

# BRIEF REPORT: DECISION TREE, NAÏVE BAYES CLASSIFIER

## Motivation:

- To know about RapidMiner tool
- Get used to different operators as well as some classifiers like; Naive bayes, Decision tree and Bayes classifier.
- Learn about model training, Pre-processing, tokenization, punctuation.

**Dataset for TASK- I :** IMDB ratings and Movie review on more than 49000 movies.

The tasks we have performed as can be seen from the images shown below.

## 1. **Data Pre-processing :**

- **Pre-processing:** It refers to the preparing (cleansing and organizing) the raw data to make it suitable for building and training a ML model.
- **Tokenization:** It is a process to split a string or text into a list of tokens.
- **Punctuation Removing:** In this process we replace certain words.

for example

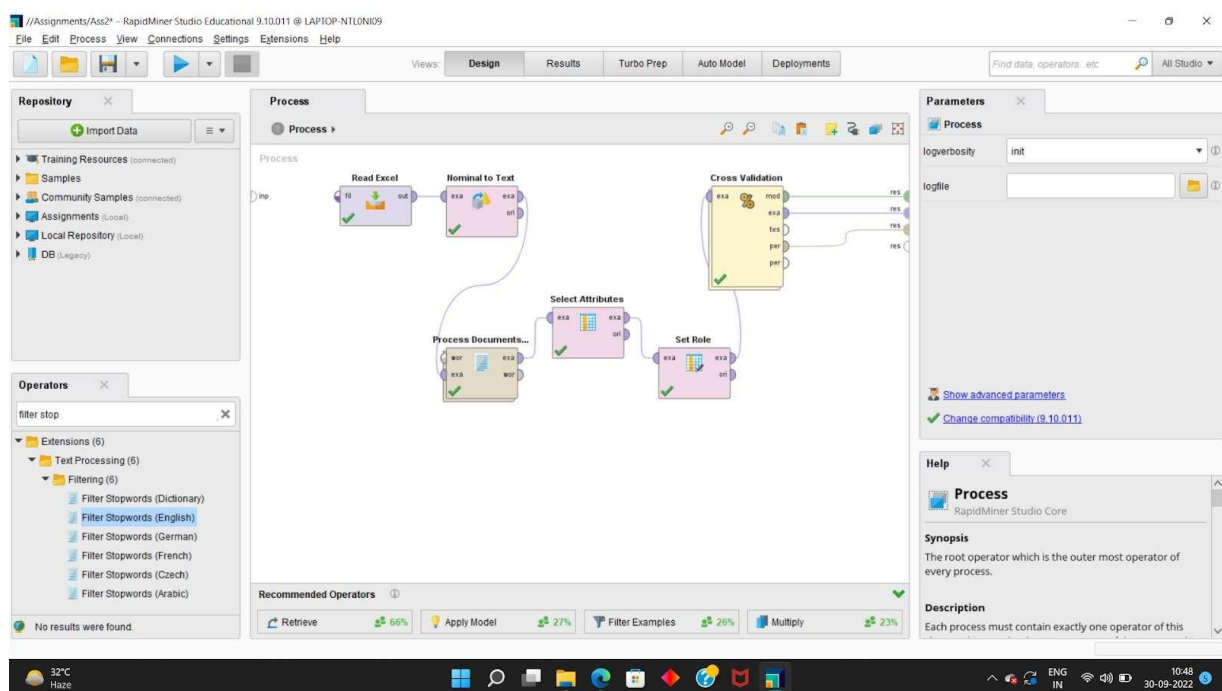
I'm = I am

Let's = Let us

I've = I Have

- **Frequency Vectorization** : In ML vectorization is a step to extract the feature. The idea is to get some distinct feature out of the text for the model to train on by converting text to numerical vectors.
- **Decision Tree Vs Naive bayes classifier:**

Decision tree is a discriminative model whereas naive bayes is a generative model. Decision tree is more easy and flexible. In decision tree pruning may neglect some values and can lead the accuracy for a toss.

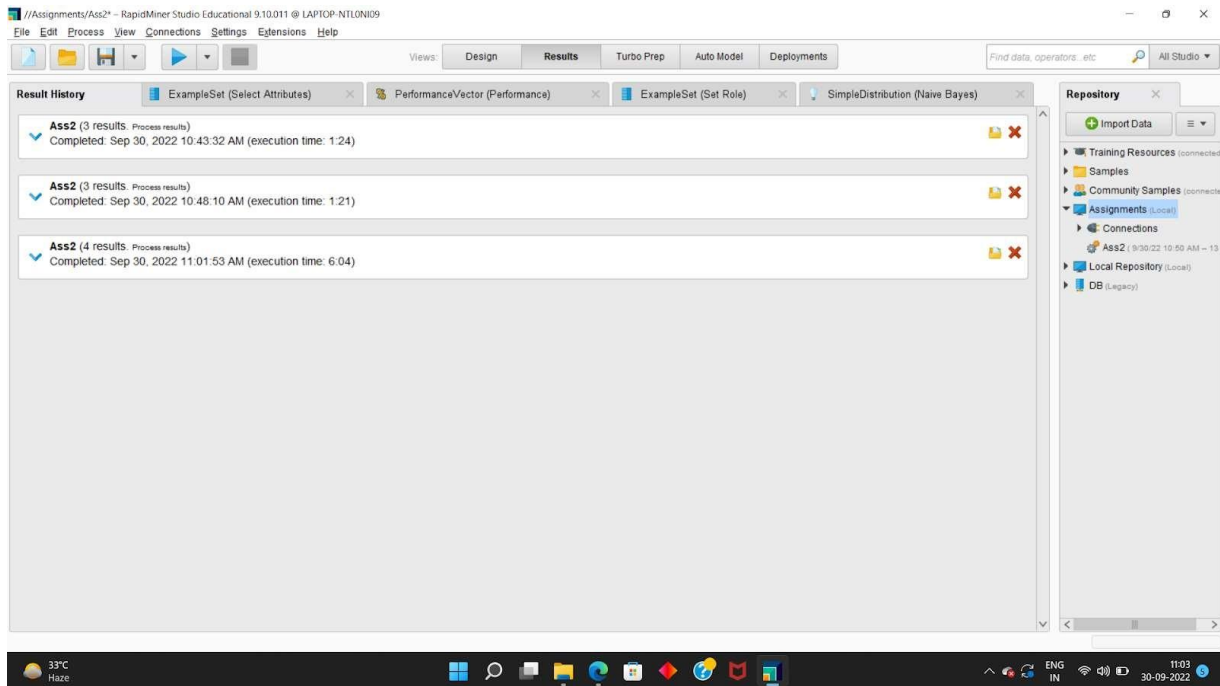


In the following task we reduced the number of examples because we have face some error regarding the entire data analysis which was approximately 49500. So we performed the task over 5000 examples.

It helped us to reduce the processing time and to resolve the error.

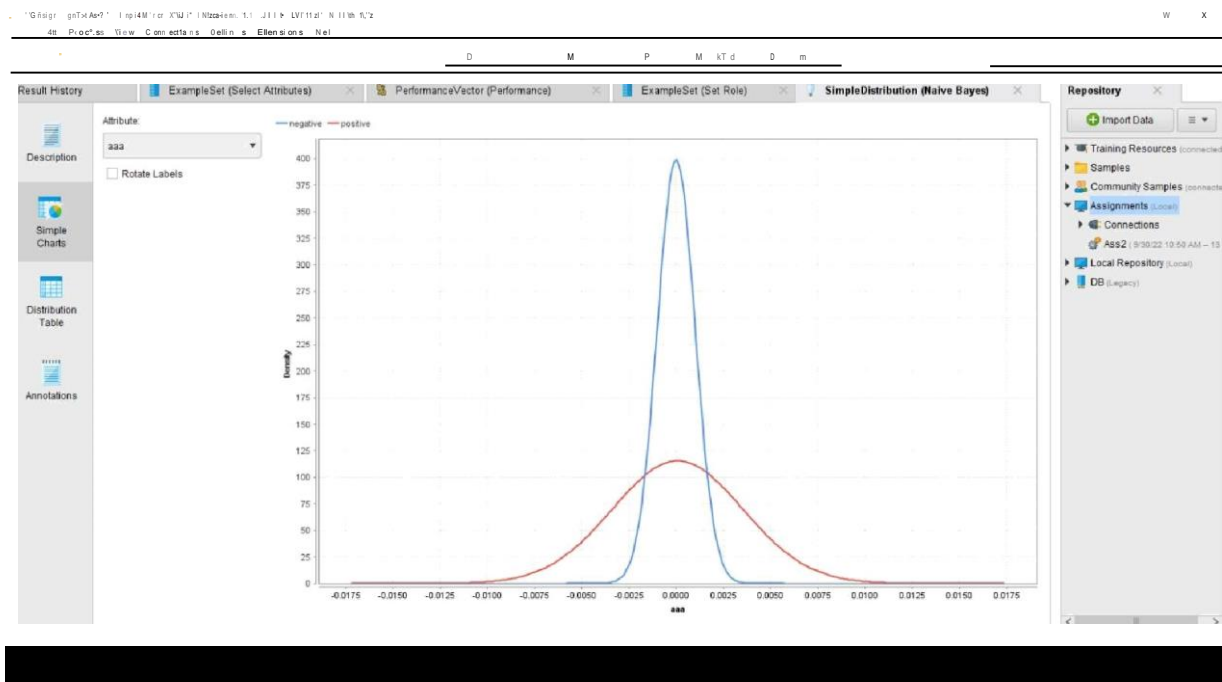
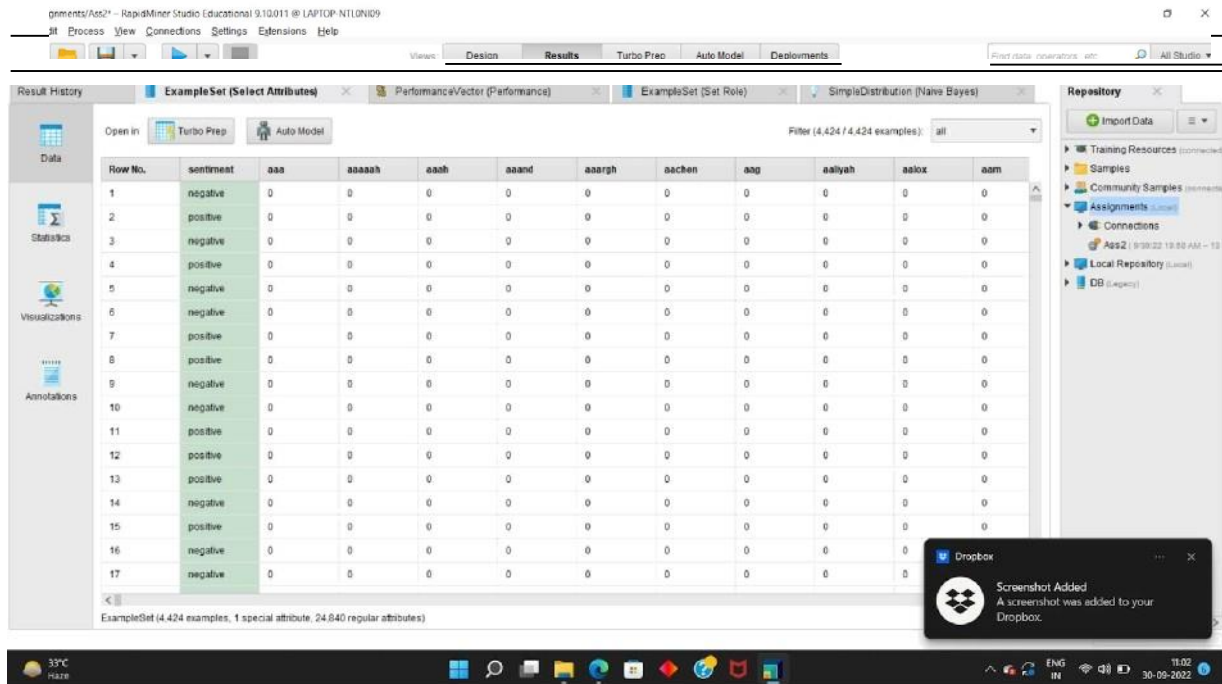
Using the filter operator we are able to distinguish between Positive and Negative review.

The processing time for about 49000 examples was around 6.5 minutes and error occurred when applying the Set-Role operator.



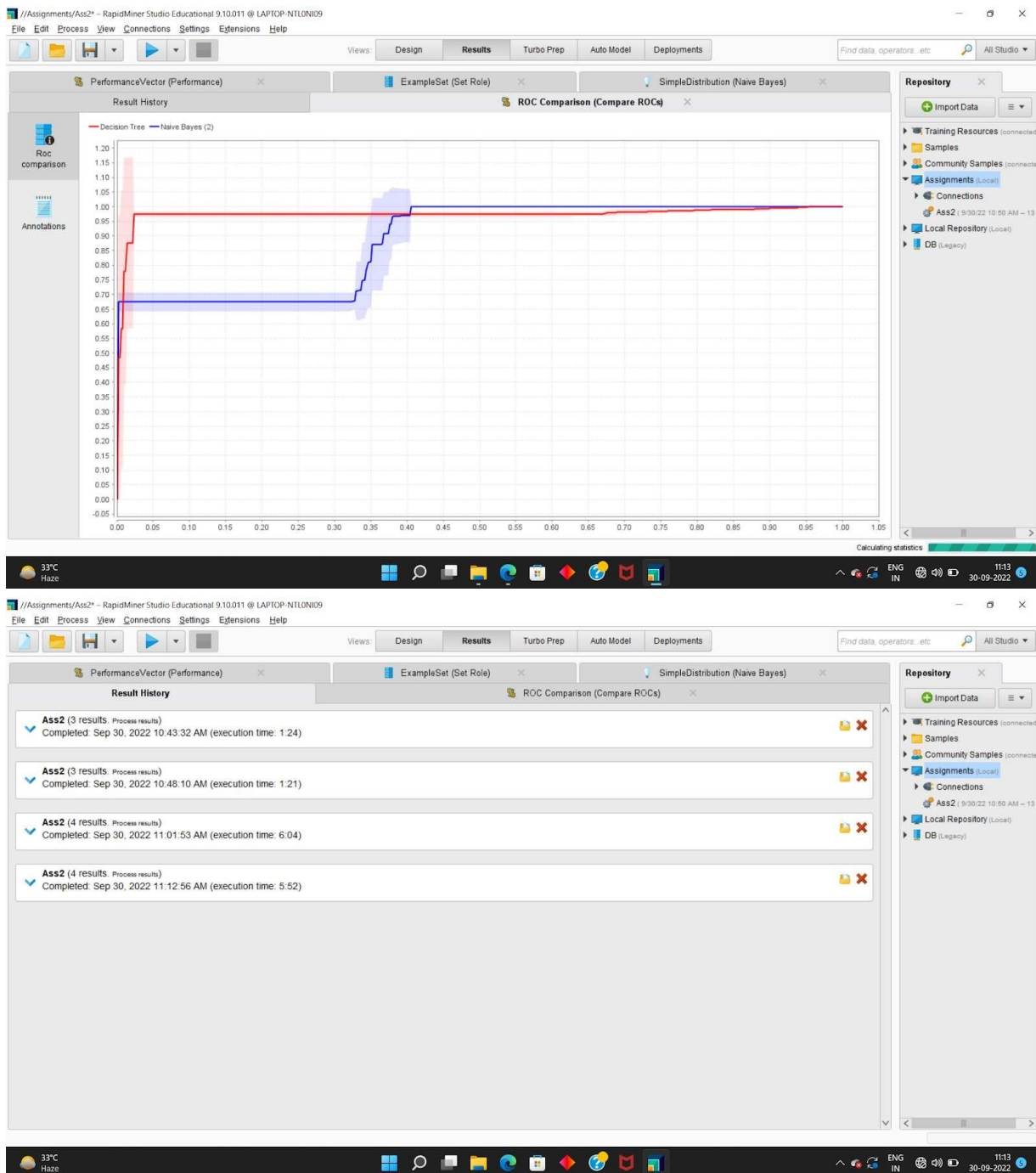
Without applying the filter operator the task performed smoothly.





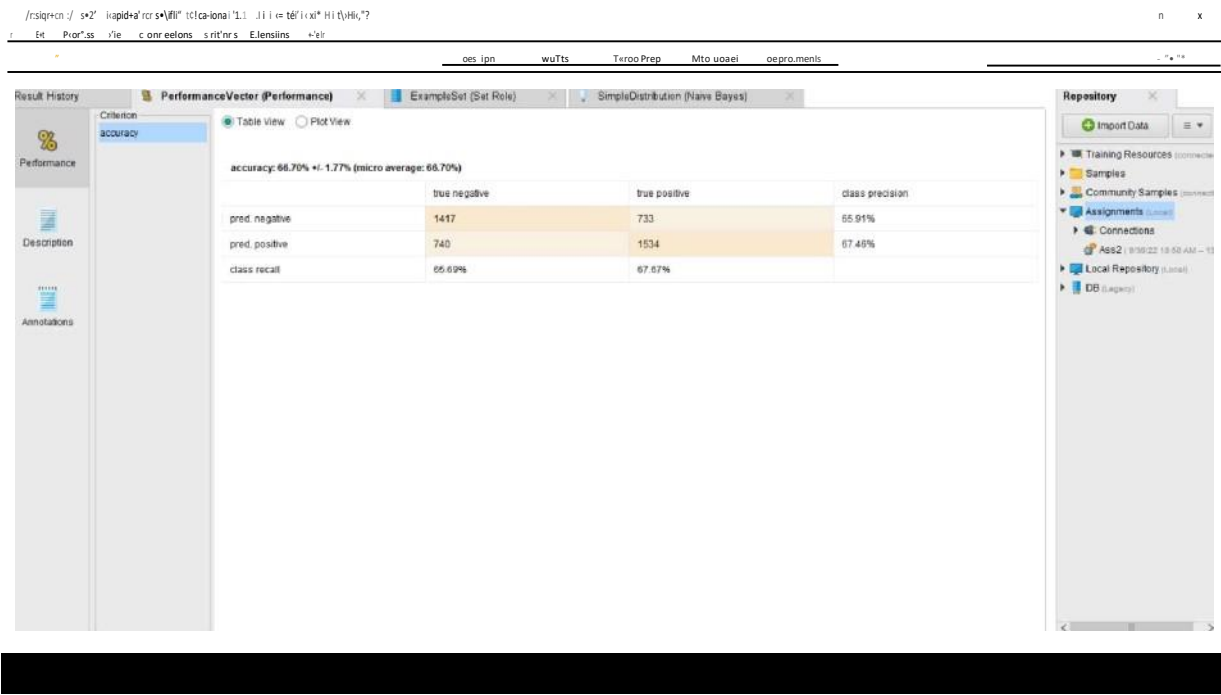
Here as seen from the figure curve of Naive bayes, the classifier distinguishes between positive and negative attributes. Standard deviation for positive attributes(Red) is greater than negative one(Blue).

## Comparing Performances:

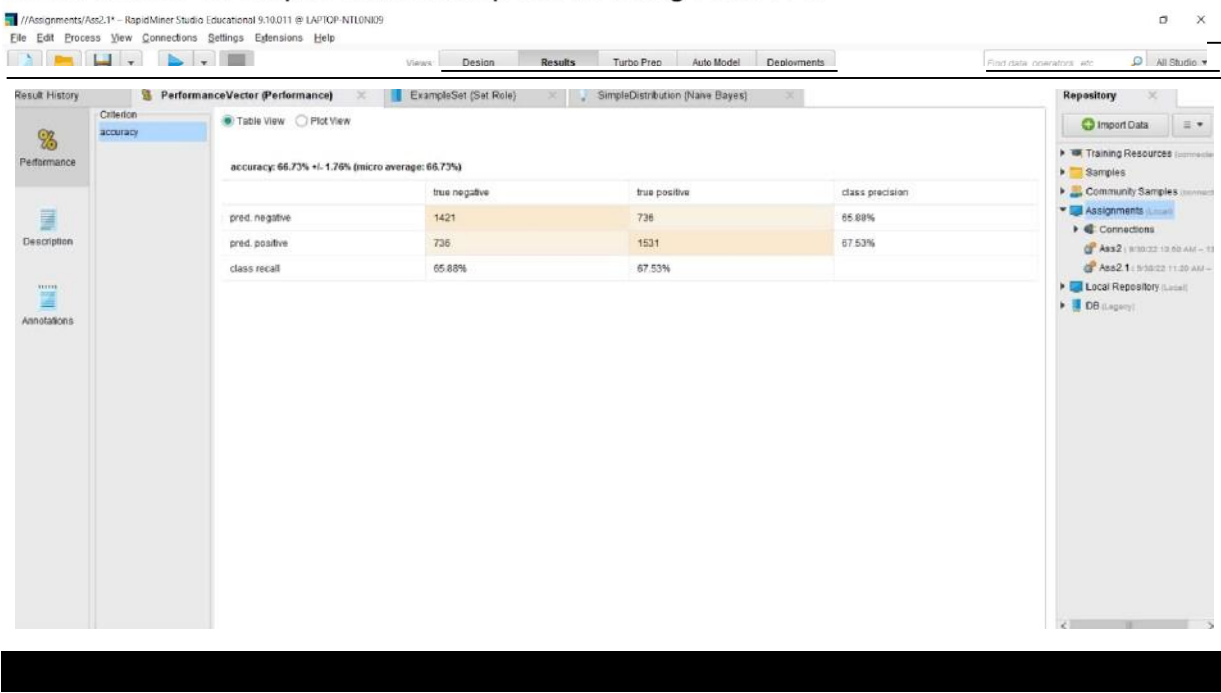


ROC is shown for decision trees and Naive Bayes classifiers. As from the ROC we can predict that Naive classifier maps more false positives. Whereas decision trees predict good for true positives. So, Decision tree is a good classifier as compared to naive bayes classifier w.r.t performance.

Model trained on corpus with stop words being removed.



Model trained on corpus without stop words being removed.



Although there is not much difference between above two images, but still the best result is shown by the model when stop words were being removed.

**Accuracy:** 66.73%( Note: we have taken less number of examples because more examples were taking more processing time)

**Recall:** 65.88%

**Precision:** 65.88%

## Task 2: Banknote authentication

The dataset which we received was in document format(.txt). We converted it in spreadsheet format(.csv)

Some modifications were done like addition of attributes heading of each column.

The attributes were,

1st attribute- Variance

2nd attribute- Skewness

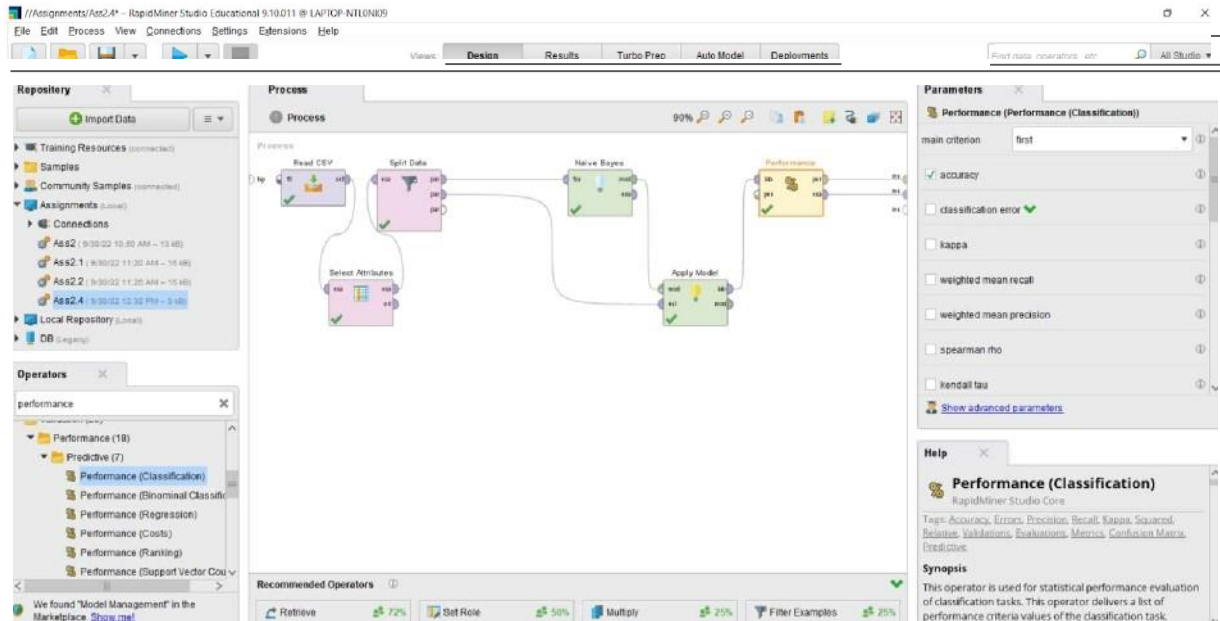
3rd attribute- Curtosis

4th attribute- Entropy

The Sth column was the class. In Rapidminer we converted the class attribute as 'label' as there were only two values '0' and '1'



The following image will show the set of operators we use.



## 1. Model Training:

We trained the model through a naive bayes classifier. The following image shows the distribution of results with a set of operators.

SimpleDistribution (Naive Bayes)

**Description**

Distribution model for label attribute class

**Simple Charts**

Class 0 (0.555)  
4 distributions

Class 1 (0.445)  
4 distributions

**Distribution Table**

**Annotations**

**Repository**

- Import Data
- Training Resources (connected)
  - Samples
  - Community Samples (connected)
  - Assignments (Local)
    - Ass2 (9/30/22 10:50 AM - 13 MB)
    - Ass2.1 (9/30/22 11:20 AM - 15 MB)
    - Ass2.2 (9/30/22 11:28 AM - 15 MB)
    - Ass2.4 (9/30/22 12:02 PM - 3 MB)
  - Local Repository (Local)
    - DB (Legacy)

SimpleDistribution (Naive Bayes)

**Process**

Process flow diagram showing the following steps:

- Read CSV
- Split Data
- Naive Bayes
- Apply Model

**Parameters**

**Split Data**

partitions: Edit Enumeration (2...)

sampling type: automatic

**Help**

**Split Data**

RapidMiner Studio Core

Tags: Divide, Separate, Part, Training, Testing, Samples, Subsets, Partitions, Sampling

**Synopsis**

This operator produces the desired number of subsets of the given ExampleSet. The ExampleSet is partitioned into subsets according to the specified relative sizes.

[Jump to Tutorial Process](#)

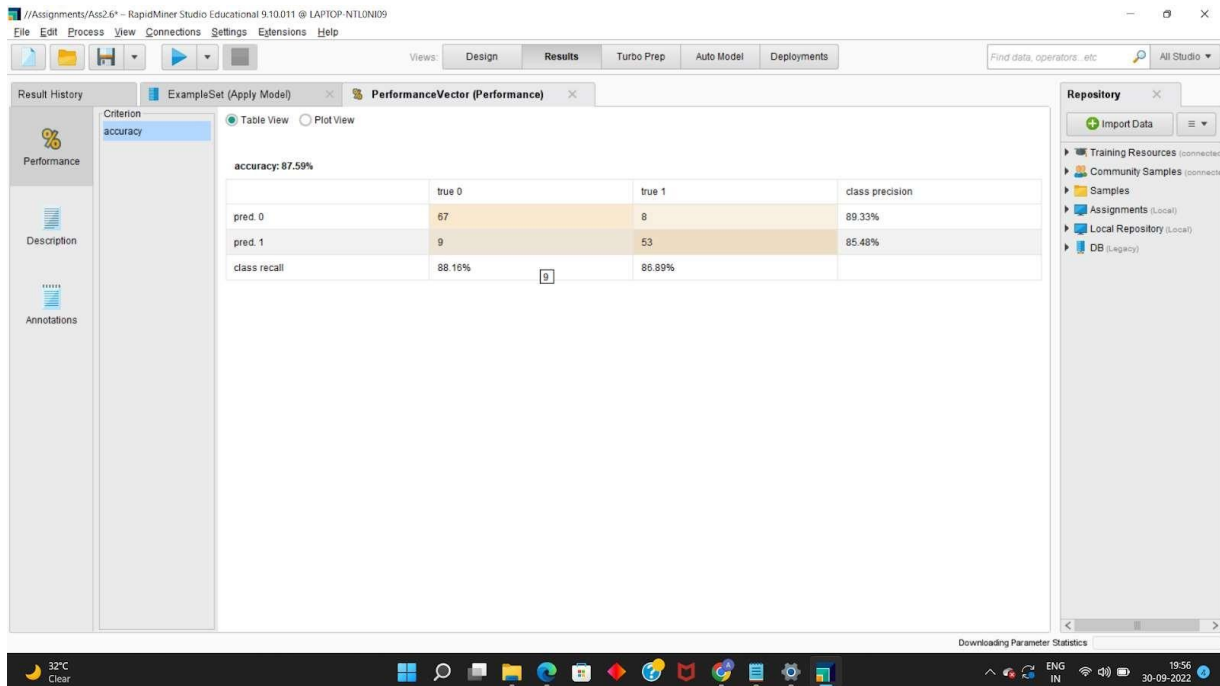
**Recommended Operators**

- Retrieve: 72%
- Set Role: 48%
- Select Attributes: 30%
- Performance (CL...: 26%

We applied 4 splits

1. 0.6 and 0.4
2. 0.7 and 0.3
3. 0.8 and 0.2
4. 0.9 and 0.1

The best results were obtained for 90% and 10% splitting, The accuracy was maximum song all 4. The accuracy was around 87%. The other three splits were showing accuracy around 84%.

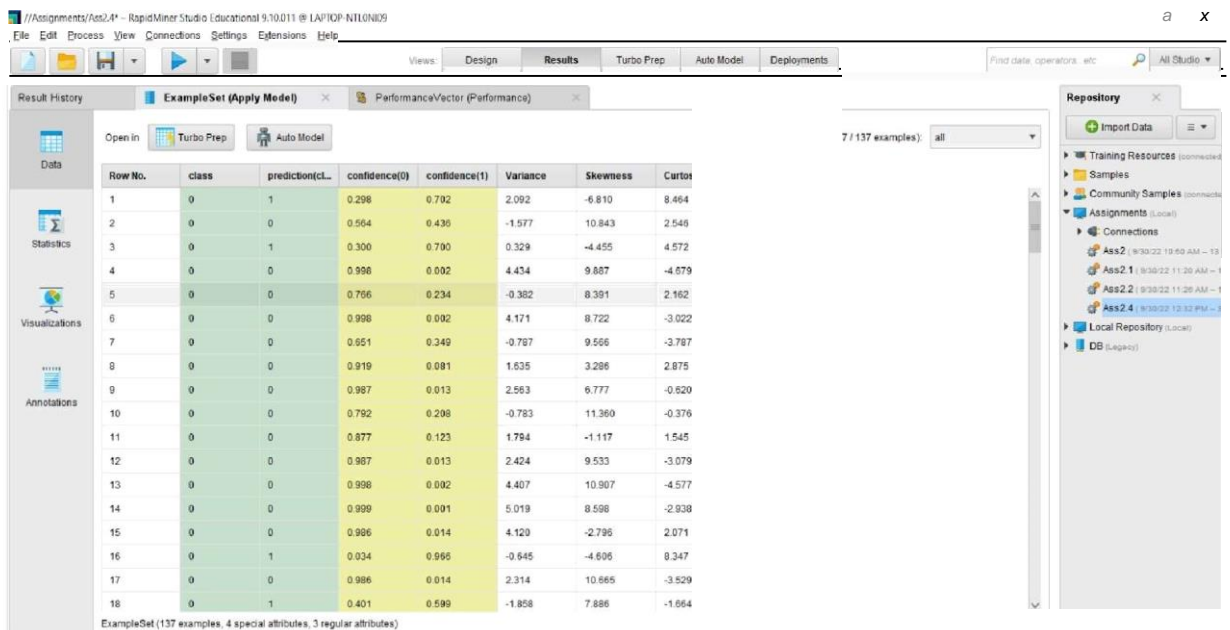
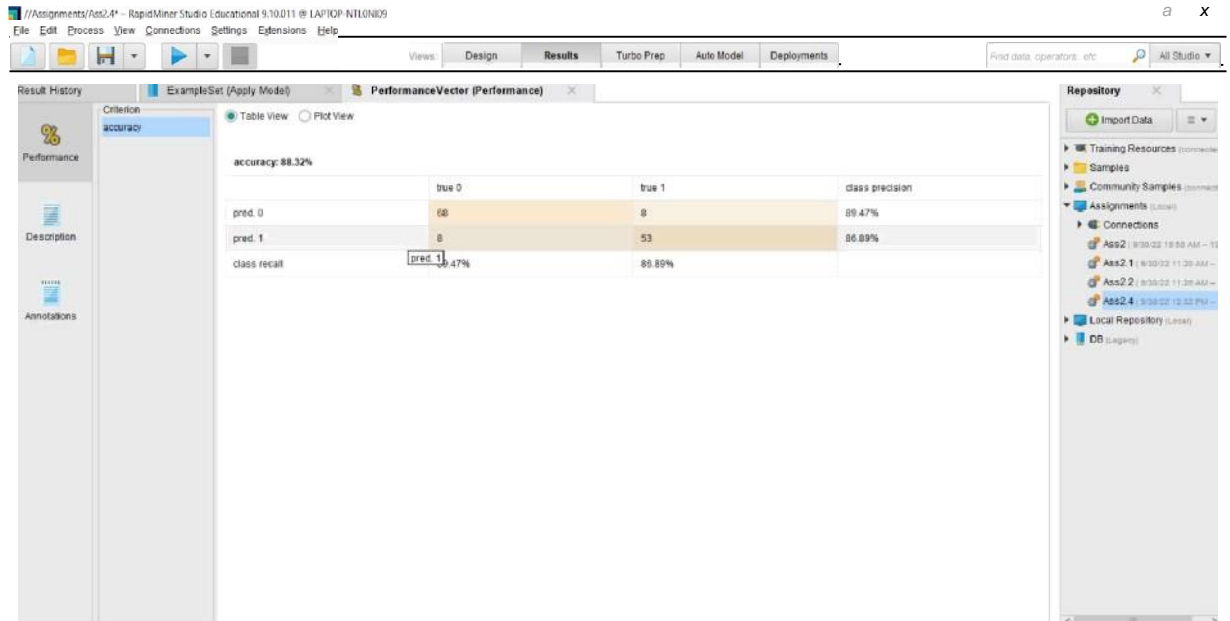


We carried further predictions with 90% and 10% split data.

To select the attributes, we 'Select attributes' operator. We changed the attribute filter to 'subtype' and then did the following activities.

### 1. Training with 1st three predictors.(Variance,Skewness,Curtosis)

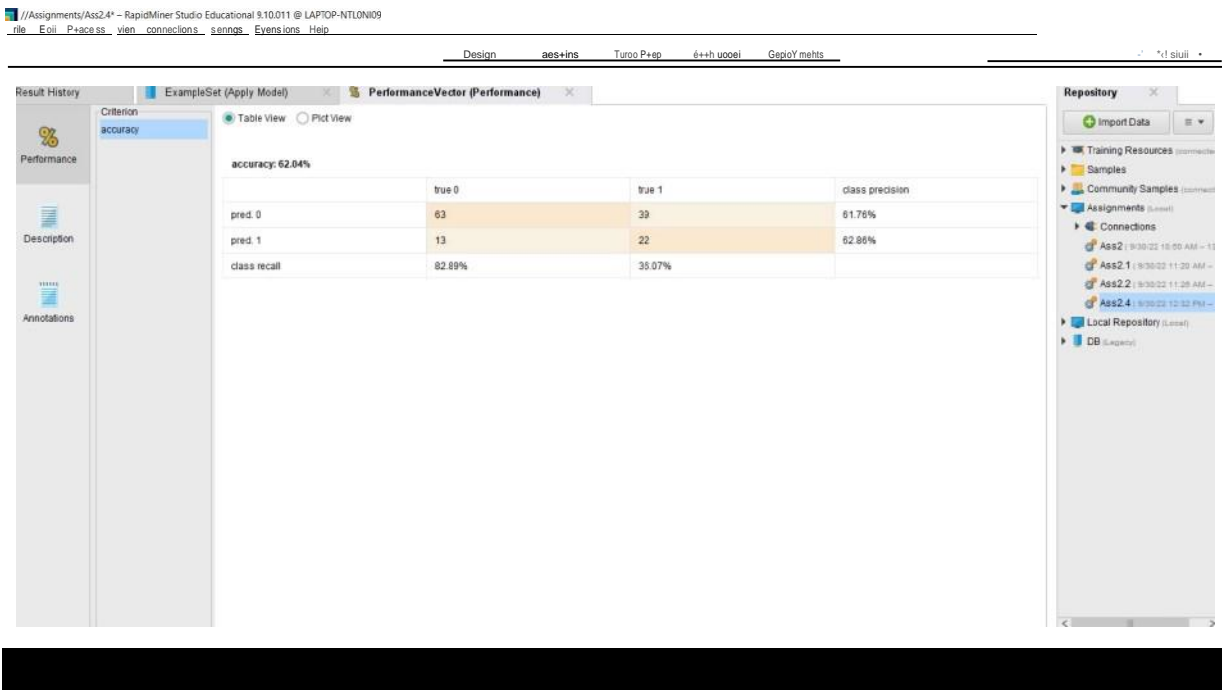
The following shows the result



The accuracy for 1<sup>st</sup> three predictors is 88.32%

## 2. Training with last 3 predictors (Skewness, Curtosis, Entropy):

The following image shows the result



The accuracy for last 3 predictors is 62.04%

Training with 1 2 and 4 predictors (Variance, Skewness, Entropy):

The following image shows the result

\\Assignments\Ass2.4\* - RapidMiner Studio Educational 9.10.011 @ LAPTOP-NTLON09  
File Edit Process View Connections Settings Extensions Help

De on +results Turbo Prep All o

Result History ExampleSet (Apply Model) PerformanceVector (Performance)

Open in Turbo Prep Auto Model

Filter (137 / 137 examples): all

Row No.	class	prediction(class)	confidence(0)	confidence(1)	Variance	Skewness	Entropy
1	0	0	0.670	0.330	2.092	-6.810	-0.602
2	0	1	0.460	0.540	-1.577	10.843	-2.936
3	0	1	0.316	0.684	0.329	-4.455	-0.989
4	0	0	0.998	0.002	4.434	9.887	-3.748
5	0	0	0.664	0.336	-0.382	8.391	-3.740
6	0	0	0.998	0.002	4.171	8.722	-0.597
7	0	0	0.574	0.426	-0.787	9.586	-7.503
8	0	0	0.899	0.101	1.635	3.286	0.087
9	0	0	0.979	0.021	2.563	6.777	0.386
10	0	0	0.643	0.357	-0.783	11.360	-7.050
11	0	0	0.823	0.177	1.794	-1.117	-0.261
12	0	0	0.982	0.018	2.424	9.533	-2.775
13	0	0	0.998	0.002	4.407	10.907	-4.427
14	0	0	0.999	0.001	5.019	8.596	-1.281
15	0	0	0.981	0.019	4.120	-2.796	0.674
16	0	1	0.120	0.872	-0.645	-4.606	-2.710
17	0	0	0.981	0.019	2.314	10.665	-4.767
18	0	1	0.296	0.704	-1.858	7.886	-1.838

ExampleSet (137 examples, 4 special attributes, 3 regular attributes)

Repository

- Import Data
- Training Resources (connected)
- Samples
- Community Samples (connected)
- Assignments (Local)
  - Ass2 (9/30/22 10:55 AM - 11:00 AM)
  - Ass2.1 (9/30/22 11:20 AM - 11:25 AM)
  - Ass2.2 (9/30/22 11:28 AM - 11:33 AM)
  - Ass2.4 (9/30/22 12:32 PM - 12:37 PM)
- Local Repository (Local)
- DB (Legacy)

\\Assignments\Ass2.4\* - RapidMiner Studio Educational 9.10.011 @ LAPTOP-NTLON09  
File Edit Process View Connections Settings Extensions Help

De on +results Turbo Prep All o

Result History ExampleSet (Apply Model) PerformanceVector (Performance)

Criterion accuracy

Table View Plot View

accuracy: 89.05%

	true 0	true 1	class prediction
pred. 0	67	6	91.78%
pred. 1	9	55	85.94%
class recall	88.16%	90.16%	

Repository

- Import Data
- Training Resources (connected)
- Samples
- Community Samples (connected)
- Assignments (Local)
  - Ass2 (9/30/22 10:55 AM - 11:00 AM)
  - Ass2.1 (9/30/22 11:20 AM - 11:25 AM)
  - Ass2.2 (9/30/22 11:28 AM - 11:33 AM)
  - Ass2.4 (9/30/22 12:32 PM - 12:37 PM)
- Local Repository (Local)
- DB (Legacy)

The accuracy for 1<sup>st</sup>, 2<sup>nd</sup> and 4<sup>th</sup> predictors is 89.05%

Training with 1<sup>st</sup>, 3<sup>rd</sup> and 4<sup>th</sup> predictors (Variance, Kurtosis, Entropy):

The following image shows the result

**PerformanceVector (Performance)**

accuracy: 81.75%

	true 0	true 1	class precision
pred. 0	62	11	84.93%
pred. 1	14	50	78.12%
class recall	81.58%	81.97%	

**ExampleSet (Apply Model)**

Filter (137 / 137 examples): all

Row No.	class	prediction(class)	confidence(0)	confidence(1)	Variance	Curtosis	Entropy
1	0	0	0.688	0.312	2.092	8.464	-0.602
2	0	1	0.200	0.800	-1.577	2.546	-2.936
3	0	0	0.573	0.427	0.329	4.572	-0.989
4	0	0	0.992	0.008	4.434	-4.679	-3.748
5	0	1	0.472	0.528	-0.382	2.162	-3.740
6	0	0	0.993	0.007	4.171	-3.022	-0.597
7	0	1	0.270	0.730	-0.787	-3.787	-7.503
8	0	0	0.895	0.105	1.635	2.875	0.087
9	0	0	0.968	0.032	2.563	-0.620	0.386
10	0	1	0.368	0.632	-0.783	-0.376	-7.050
11	0	0	0.922	0.078	1.794	1.545	-0.261
12	0	0	0.946	0.054	2.424	-3.079	-2.775
13	0	0	0.992	0.008	4.407	-4.577	-4.427
14	0	0	0.998	0.002	5.019	-2.938	-1.281
15	0	0	0.994	0.006	4.120	2.071	0.674
16	0	1	0.096	0.904	-0.645	8.347	-2.710
17	0	0	0.928	0.072	2.314	-3.529	-4.767
18	0	1	0.176	0.824	-1.858	-1.664	-1.838

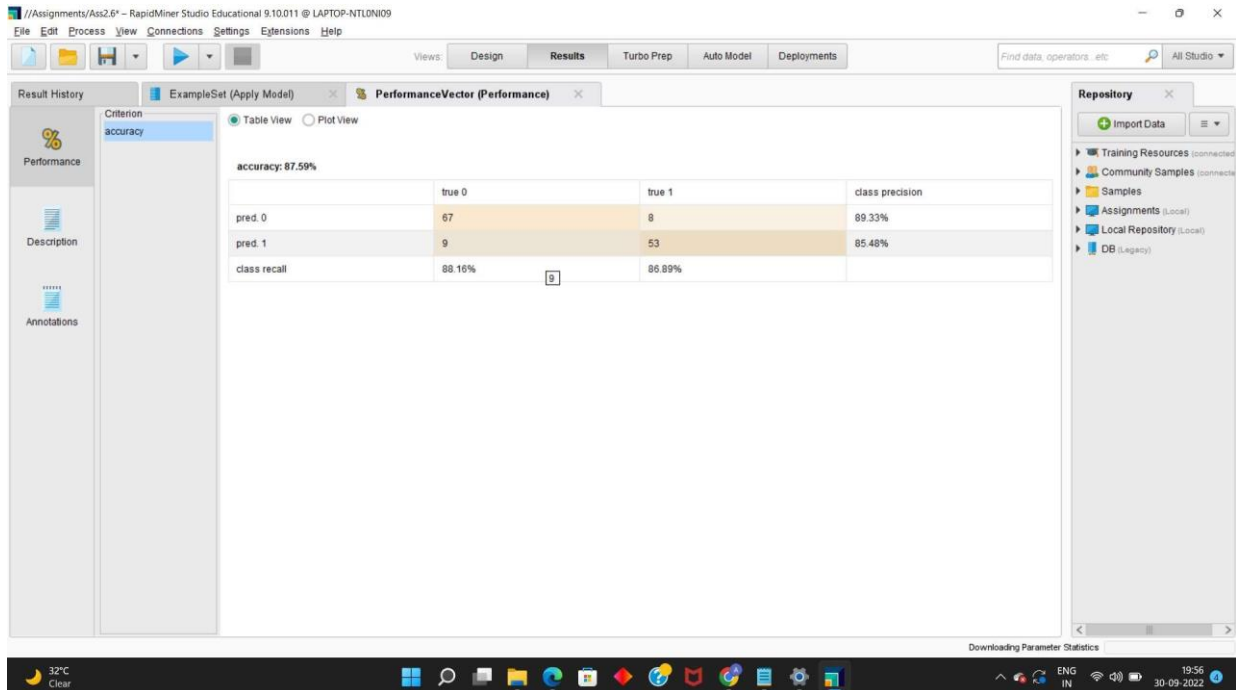
ExampleSet (137 examples, 4 special attributes, 3 regular attributes)

The accuracy is 81.75%

### 3.Comparing Performances:

From above set of models, the best model is 3rd model having predictors 1st, 2nd and 4th predictors. It shows the maximum accuracy of 89.04%. Not only that,

but if we compare with Naive Bayes classifier( having all four attributes), the 3rd model is still the best because naive bayes classifier has accuracy of 87.58%



The above image shows the result of naive bayes classifier

Therefore the accuracy,recall,precision for best model is:

Accuracy: 89.05%

Recall: 88.16%

Precision:91.78%

The result(confusion matrix) has already been shown for all the variants of model.