

Centro Brasileiro de Pesquisas Físicas



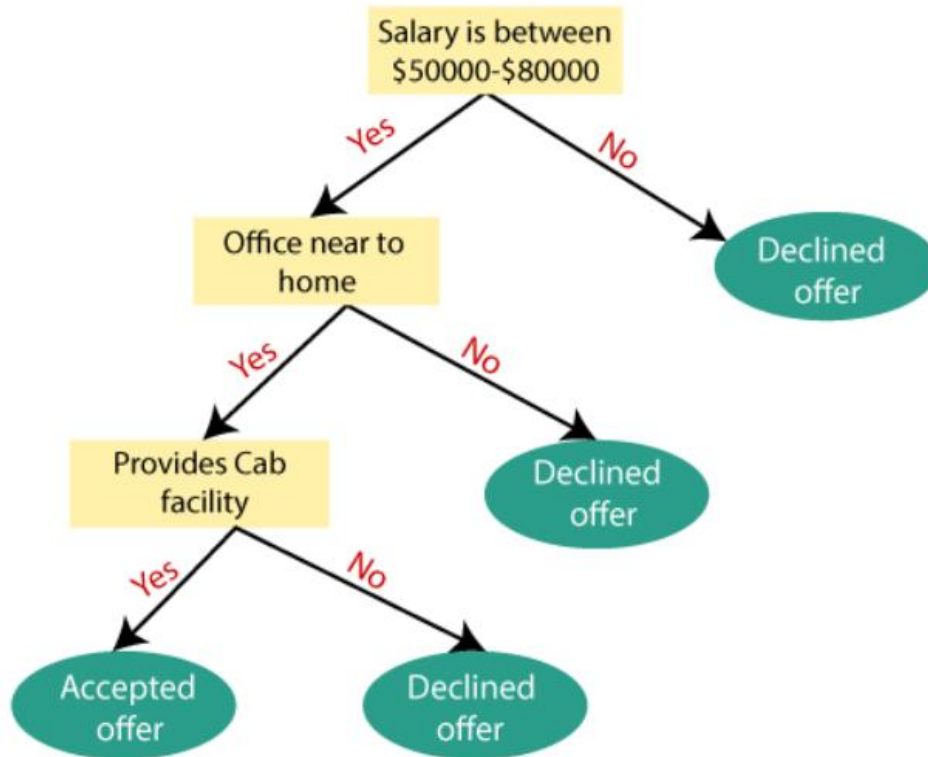
Métodos para Análise de grande volume de dados e Astroinformática

Clécio Roque De Bom – debom@cbpf.br

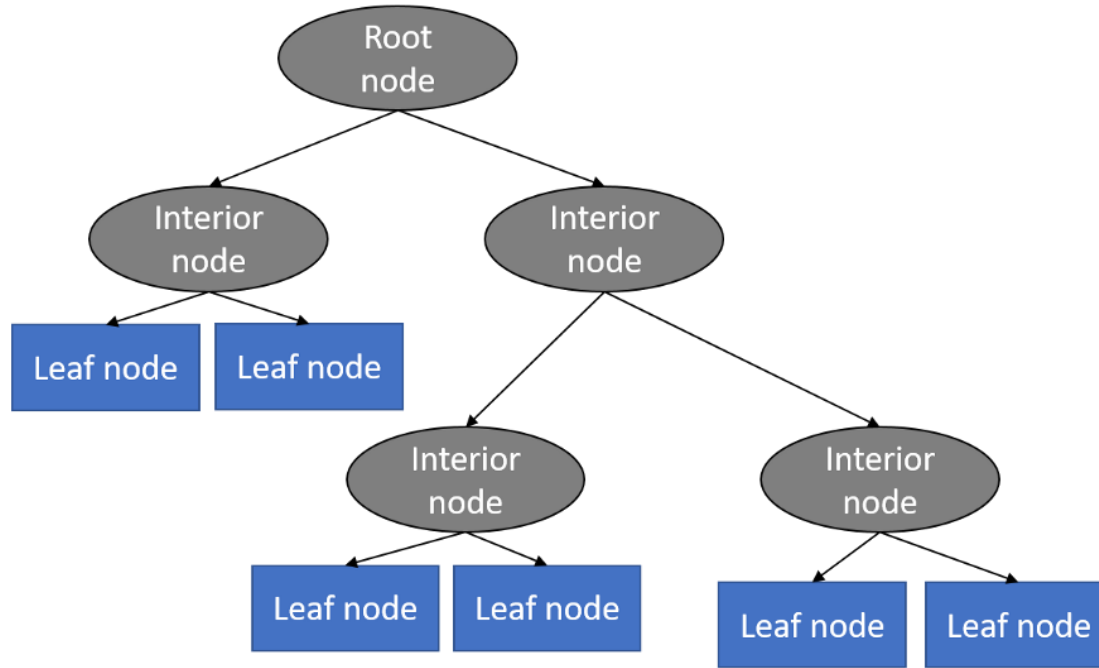
clearnightsrthebest.com



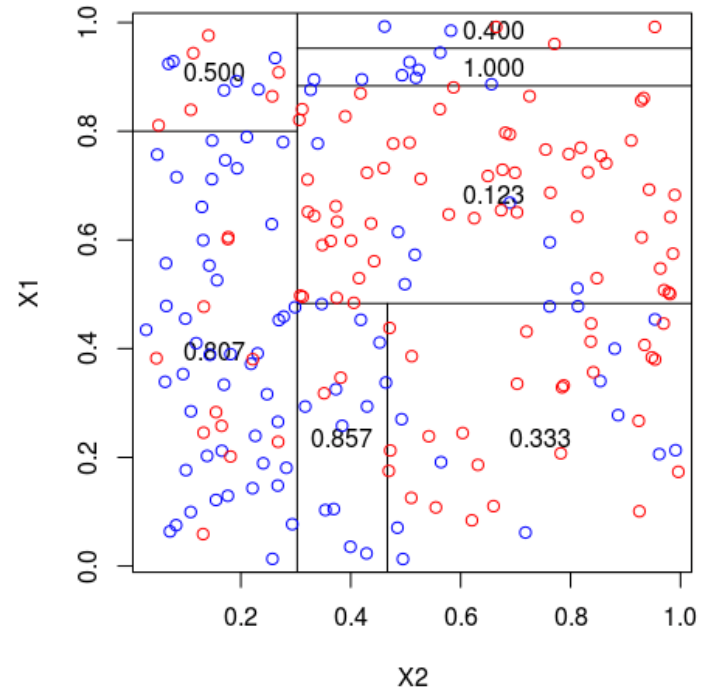
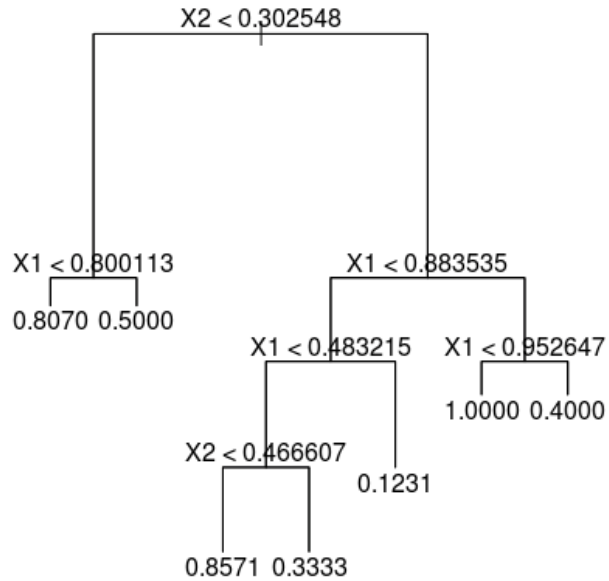
Decision Trees



Decision Trees




Decision Trees



Decision Trees

	Weight	Heart Disease
Lowest	155	No
	180	Yes
	190	No
	220	Yes
Highest	225	Yes



Decision Trees

This is the best value of your metric, This is going to be the Threshold

Weight	Heart Disease
155	No
167.5	
180	Yes
185	
190	No
205	
220	Yes
222.5	
225	Yes

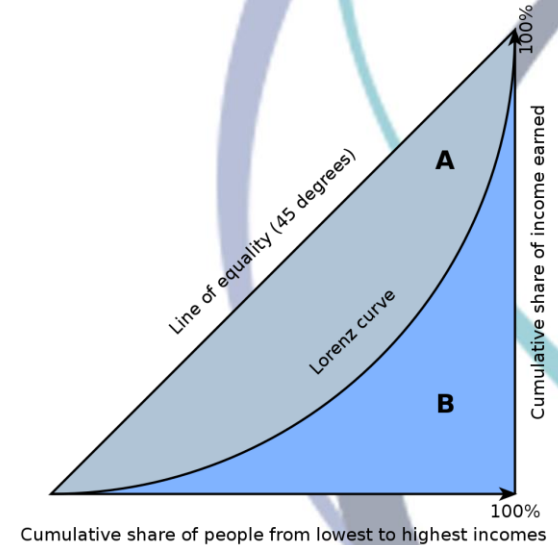
Gini impurity = 0.3

Gini impurity = 0.47

Gini impurity = 0.27

Gini impurity = 0.4

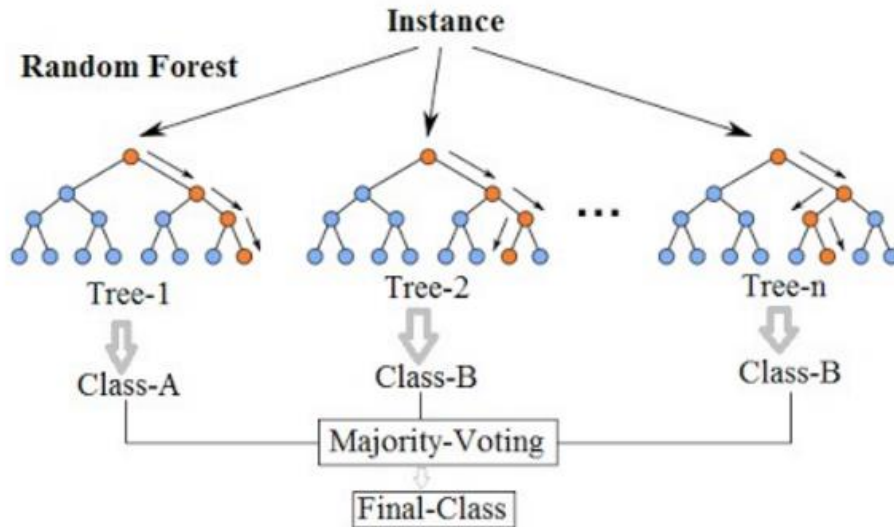
$$\text{Gini impurity} = 1 - \text{Gini}$$
$$G = 1 - \sum_{k=0}^{k=n-1} (X_{k+1} - X_k)(Y_{k+1} + Y_k)$$



Random Forest

Decision trees overfits... Decision trees are very sensible to the training set... So Random Forest

Random Forest Simplified



Now you need a set of (wealy/um) correlated trees.

How you do it?

Bagging the samples.
Bagging the features.

Centro Brasileiro de Pesquisas Físicas



Métodos para Análise de grande volume de dados e Astroinformática

Clécio Roque De Bom – debom@cbpf.br

clearnightsrthebest.com

