

Project#2

HAMZA ZAFAR

SYED JEHANDAD KAMAL

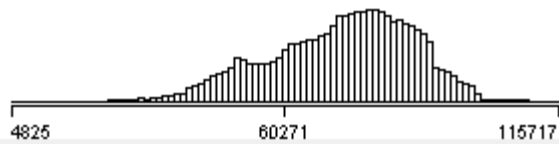
ABDUL MOIZ

Don't Get Kicked!

- Predict if a car purchased at auction is a lemon
 - There are columns 34
 - The data set is split to 60% (72983) training and 40% (48707) testing.

Nature of Data

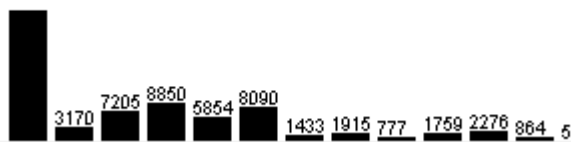
VehOdo



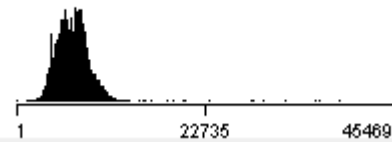
Make



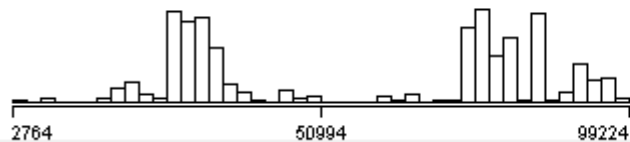
Size



VehBCost



VNZIP1



VNST



Issues with data

- Date contains many formats.
- Nominal values contain NULL Values
- Skewed data.
- Numeric data may contain outliers.

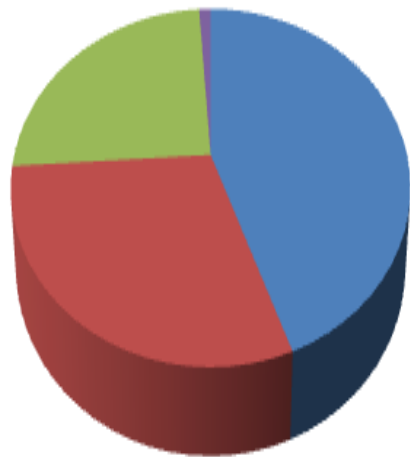
Data Treatment

- Missing Values
 - Most Frequently and mean
- DATE:
 - Convert to same format (dd-mmm-yyyy)
 - drill down/ drill up approach is used to see data more deeply i.e. year to month to day.
- NULL:
 - Replace by most frequently values
 - Check the weight of 1 in null values and values having more weight to the null are kept as null.

NULL TREATMENT

Wheel Type

1

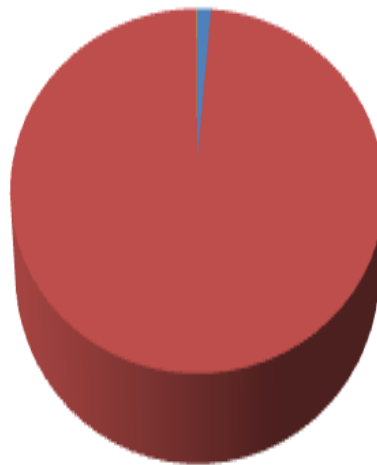


WheelType ▾

- Alloy
- Covers
- NULL
- Special

Prime Unit

1



PRIMEUNIT ▾

- NO
- NULL
- YES

Acguart

1

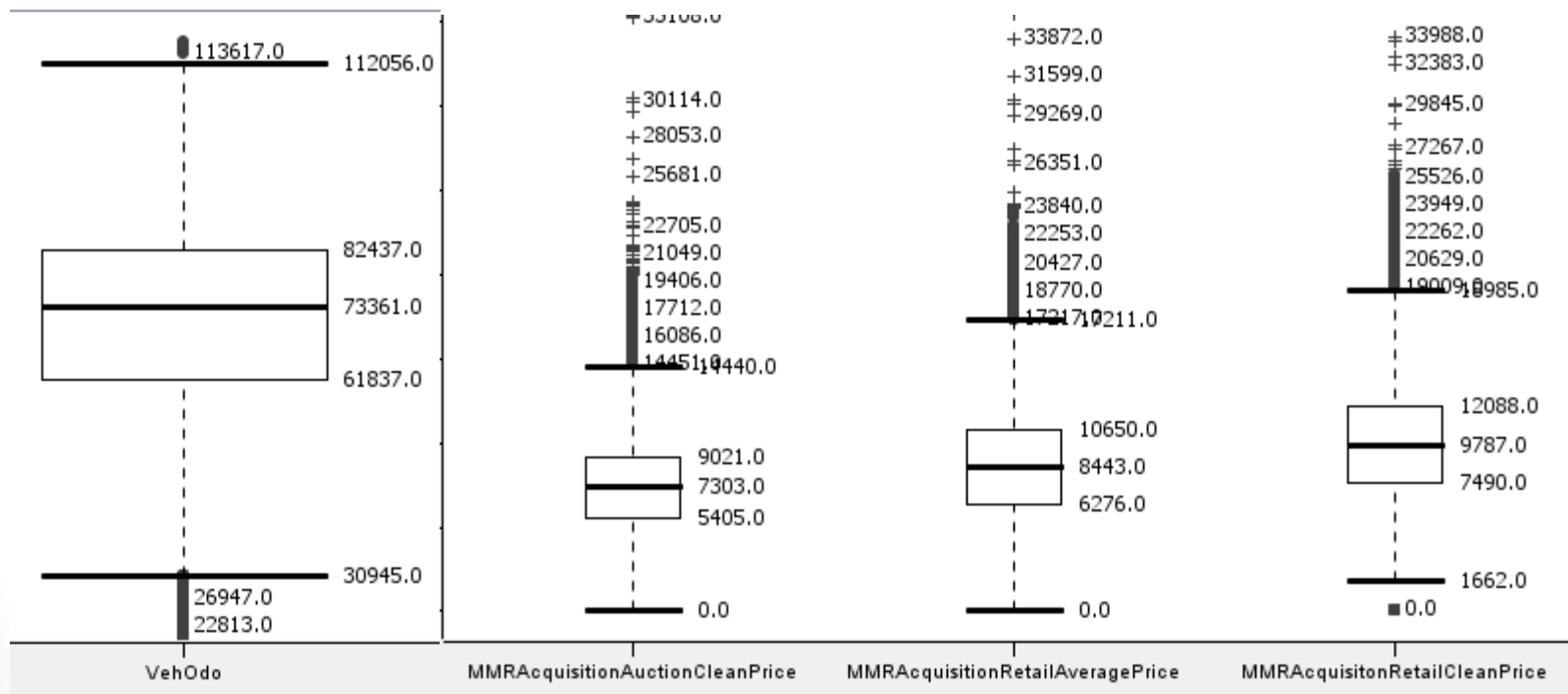


AUCGUART ▾

- GREEN
- NULL
- RED

Data Treatment

- Outliers:

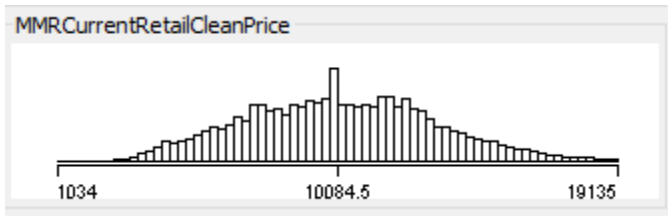
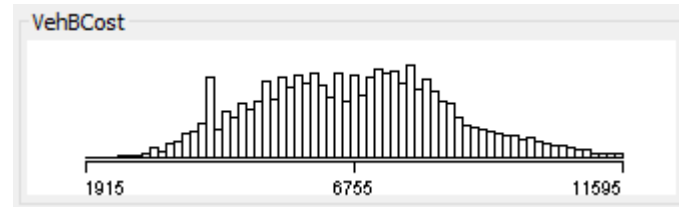
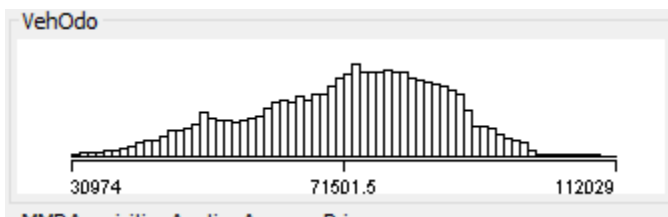


Transformation

COLUMN	VALUE	TRANSFORMED
Transmission	Auto	A
	Manual	M
Wheel type	Alloy	A
	Null	N
	Special	S
	Cover	C
Primeunit	Yes	Y
	No	N
	Null	NI
Aucguart	Green	G
	Red	R
	Yellow	Y
	Null	N
topThreeName	chrysler	Ch
	Other	O
	Ford	F

After Treatment

-



Feature Selection

- Wheelid,Refid,buyerNo,Vnzip removed..
- Backward Feature Elemination
- InfoGain
- Chisquare

Binning

- fuzzy c cluster is used
- Auto Binning

*0: Null Treatment

*0:3 - Feature Elimination X

Replace this node by your
favourite model learner

Naive Bayes Learner

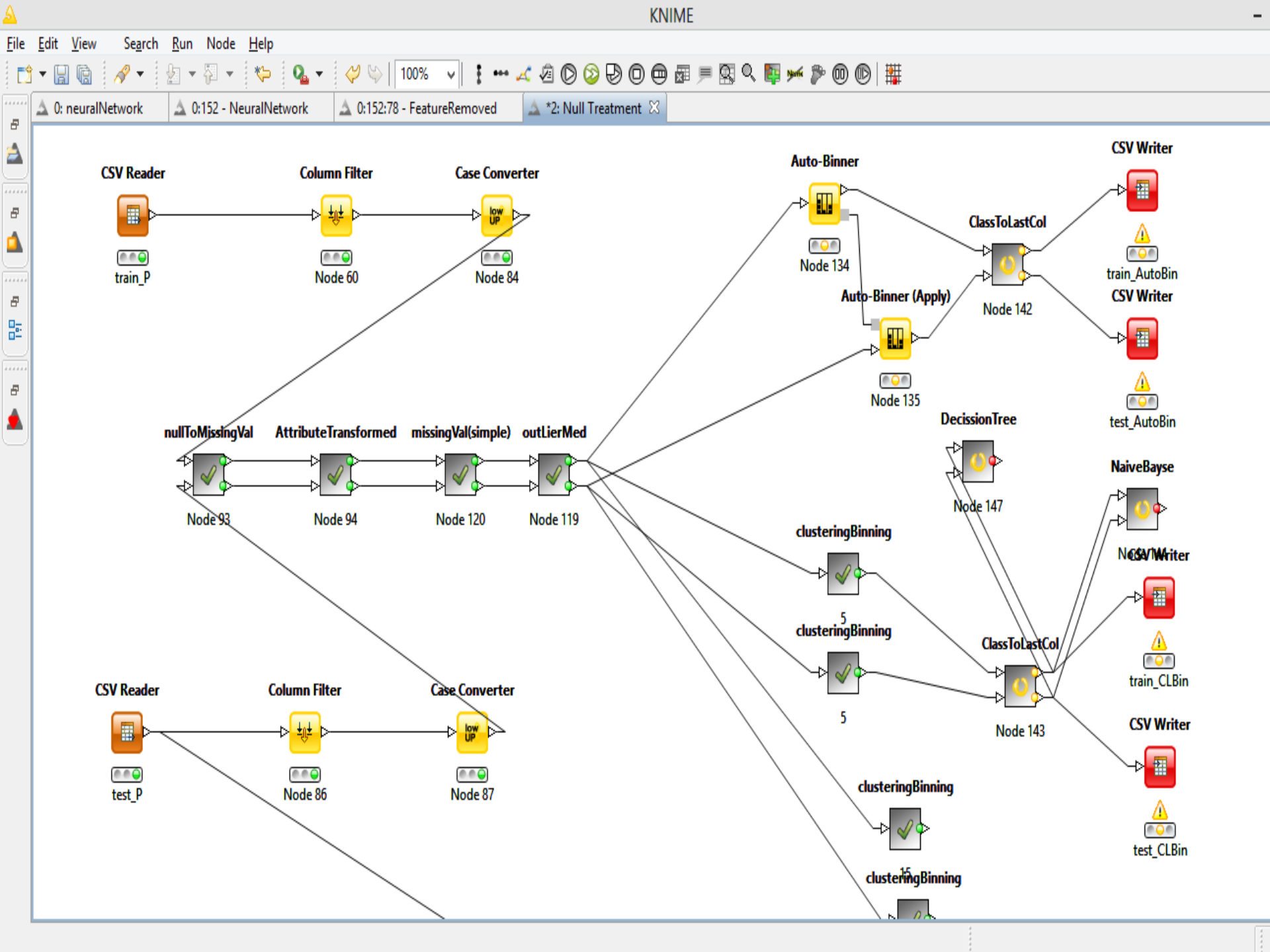


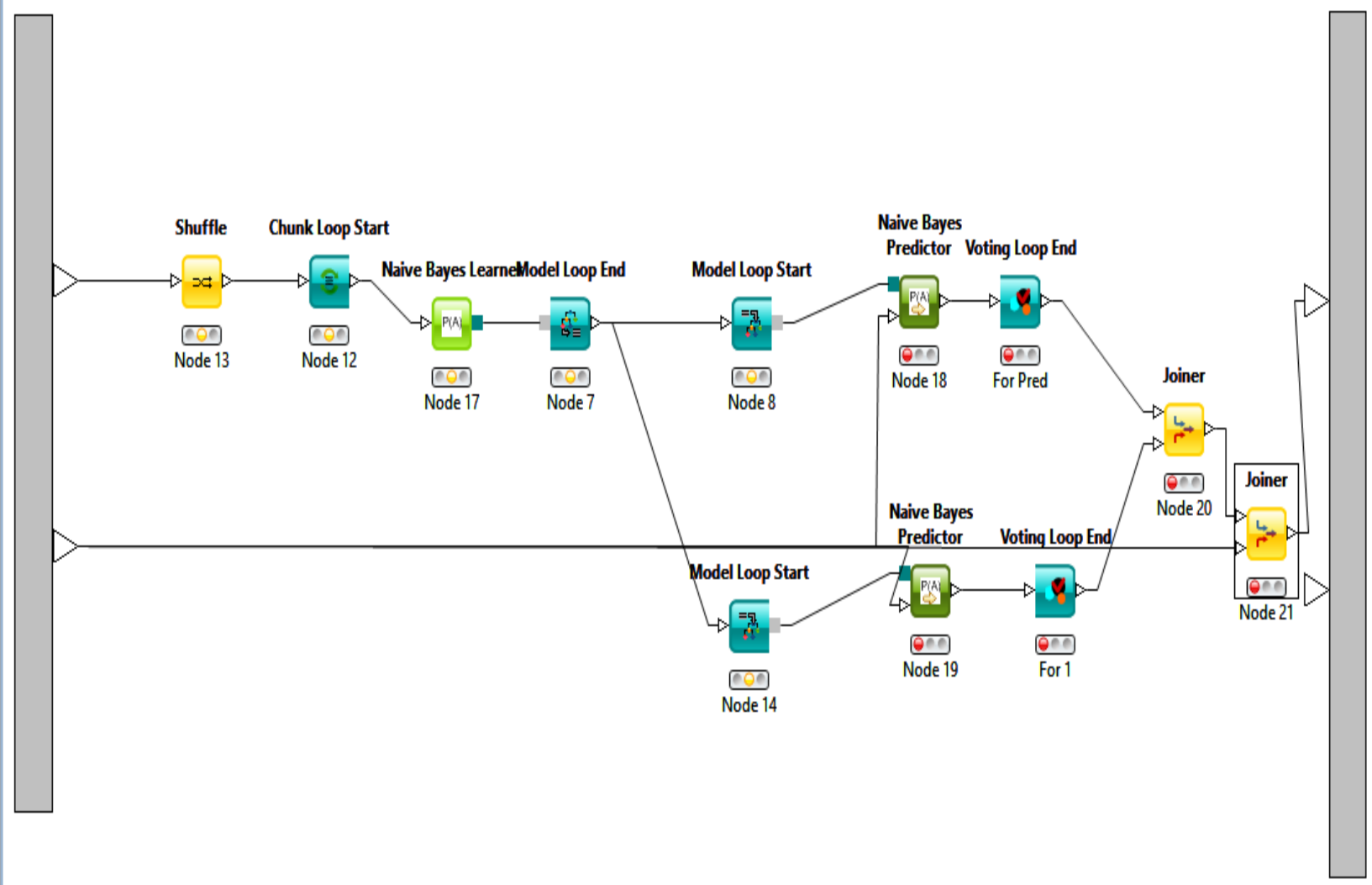
Replace this node by the
appropriate predictor

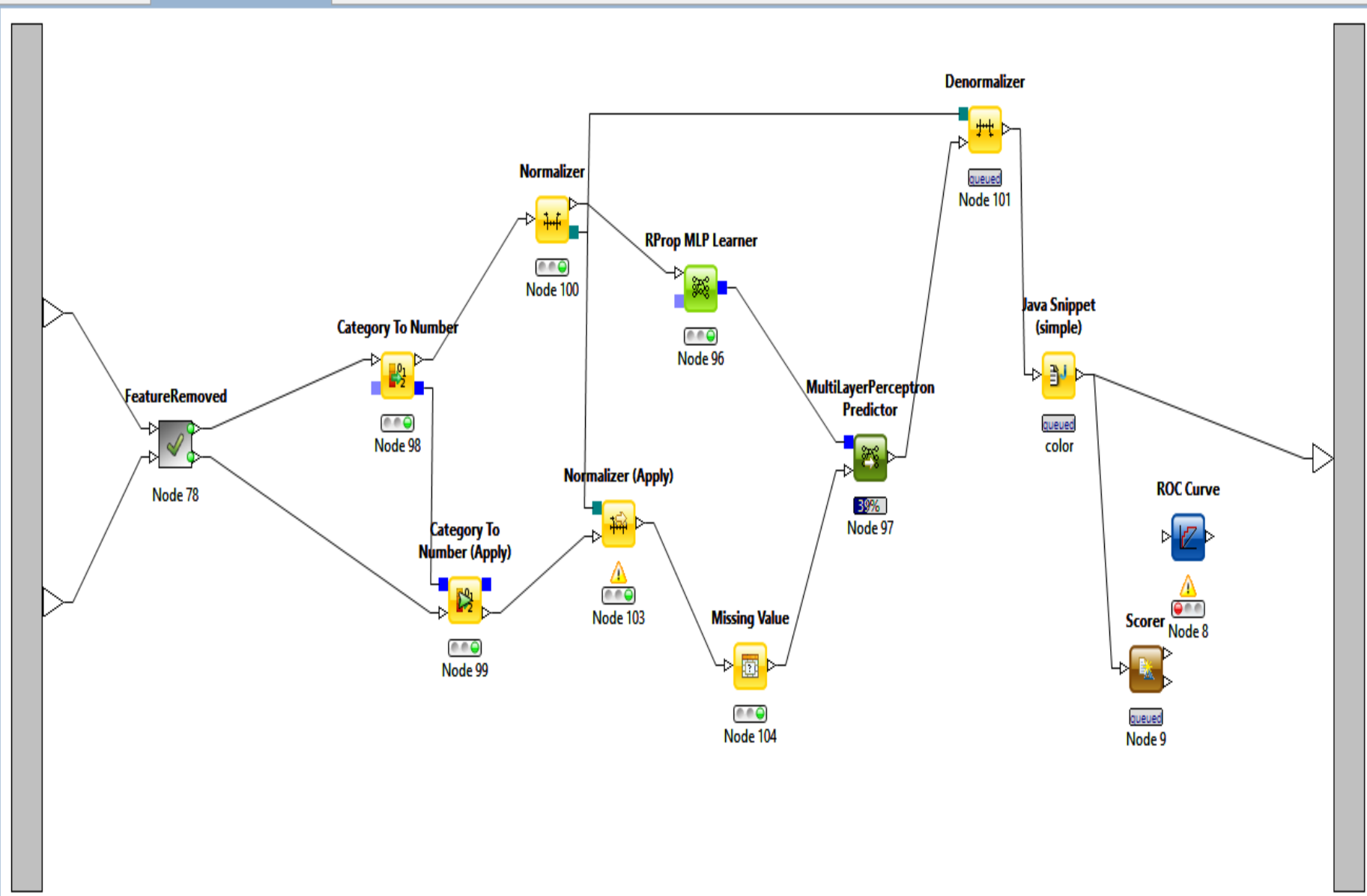
Naive Bayes
PredictorBackward Feature
Elimination Start (1:1)

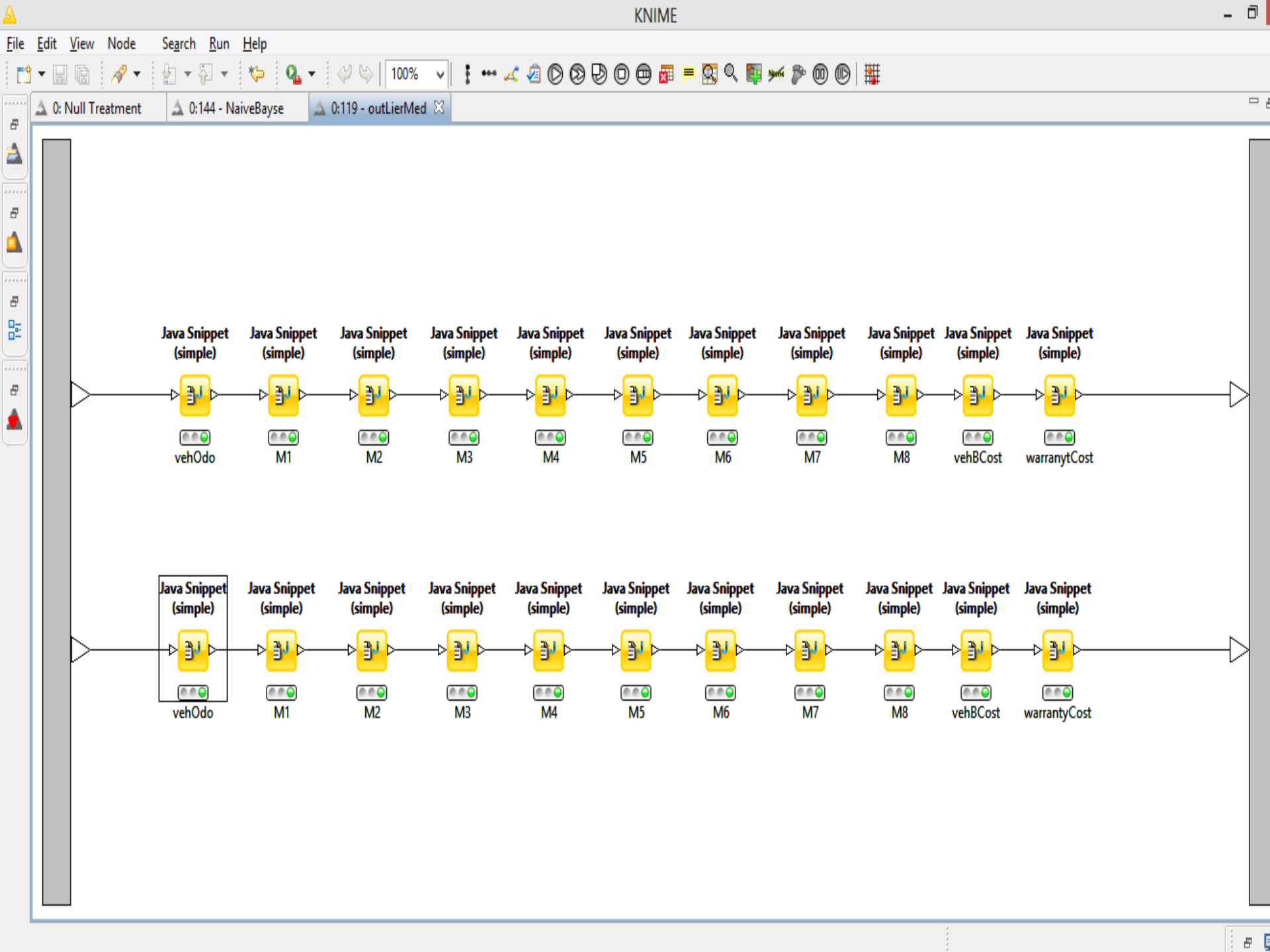
Partitioning

Backward Feature
Elimination EndBackward Feature
Elimination Filter









RESULT COMPARISION

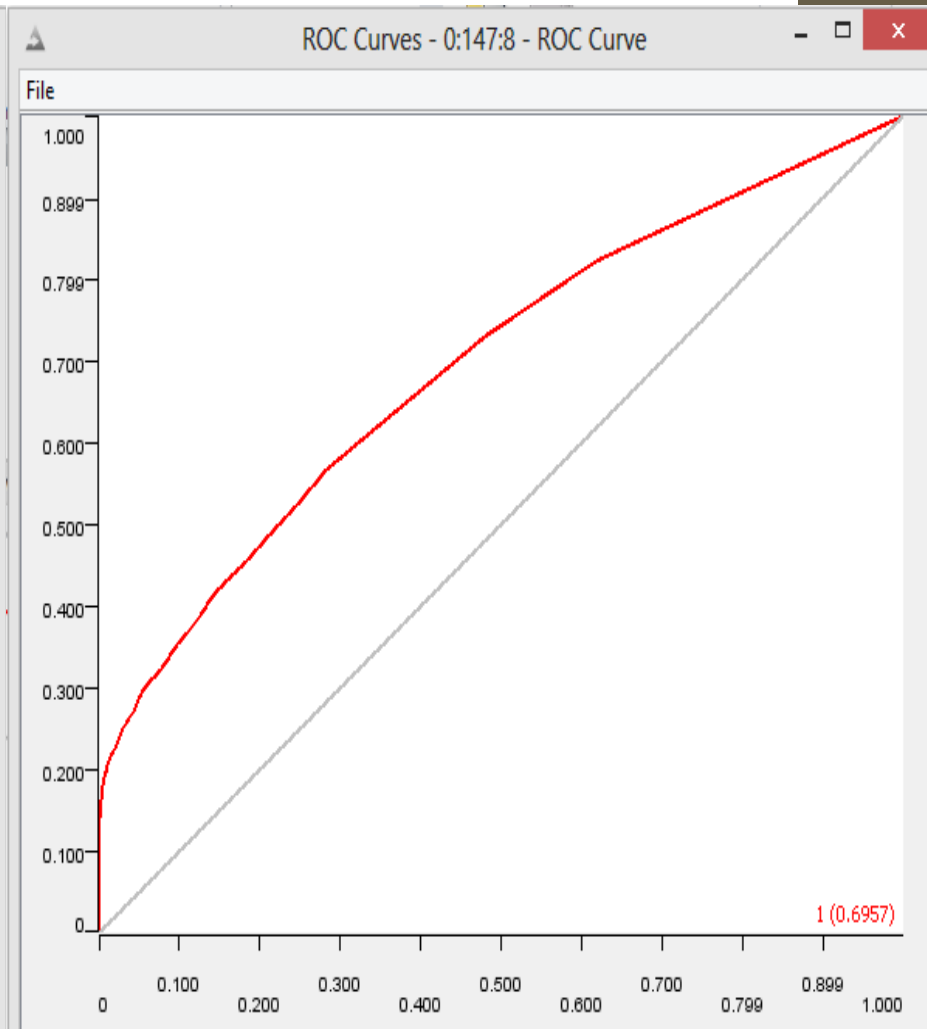
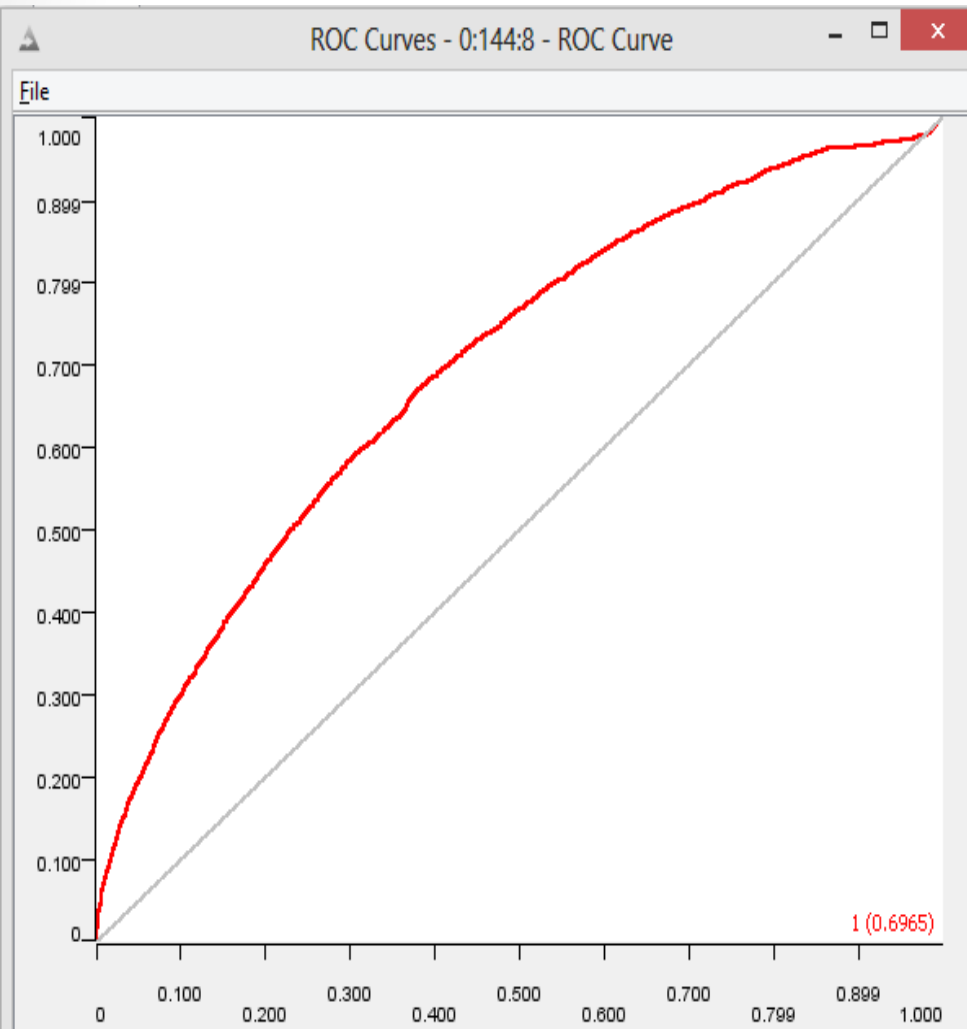
NAÏVE BAYSE	CLUSTER BINNING		AUTO-BINNER	
	ROC	FMEASURE	ROC	FMEASURE
Simple	0.696	0.318	0.693	0.317
Bagging	0.653	0.308	0.66	0.306
Boosting	0.694	0.315	0.693	0.317
DECISSION TREE				
Simple	0.597	0.318	0.603	0.322
Bagging	0.595	0.344	0.601	0.344
Boosting	0.607	0.316	0.624	0.317
Tree Ensembler	0.667	0.318	0.697	0.337

For Evaluation of model train data is spilted in 2/3 as training and 1/3 for testing

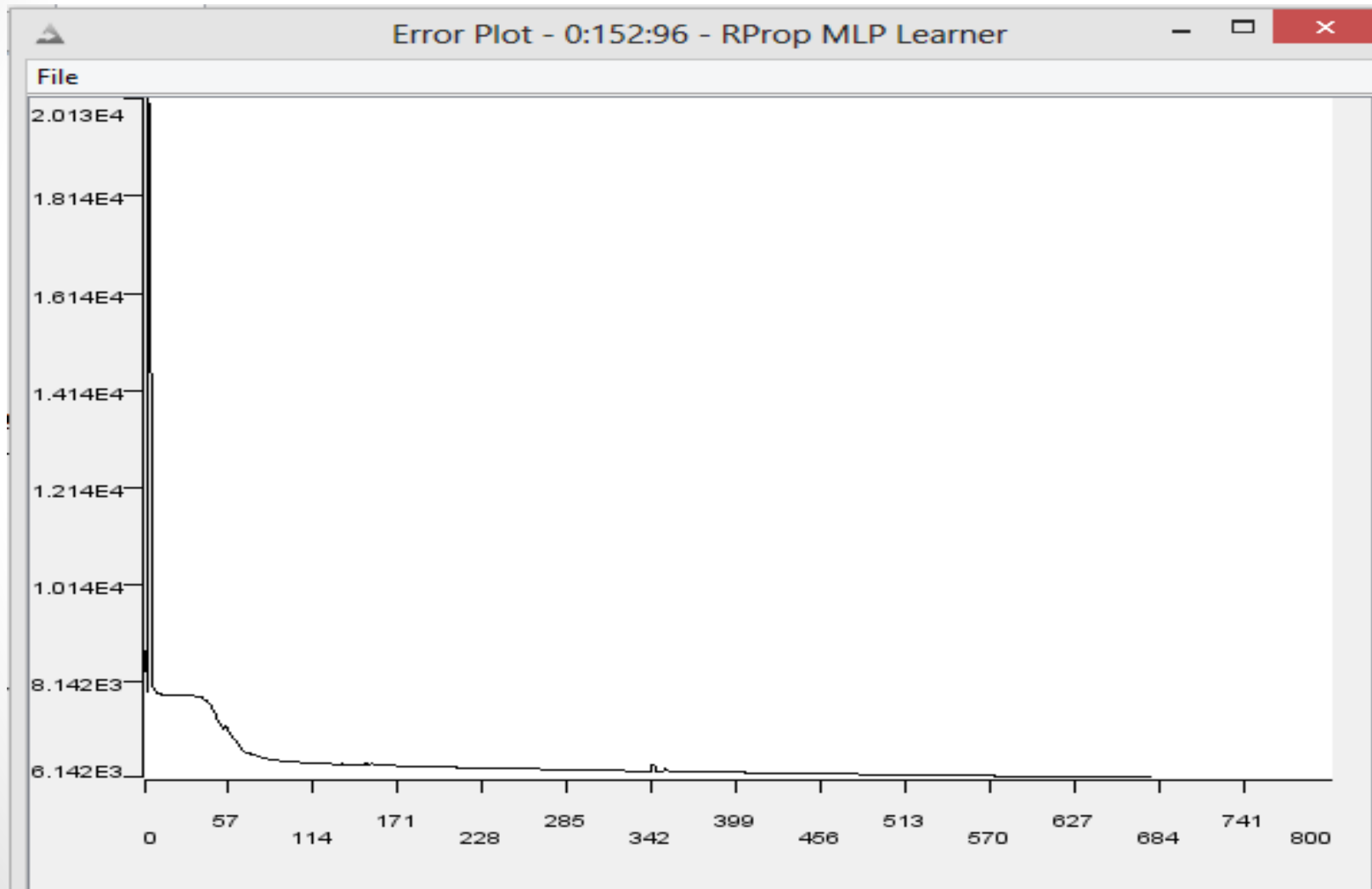
ROC

- NAÏVE BAYSE

DECISION TREE



SSE Neural Network



Gini

Learner	Gini
Naïve Bayse (Simple) (5bins)	0.12374(446)
Decission Tree Ensemebler (50 Models,All columns)	0.09789
Neural Network (800it,10hidenLayer,25neurons)	0.1053

Conclusion

- The Decision Tree with Ensembling Meta Node gave the best Roc and fmeasure when train data is spited in 2/3 as training and 1/3 as testing
- Naïve Bayes with Simple Learner gave the best Gini i.e. 0.12374(It scored 446 on leaderboard)
- Neural Network gave 0.1053, and may be improved with more numbers of iterations.