

Subject

ROC graphs (Receiver Operating Characteristics) with TANAGRA.

ROC graphs enable to compare two or more supervised learning algorithms, they have properties that make them especially useful for domains with skewed class distribution and unequal classification error costs.

An ROC graph depicts relative trade-offs between true positives rate and false positives rate. It needs continuous output of classifier, an estimate of an instance's class membership probabilities. In fact, a "score", a numeric value that represents the degree to which an instance is a member of a class is sufficient.

AUC (Area Under Curve) reduces ROC performances to a single scalar value, which enables to compare several classifiers: this area is equivalent to the probability that the classifier will rank a randomly chosen positive instance higher than a randomly chosen negative instance.

In this tutorial, we compare linear discriminant analysis (LDA) and support vector machine (SVM) on a heart-diseases detection problem.

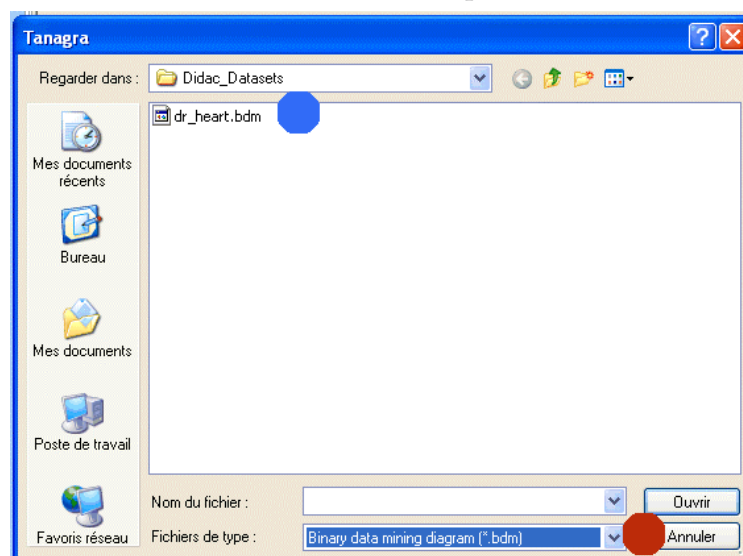
Dataset

HEART dataset; detect heart-disease from patient's characteristics.

ROC graphs

Download the dataset

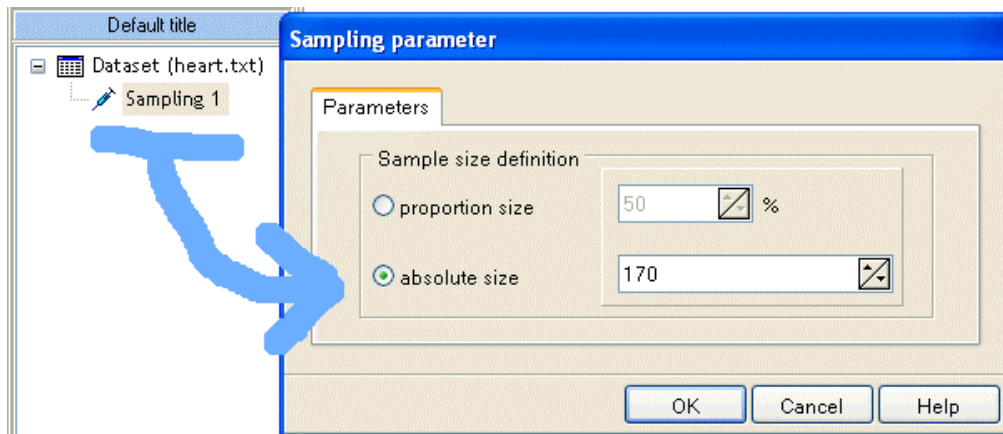
Open the DR_HEART.BDM dataset, click on File / Open menu.



Define training and test set

In order to obtain an unbiased estimate of the performances of the classifiers, we will divide the whole dataset into two parts: 170 instance as a training set; 100 instances for the evaluation of the classifiers.

Add the SAMPLING component in the stream diagram, set the following parameters.



Feature construction

LDA and SVM cannot handle discrete attributes, we must transform them.

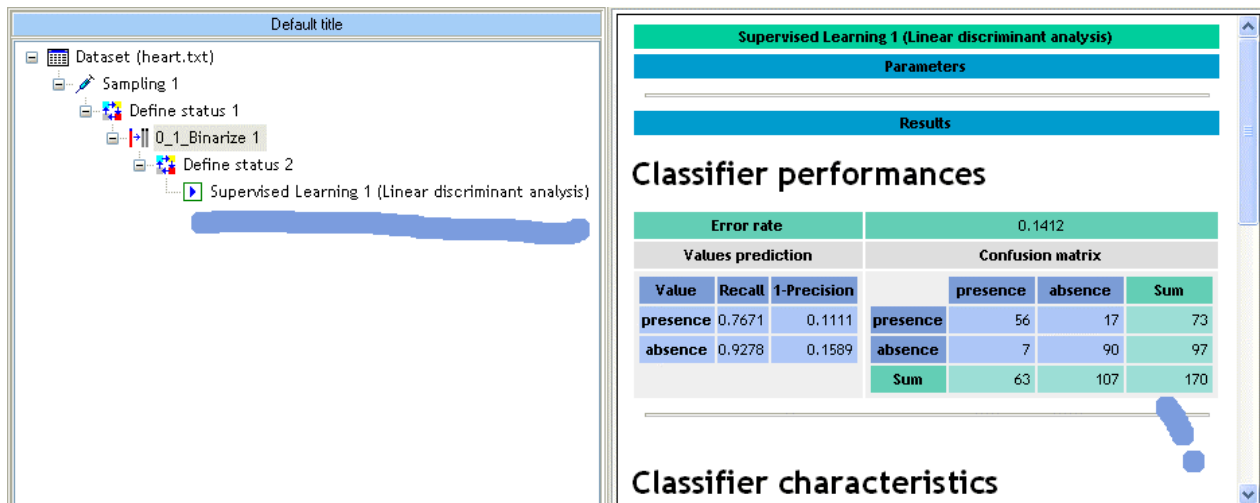
Add the 0_1_BINARIZE component after having selected all discrete attributes (except the class attribute COEUR).

Attribute binarization	
Source att	New attributes
sexe	(sexe_masculin_1)
type_douleur	(type_douleur_D_1,type_douleur_C_1,type_douleur_B_1)
sucre	(sucre_A_1)
electro	(electro_C_1,electro_A_1)
engine	(engine_non_1)
vaisseau	(vaisseau_D_1,vaisseau_A_1,vaisseau_B_1)

LDA

We can run the LDA learning algorithm now. Set as INPUT all continuous attributes and as TARGET the COEUR attribute, add the SUPERVISED LEARNING component.

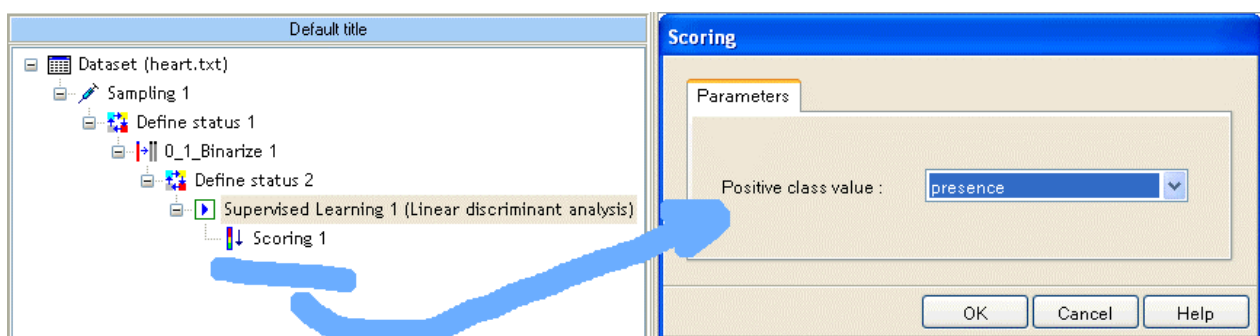
We obtain the following results; we see that we have built the prediction model on 170 instances.



SCORING for LDA

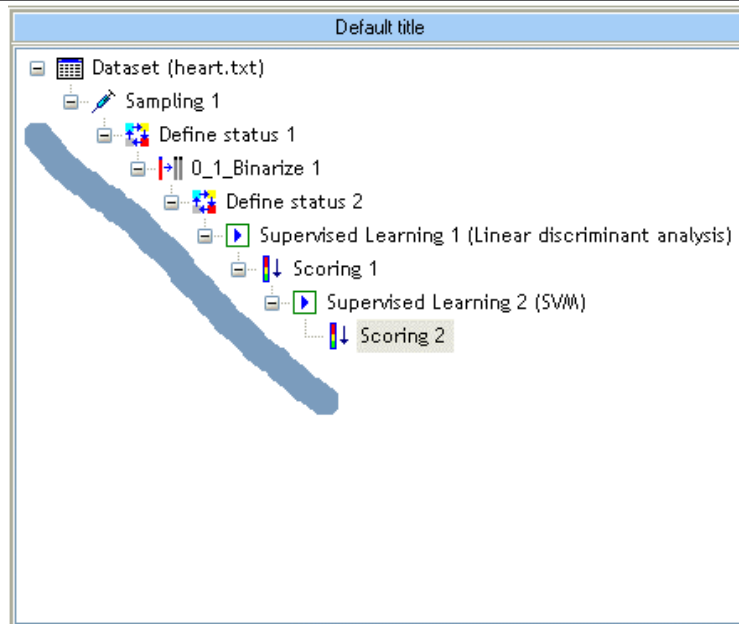
We must now compute the "score" of each instance for the positive class value on the whole dataset.

Add the SCORING component, and set as positive the "PRESENCE" value.



SCORING for SVM

Add also the SVM component in the stream diagram and build the "score" value of each instance.



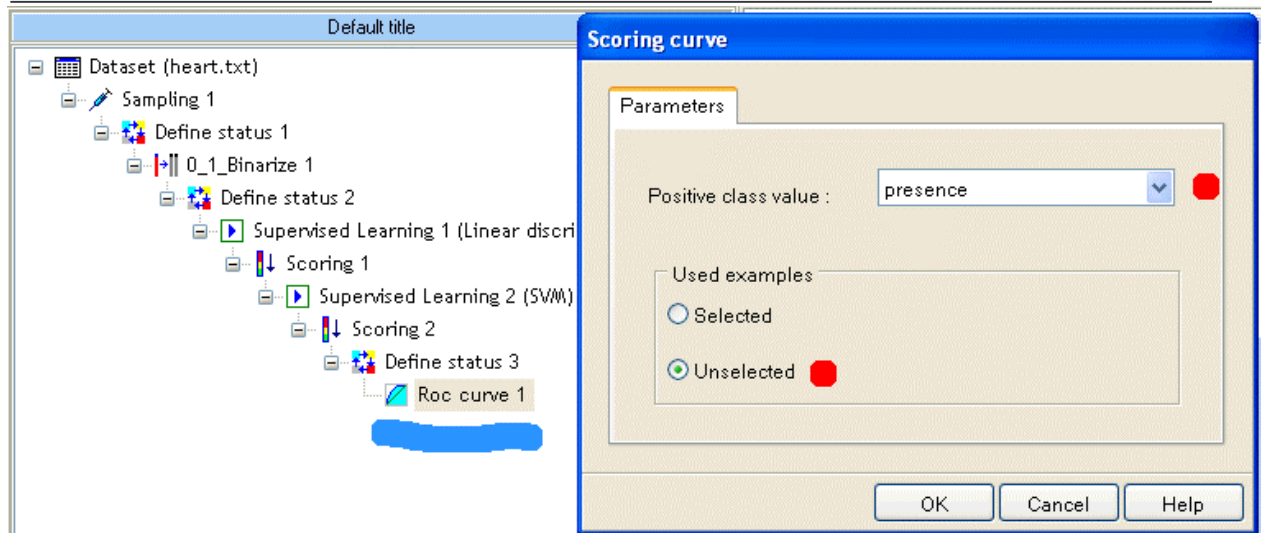
ROC graphs table

In order to build the ROC graphs, we must specify the TARGET attribute and the SCORE attributes.

Add the DEFINE STATUS component in the diagram, set COEUR as TARGET, SCORE_1 (from LDA) and SCORE_2 (from SVM) as INPUT. We note that we can add other score attributes, for instance a score which is provided by a domain expert.

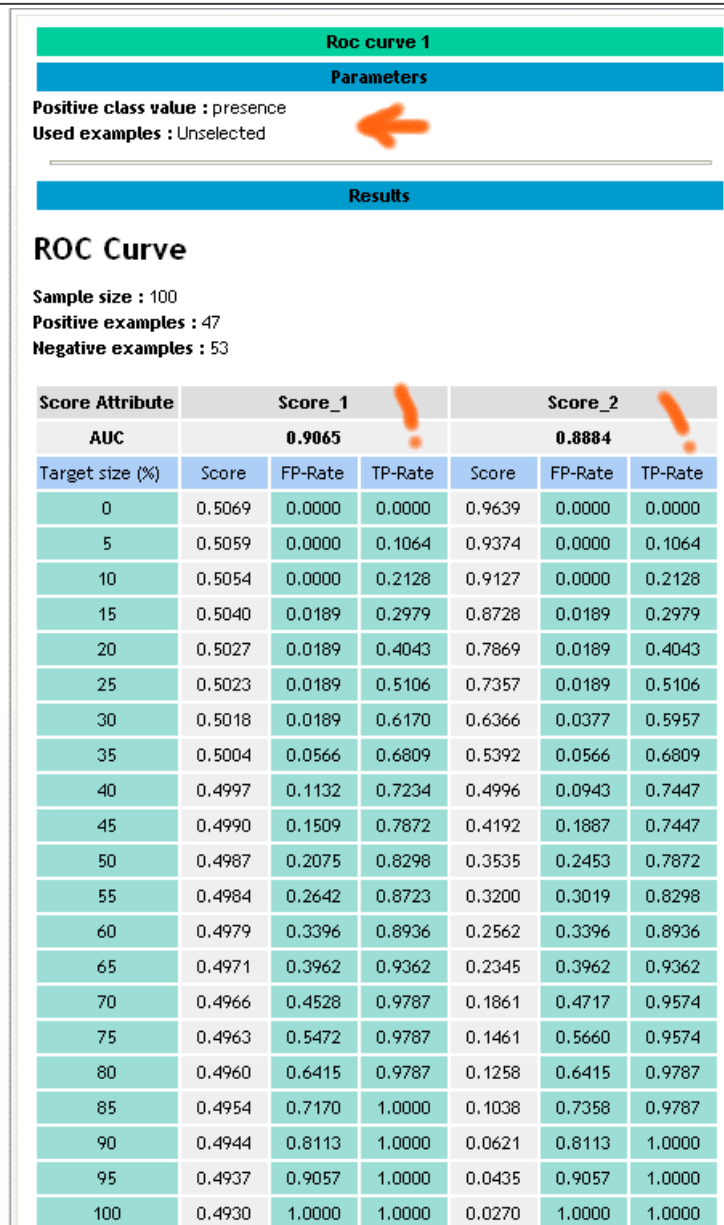
vaisseau				
coeur	yes	-	-	●
sexe_masculin_1	-	-	-	
type_douleur_D_1	-	-	-	
type_douleur_C_1	-	-	-	
type_douleur_B_1	-	-	-	
sucrer_A_1	-	-	-	
electro_C_1	-	-	-	
electro_A_1	-	-	-	
angine_non_1	-	-	-	
vaisseau_D_1	-	-	-	
vaisseau_A_1	-	-	-	
vaisseau_B_1	-	-	-	
pred_SpvInstance_1	-	-	-	
Score_1	-	yes	-	●
pred_SpvInstance_2	-	-	-	
Score_2	-	yes	-	●

Set the right parameter of the ROC component: the positive class value; the ROC graphs must be computed on the test set.



We obtain the following results:

- AUC.
- For each target size (True Positive + False Positive), we have the true positive rate and the false positive rate.
- Generally, score values are not comparable for different methods.



For the heart disease problem, we see that LDA and SVM have similar performances (because we have randomly selected the training and test set, you obtain a little different results).

ROC graphs with a spreadsheet

In order to draw the ROC graph, you can copy the results in a spreadsheet and build a scatter plot.

Click on the COMPONENT / COPY RESULTS menu and paste the grid in a spreadsheet, to build the graph is easy.

