**School of Built Environment, Engineering and Computing**

**Intelligent Systems & Machine Learning**
**Dr. Gopal Jamnal**
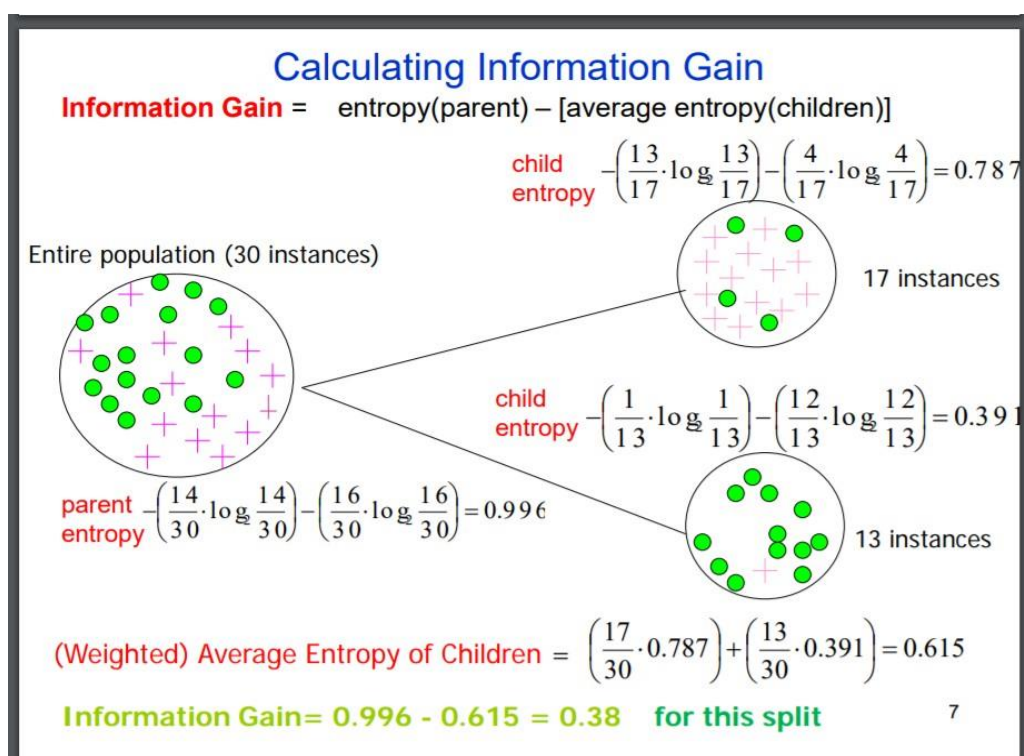**Week 3 – Decision Trees and Random Forest**

In this section, we'll be working on tree-based models such as decision tree and random forest. Open rapid miner and import example dataset to perform lab tasks.

**Supporting Tasks**
Read the RapidMiner Documentation on Decision Trees
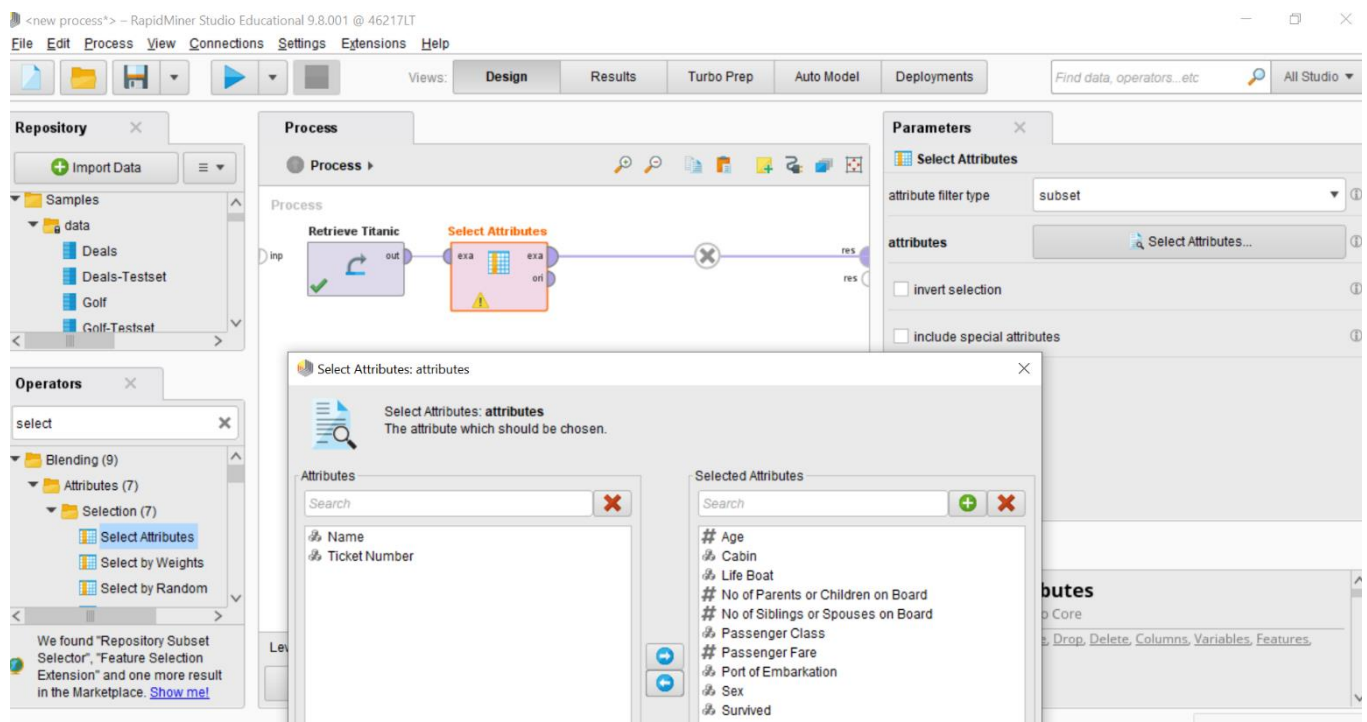
**Information Gain**
Given the following dataset, calculate the Information Gain of the first node split.



## Calculating Information Gain

**Information Gain** = entropy(parent) − [average entropy(children)]

Entire population (30 instances)

child entropy $-\left(\dfrac{13}{17}\cdot\log\dfrac{13}{17}\right)-\left(\dfrac{4}{17}\cdot\log\dfrac{4}{17}\right)=0.787$

17 instances

child entropy $-\left(\dfrac{1}{13}\cdot\log\dfrac{1}{13}\right)-\left(\dfrac{12}{13}\cdot\log\dfrac{12}{13}\right)=0.391$

13 instances

parent entropy $-\left(\dfrac{14}{30}\cdot\log\dfrac{14}{30}\right)-\left(\dfrac{16}{30}\cdot\log\dfrac{16}{30}\right)=0.996$

(Weighted) Average Entropy of Children $=\left(\dfrac{17}{30}\cdot0.787\right)+\left(\dfrac{13}{30}\cdot0.391\right)=0.615$

**Information Gain= 0.996 - 0.615 = 0.38** for this split
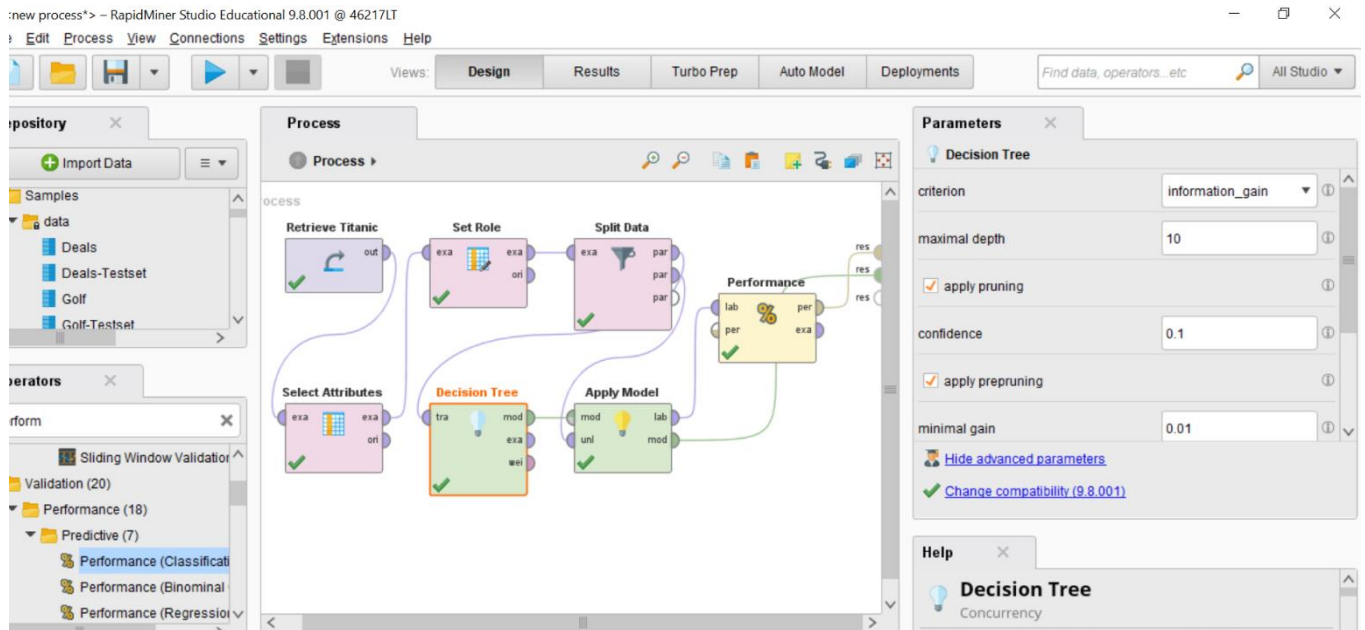
7

## Decision Tree

**Step 1:** first import the sample dataset and use select attributes operator to filter data. As you can see, Name and Ticket Number are filtered out from selected attributes.



**Step 2**: For the target value, you'll set role to predict survived attribute. Once set role is defined, you'll need to split data into train and test set.

**Task:** you need to perform the task you did in week 2 tutorial - set role and split data. You can retrospectively revise the concept for these two operator by checking week 2 tutorial sheet.

**Step 3**: You'll be using decision tree as first classifier with parameters for split criterion>info_gain, pre-post pruning and maximum_depth = 10; as shown in the image below;
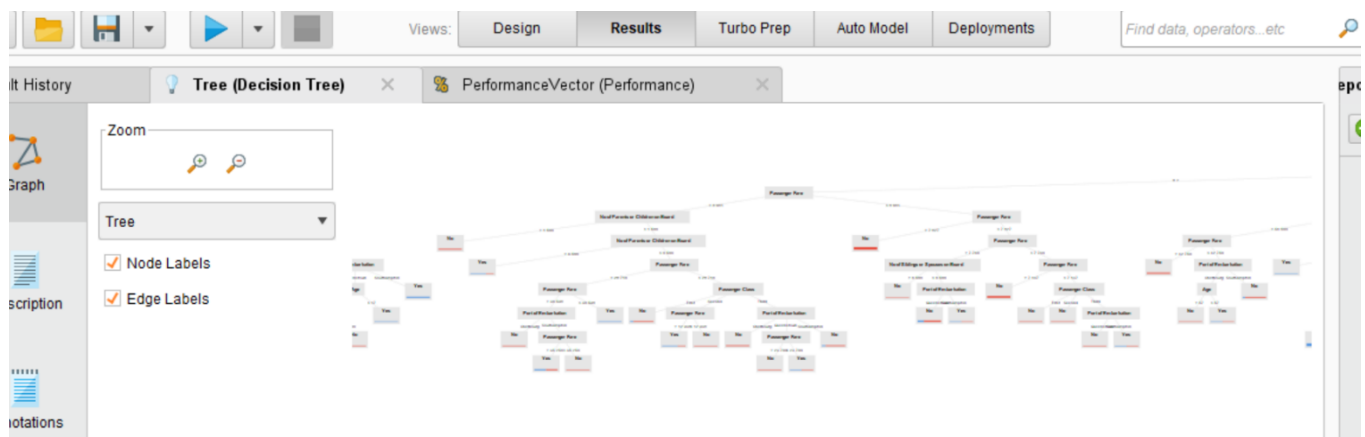
**Step 4:** Performance Vector analysis – you will be using performance(classification) to generate confusion matrix and accuracy score.



accuracy: 76.08%

| | true Yes | true No | class precision |
|---|---|---|---|
| pred. Yes | 95 | 39 | 70.90% |
| pred. No | 55 | 204 | 78.76% |
| class recall | 63.33% | 83.95% | |

63.33%

**Random Forest**

**Note:** same step for 1,2 and 4 will be used from random forest

Step 1: follow the same previous step from decision tree.

Step 2: follow the same previous step from decision tree.

**Step 3:** In this step, you'll be using random forest operator and would set parameters for number of tree=100, split criterion=info_gain, pre-post pruning abd maximum_depth=10 as shown in image below;

follow the same previous step from decision tree.



**Step 4:** Analyze confusion matrix and accuracy score. Is it better than decision tree ?

accuracy: 85.75%

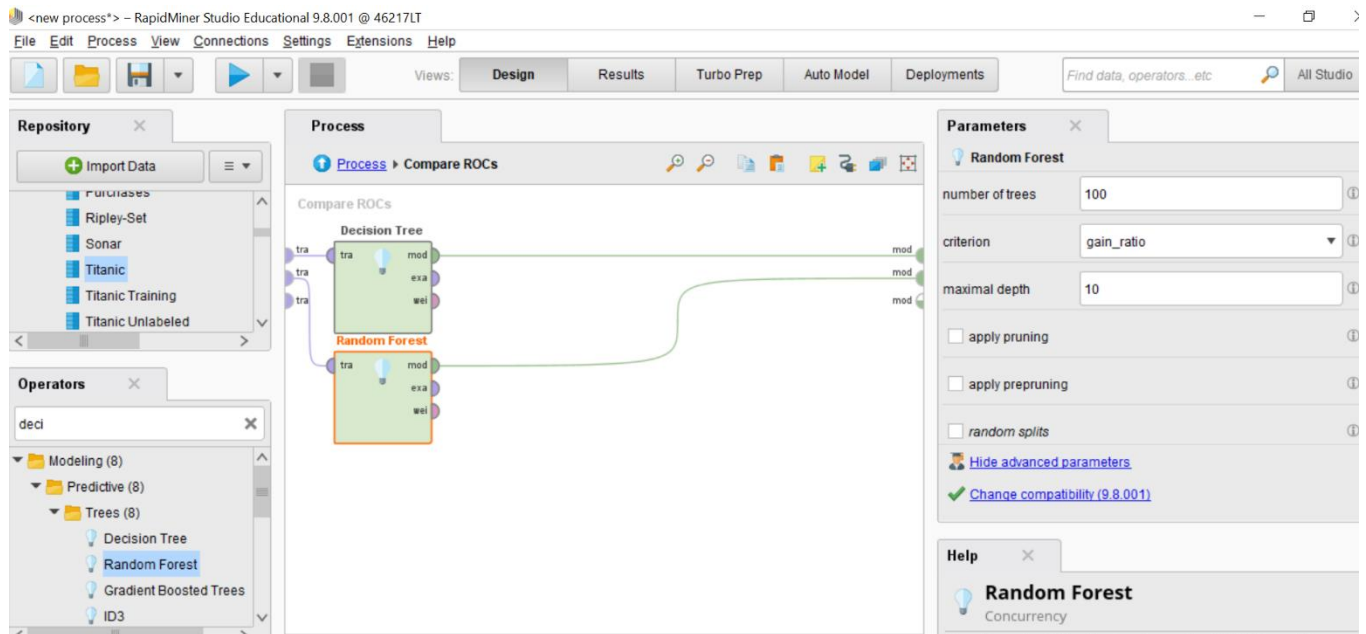|  | true Yes | true No | class precision |
|---|---|---|---|
| pred. Yes | 99 | 5 | 95.19% |
| pred. No | 51 | 238 | 82.35% |
| class recall | 66.00% | 97.94% |  |

**Compare Decision Tree and Random Forest Performance with ROC**

In the tutorial, we'll be using two different models named, decision tree and random forest. For the performance evaluation of both models, we'll compare roc curves. ROC curve is a graphical plot of the sensitivity, or true positive rate, vs. false positive rate (one minus the specificity or true negative rate), for a binary classifier system as its discrimination threshold is varied. The ROC can also be represented equivalently by plotting the fraction of true positives out of the positives (TPR = true positive rate) vs. the fraction of false positives out of the negatives (FPR = false positive rate). ROC curves are calculated by first ordering the classified examples by confidence. Afterwards all the examples are taken into account with decreasing confidence to plot the false positive rate on the x-axis and the true positive rate on the y-axis.
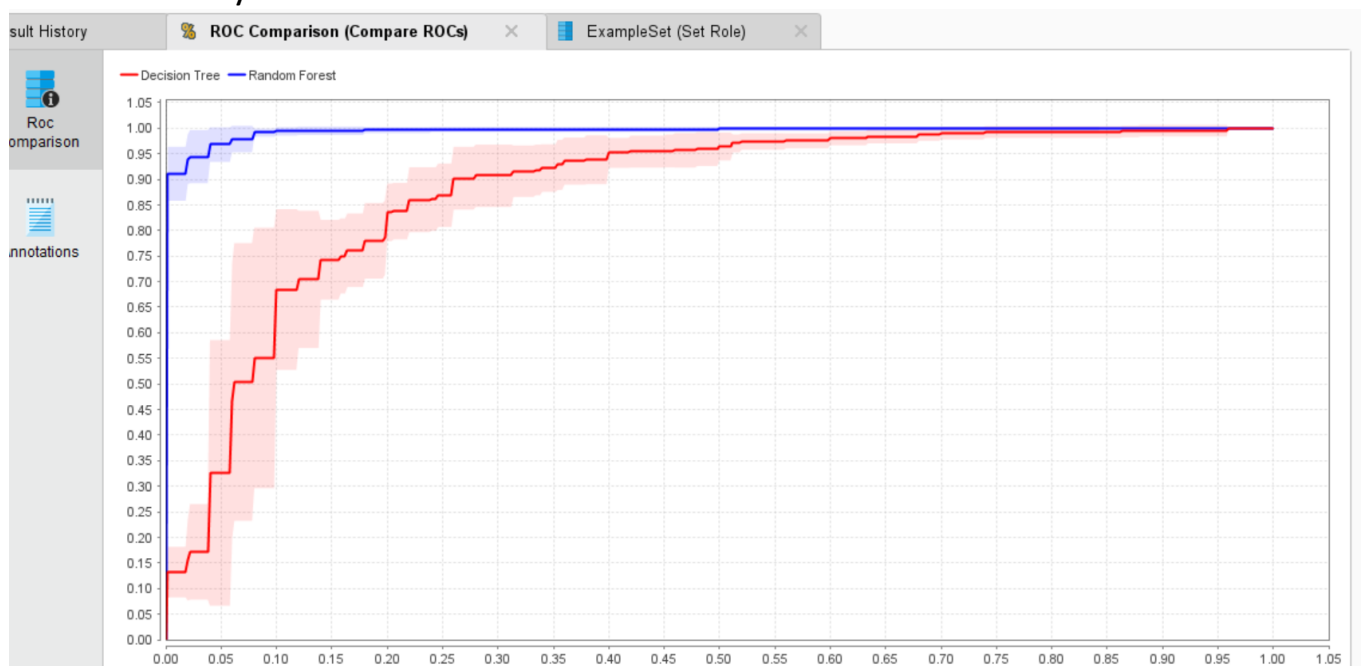
Step 3: you'll be using compare ROCs operator for decision tree and random forest comparison.



Inside ROC, you will be using two classifiers as shown below;

Compare ROC curve and analyze which one is better to predict true positive(tp) cases correctly?



**Exercise:** for the practice and gain deeper understanding, look for other sample dataset and try to interpret their performance with confusion matrix and roc curve for true positive and false positive rates.