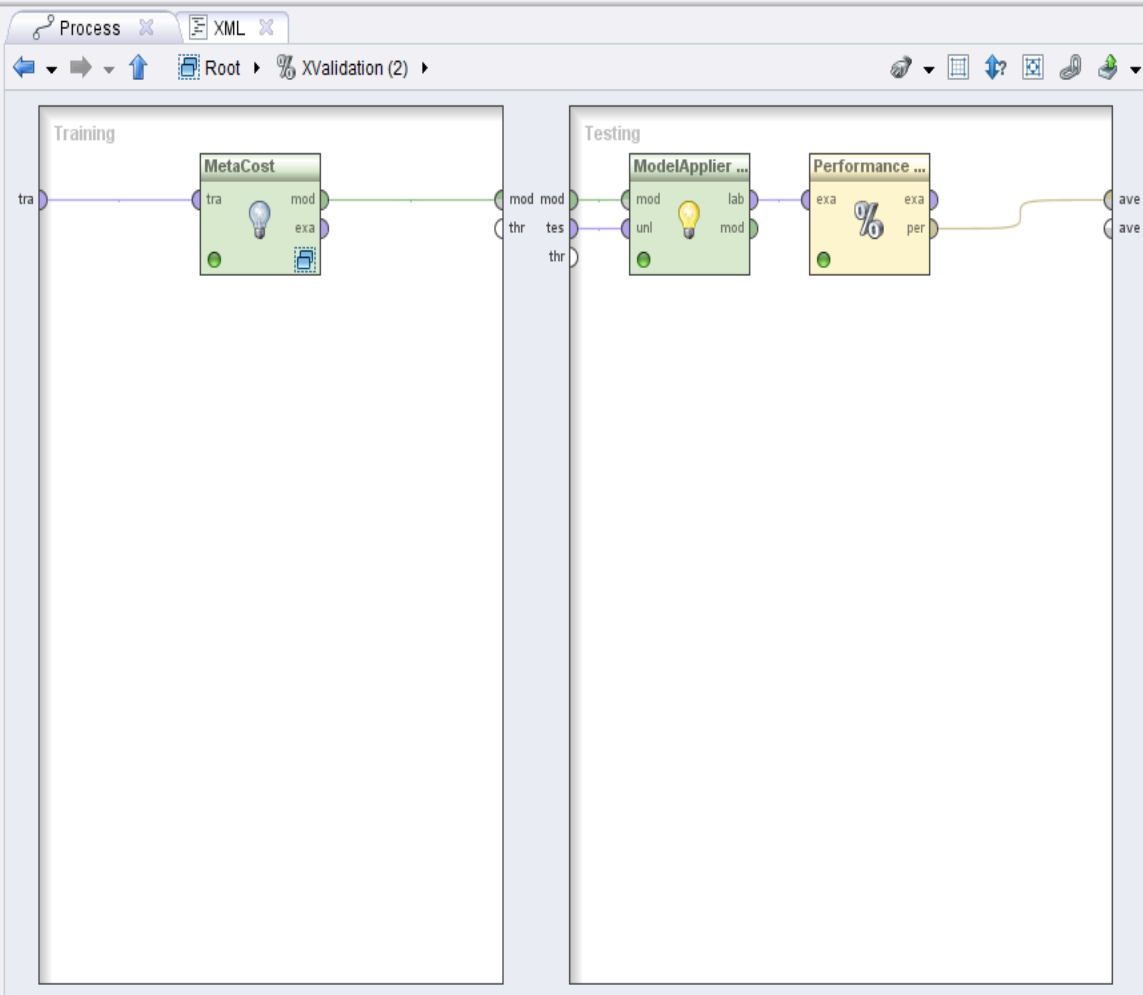
Operators

Search

- Export (18)
- Data Transformation (114)
- Modeling (262)
 - Classification and Regression (168)
 - Lazy Modeling (2)
 - Bayesian Modeling (2)
 - Tree Induction (13)
 - Decision Tree
 - Decision Tree (Multiway)
 - Decision Tree (Weight-Based)
 - ID3
 - CHAID
 - Decision Stump
 - Random Tree
 - Random Forest
 - Decision Tree (Parallel)
 - Decision Stump (Parallel)
 - ID3 (Parallel)
 - ID3 Numerical (Parallel)
 - Decision Tree (Weight-Based)

Repositories

- Samples (none)
- DB
- Local Repository (Aler)
 - data (Aler)
 - oil (Aler - v1, 11/27/13 5:39 PM - 369 kB)
 - oilr (Aler - v1, 11/27/13 6:56 PM - 321 kB)
 - processes (Aler)



Parameters

Context

XValidation (2) (X-Validation)

- average performances only
- leave one out

number of validations: 10

sampling type: stratified sampling

- use local random seed
- parallelize training
- parallelize testing

Compatibility level: 5.1.002

Problems

Log

One potential problem

Message	Fixes	Location
Parameter 'repository entry' accesses a repository by name (//Local Repo...	No quick fix available	Retrieve

Help

Comment

X-Validation (RapidMiner Core)

Synopsis

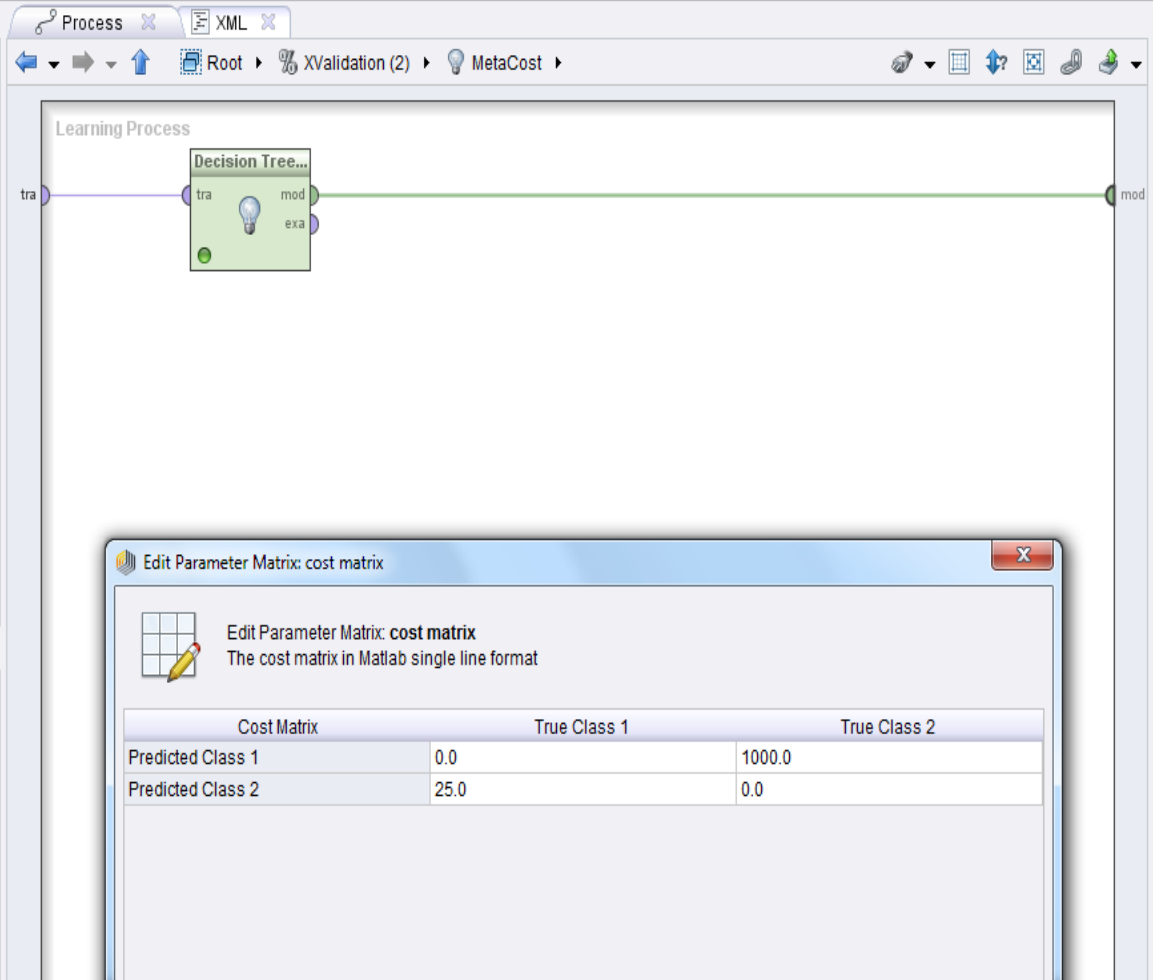
This operator performs a cross-validation in order to estimate the statistical performance of a learning operator (usually on unseen data sets). It is mainly used to estimate how accurately

Operators

- Export (18)
- Data Transformation (114)
- Modeling (262)
 - Classification and Regression (168)
 - Lazy Modeling (2)
 - Bayesian Modeling (2)
 - Tree Induction (13)
 - Decision Tree
 - Decision Tree (Multiway)
 - Decision Tree (Weight-Based)
 - ID3
 - CHAID
 - Decision Stump
 - Random Tree
 - Random Forest
 - Decision Tree (Parallel)
 - Decision Stump (Parallel)
 - ID3 (Parallel)
 - ID3 Numerical (Parallel)
 - Decision Tree (Weight-Based)

Repositories

- Samples (none)
- DB
- Local Repository (Aler)
 - data (Aler)
 - oil (Aler - v1, 11/27/13 5:39 PM - 369 kB)
 - oilr (Aler - v1, 11/27/13 6:56 PM - 321 kB)
 - processes (Aler)



Parameters Context

MetaCost

cost matrix

use subset for training

iterations

sampling with replacement

use local random seed

parallelize learning process

Edit Parameter Matrix: cost matrix

The cost matrix in Matlab single line format

Cost Matrix	True Class 1	True Class 2
Predicted Class 1	0.0	1000.0
Predicted Class 2	25.0	0.0

Help Comment

MetaCost (RapidMiner Core)

Synopsis

This metaclassifier makes its base classifier cost-sensitive by using the given cost matrix to compute label predictions according to classification costs.

Description

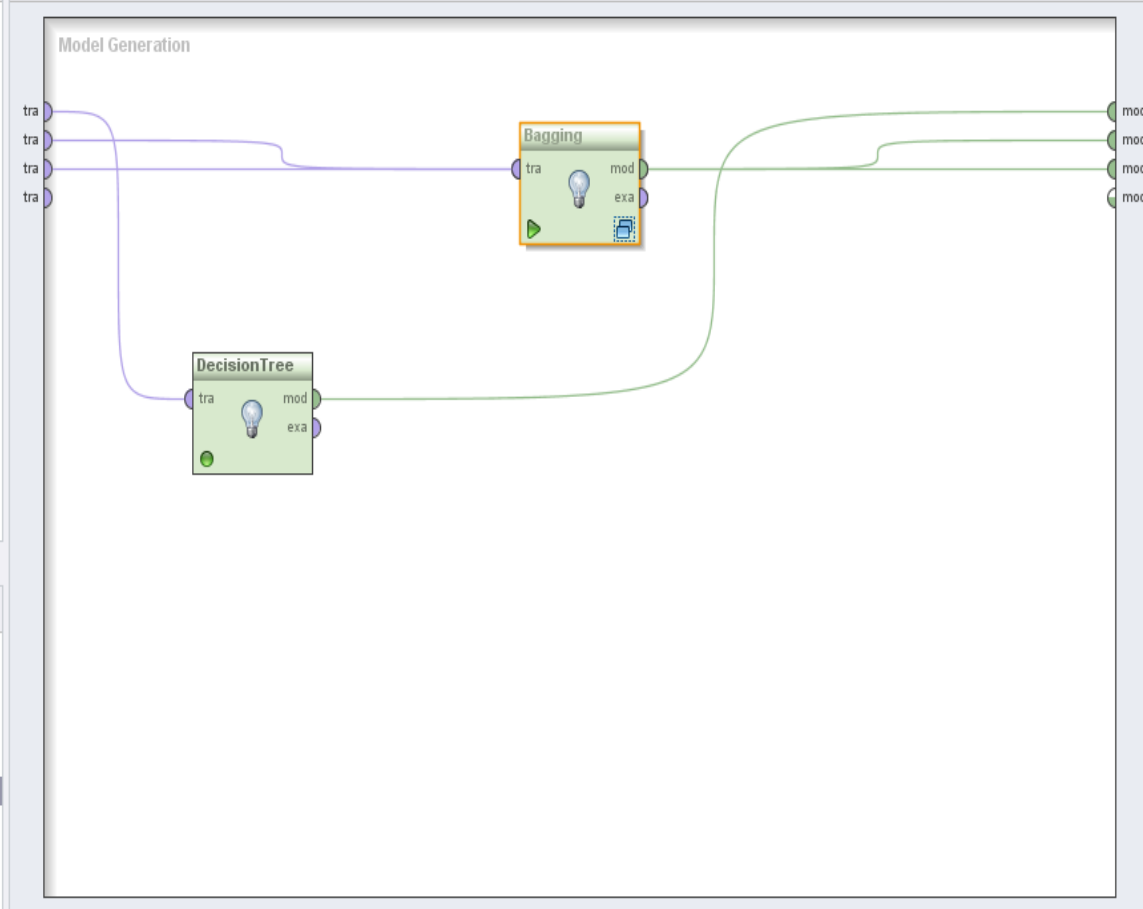
7. COMPARACIÓN DE CURVAS ROC

- Usaremos la File -> Template “compare ROCs”
- Incluiremos decision tree y bagging con decision trees

Modeling (1)
Classification and Regression (1)
Tree Induction (1)
Random Forest

Repositories

- Samples (none)
- DB
- Local Repository (Aler)
 - data (Aler)
 - oil (Aler - v1, 11/27/13 5:39 PM - 369 kB)
 - oilr (Aler - v1, 11/27/13 6:56 PM - 321 kB)
 - processes (Aler)



Bagging

sample ratio: 0.9

iterations: 10

average confidences

use local random seed

parallelize learning process

Help Comment

Bagging (RapidMiner Core)

Synopsis

Bootstrap aggregating (bagging) is a machine learning ensemble meta-algorithm to improve classification and regression models in terms of stability and classification accuracy. It also reduces

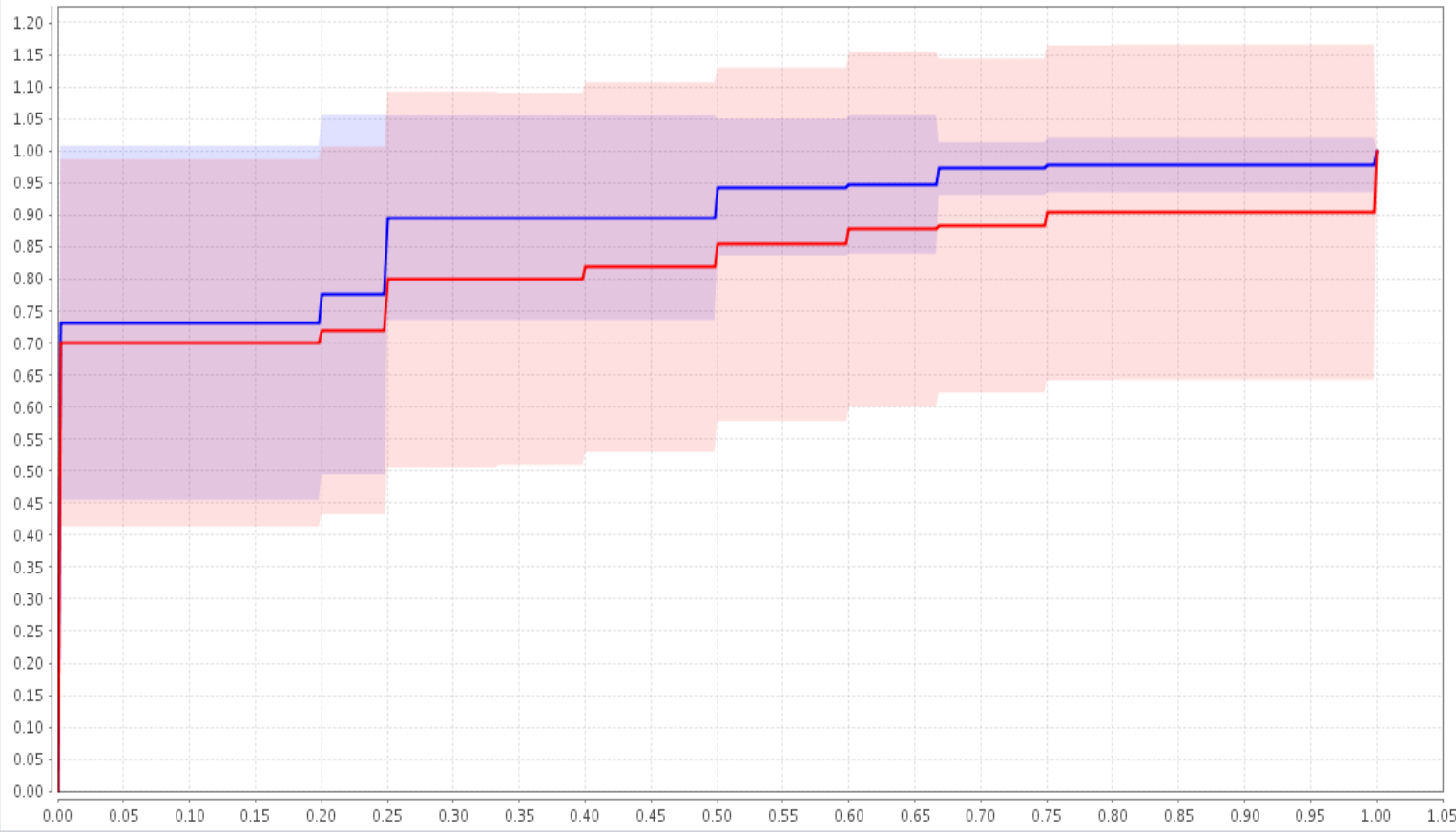
Problems Log

One potential problem

Message	Fixes	Location
Parameter 'repository entry' accesses a repository by name (//Local Repo...	No quick fix available	Retrieve

ROC Comparison Annotations

Bagging DecisionTree



Repositories

- Samples (none)
- DB
- Local Repository (Aler)
 - data (Aler)
 - oil (Aler - v1, 11/27/13 5:39 PM - 369 kB)
 - oilr (Aler - v1, 11/27/13 6:56 PM - 321 kB)
 - processes (Aler)

Log

Nov 27, 2013 9:26:39 PM INFO: Executing process concurrently: Model Generation
Nov 27, 2013 9:26:41 PM INFO: Executing process concurrently: Model Generation
Nov 27, 2013 9:26:44 PM INFO: Executing process concurrently: Model Generation
Nov 27, 2013 9:26:47 PM INFO: Saving results.
Nov 27, 2013 9:26:47 PM INFO: Process finished successfully after 24 s

System Monitor

