





## LOGISTICS REGRESSION STEPWISE

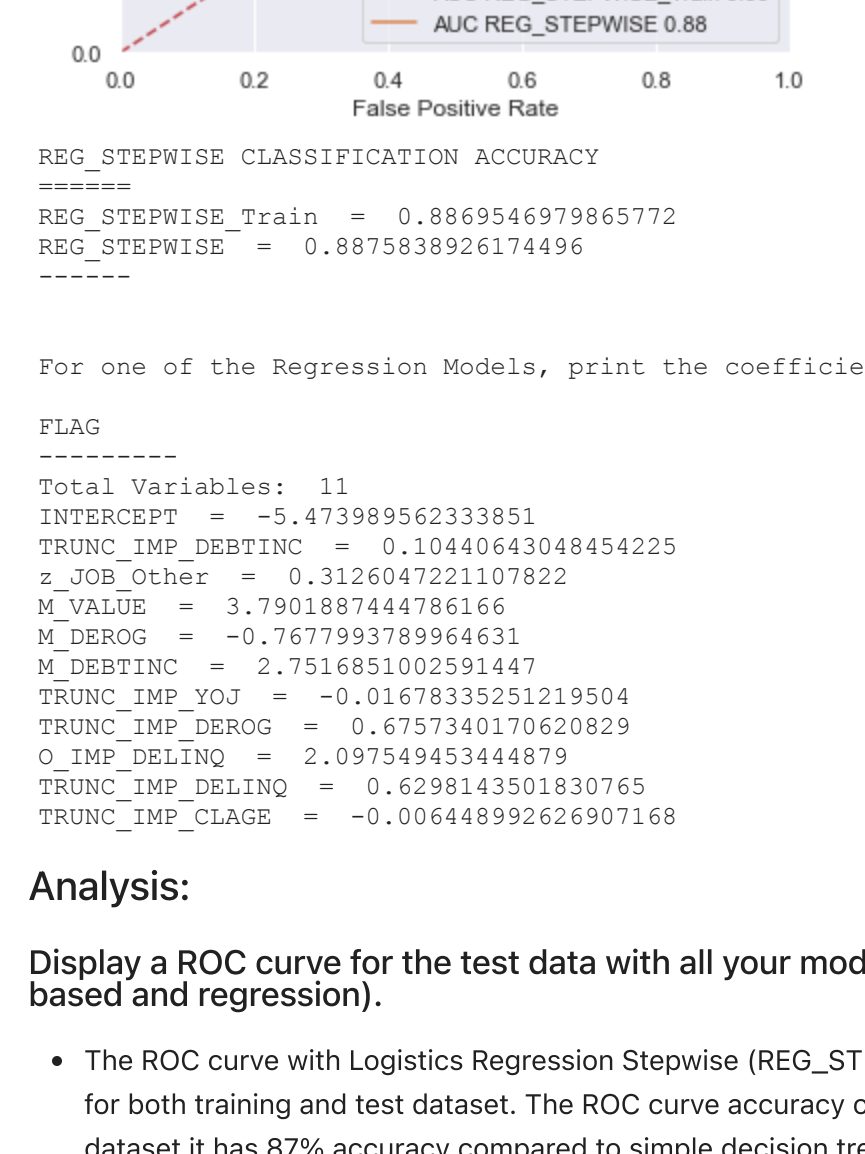
```
WHO = "REG_STEPWISE"

CLM = LogisticRegression( solver='newton-cg', max_iter=1000 )
CLM = CLM.fit( X_train, Y_train[ TARGET_F ] )

TRAIN_CLM = getProbaAccuracyScores( WHO + " Train", CLM, X_train, Y_train[ TARGET_F ] )
TEST_CLM = getProbaAccuracyScores( WHO, CLM, X_test, Y_test[ TARGET_F ] )

print_ROC_Curve( WHO, { TRAIN_CLM, TEST_CLM } )
print_Accuracy( WHO + " CLASSIFICATION ACCURACY", { TRAIN_CLM, TEST_CLM } )

print("For one of the Regression Models, print the coefficients. Do the variables make sense?")
REG_STEP_CLM_COEF = getCoefLogit( CLM, X_train )
REG_STEP_CLM = TEST_CLM.copy()
```



```
REG_STEPWISE CLASSIFICATION ACCURACY
=====
REG_STEPWISE_Train    = 0.8869546979865772
REG_STEPWISE    = 0.8875838926174496
-----
```

For one of the Regression Models, print the coefficients. Do the variables make sense?

```
FLAG
-----
Total Variables: 11
INTERCEPT    = -5.473989562338851
TRUNC_IMP_DEBTINC    = 0.10440643048454225
Z_VALUE_Other    = 0.312604722107822
M_VALUE    = 3.790187444786166
M_DEROG    = 0.7671993789964631
W_DEBTINC    = 2.751489310029391447
TRUNC_IMP_YOY    = -0.01678335251219504
TRUNC_IMP_DEROG    = 0.6757340170620629
O_IMP_DELINQ    = 2.0975494953444879
TRUNC_IMP_DELINQ    = 0.6298143501830765
TRUNC_IMP_CLAGE    = -0.006448992626907168
```

### Analysis:

Display a ROC curve for the test data with all your models on the same graph (tree based and regression).

- The ROC curve with Logistics Regression Stepwise (REG\_STEPWISE) has accuracy rate of 88% for both training and test dataset. The ROC curve accuracy on Regression Stepwise on test dataset it has 87% accuracy compare to simple decision tree based model has 83%. The forward stepwise used variables that trees liked. I would recommend going with Regression Stepwise model instead of tree based model which compares these two model for this particular dataset.

For one of the Regression Models, print the coefficients. Do the variables make sense? If not, what would you recommend?

- I printed the coefficients for Logistics Regression variables. All the variables mentioned below do make sense
- TRUNC\_IMP\_DEBTINC = 0.1044 (The positive value of Debt Income ratio shows more people have defaulted loans if they have higher Debt Income ratio.)
- M\_VALUE = 3.790 (People who had missing home value which was fixed after really high chance of default loans since it home equity insurance company to trust people if their reasoning of all is Home Improvement or Consolidating debt and value of the house is missing.)
- M\_DEROG = -0.767 (If they lower delinquencies report they less chance of having defaulted loans)
- M\_DEBTINC = 2.751 (Similarly if the debt income ratio value is missing they have more chance of having defaulted loans.)
- TRUNC\_IMP\_YOY = -0.0167 (This results shows that these people have higher values in Year on Job (YOJ) which means they are more stable and assuming less riskier job the probability of having default loans decreases as results shows TRUNC\_IMP\_YOY value close to 0.)
- TRUNC\_IMP\_DEROG, TRUNC\_IMP\_DELINQ, and O\_IMP\_DELINQ. (These variables do make sense as we can see these variables have positive values means that they higher higher derogatory marks on credit records, and outliers in the delinquencies reports might have shown these people have default loans could be because of bankruptcy.
- TRUNC\_IMP\_CLAGE = -0.0064 (Again, I would not recommend this variables based on the results its fluctuating a lot there is almost 50% percent probability of person having either default loans or not. If you have had credit for a long time, you are considered less risky than a new high school student that doesn't mean they cannot have default loans.

## Logistics Regression ROC Curve and Accuracy for ALL Models

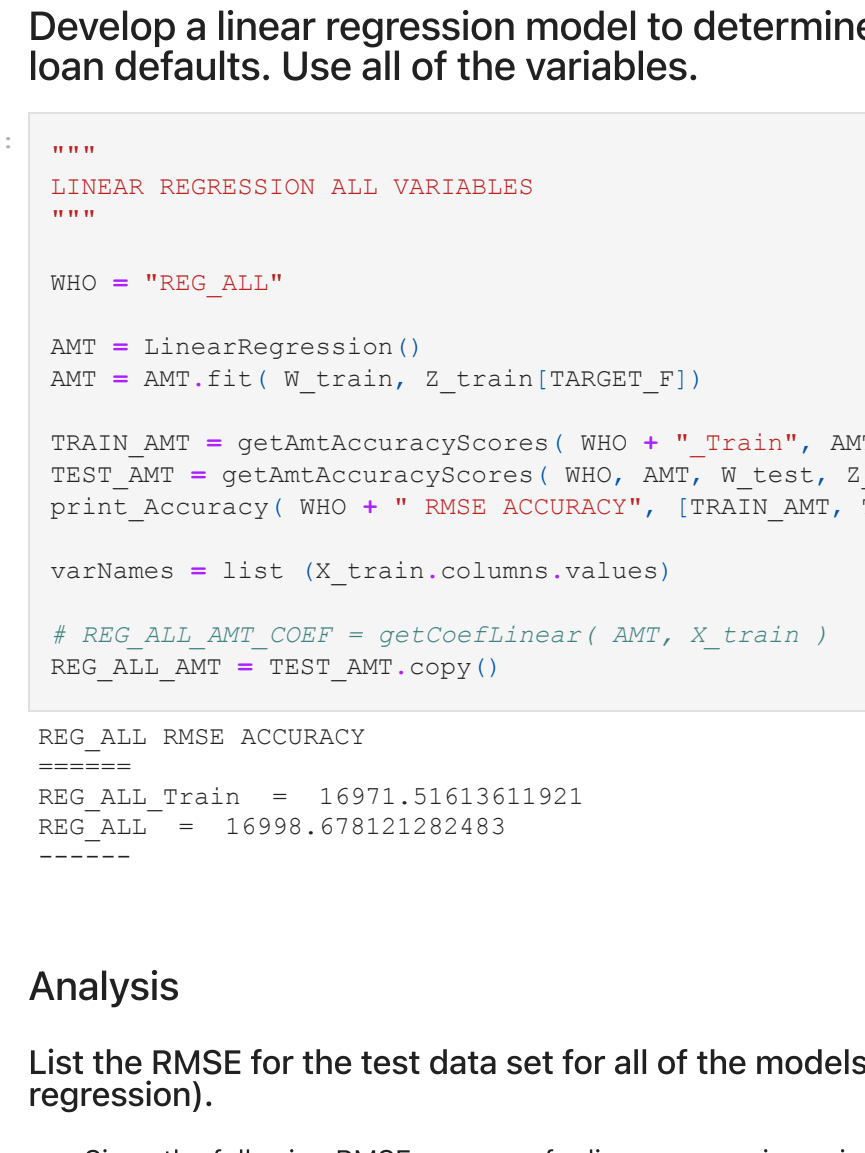
```
In [21]: WHO = "ALL MODELS"

#Classifier Models
#Creating a list of all the models we have created
ALL_CLM = [ TREE_CLM, RF_CLM, GB_CLM, REG_ALL_CLM, REG_TREE_CLM, REG_RF_CLM, REG_GB_CLM ]

# I am going to sort this by 4 element in the list ROC curve
ALL_CLM = sorted( ALL_CLM, key = lambda x: x[4], reverse=True )
print_ROC_Curve( WHO, ALL_CLM )

#Sorting is by the accuracy of the model
ALL_CLM = sorted( ALL_CLM, key = lambda x: x[1], reverse=True )
print_Accuracy( WHO + " ALL CLASSIFICATION ACCURACY", ALL_CLM )

=====
REG_ALL_RMSE ACCURACY
=====
RF    = 0.9144295302013423
GB    = 0.90352348993228859
REG_TREE    = 0.8873898926174496
REG_STEPWISE    = 0.8875838926174496
REG_ALL    = 0.886744966442953
TREE    = 0.889060020648564
REG_GB    = 0.8808724832214765
REG_RF    = 0.87751677852349
-----
```



```
ALL MODELS ALL CLASSIFICATION ACCURACY
=====
RF    = 0.9144295302013423
GB    = 0.90352348993228859
REG_TREE    = 0.8873898926174496
REG_STEPWISE    = 0.8875838926174496
REG_ALL    = 0.886744966442953
TREE    = 0.889060020648564
REG_GB    = 0.8808724832214765
REG_RF    = 0.87751677852349
-----
```

### Analysis:

Display a ROC curve for the test data with all your models on the same graph (tree based and regression).

- Based on the ROC accuracy curve results. The random forest tree model (RF) as highest accuracy compare to all other models. Second is Gradient boosting model (GB) at 93% I have seen these model doing fine but was surprise they are better than regression tree based models. The Regression of All (REG\_ALL) variables is 90% accuracy even though it has higher accuracy its uses all the variables in the dataset for predicting target variables which not efficient and time consuming. I would recommend using Regression Tree Based Models (REG\_TREE) which is 88% accuracy because its faster and simpler model using less variables its fine to loss some accuracy but more efficient than REG\_ALL model.

- Based on the All Classification Accuracy we can see that Random Forest and Gradient Boosting model are on the top in terms of accuracy similar to ROC results. But Regression Tree has very small margin of higher accuracy then regression all variables models which makes sense based on the reasoning mentioned above it uses less variables, not time consuming and simpler model compare to REG\_ALL model. Its surprising to see ROC curve results for simple Decision Tree Based model is lowest at 83% accuracy compare all other models but has higher classification accuracy rate of 88% compare to regression (REG\_GB) gradient boosting model and (REG\_RF) random forest model

## Linear Regression

Develop a linear regression model to determine the expected loss if the loan defaults. Use all of the variables.

```
In [22]: """
LINEAR REGRESSION ALL VARIABLES
"""

WHO = "REG_ALL"

AMT = LinearRegression()
AMT = AMT.fit( X_train, Z_train[TARGET_F] )

TRAIN_AMT = getAmtAccuracyScores( WHO + " Train", AMT, X_train, Z_train[TARGET_A] )
TEST_AMT = getAmtAccuracyScores( WHO, AMT, X_test, Z_test[TARGET_A] )
print_Accuracy( WHO + " RMSE ACCURACY", { TRAIN_AMT, TEST_AMT } )

varNames = list( X_train.columns.values )

# REG_ALL_AMT_COEF = getCoefLinear( AMT, X_train )
REG_ALL_AMT = TEST_AMT.copy()

=====
REG_ALL_RMSE ACCURACY
=====
REG_ALL_Train    = 16971.51613611921
REG_ALL    = 16998.678121282483
-----
```

### Analysis

List the RMSE for the test data set for all of the models created (tree based and regression).

- Since the following RMSE accuracy for linear regression using all the variables measuring the difference between values predicted by a model and their actual values for loss amount loan not repaid of 16,000 dollars seems to be a lot higher compare to tree based models. We need to have our models to be less complicated as simple as possible. But its decent to start off by using all variables in the regression model. I would only recommend using REG\_ALL variables for regression model as a starting point to see the RMSE accuracy and how model behaves when we use all variables.

Develop a linear regression model to determine the expected loss if the loan defaults. Use the variables that were selected by a DECISION TREE.

```
In [23]: WHO = "REG_TREE"

AMT = LinearRegression()
AMT = AMT.fit( X_train[vars_tree_amt], Z_train[TARGET_A] )

TRAIN_AMT = getAmtAccuracyScores( WHO + " Train", AMT, X_train[vars_tree_amt], Z_train[TARGET_A] )
TEST_AMT = getAmtAccuracyScores( WHO, AMT, X_test[vars_tree_amt], Z_test[TARGET_A] )
print_Accuracy( WHO + " RMSE ACCURACY", { TRAIN_AMT, TEST_AMT } )

varNames = list( X_train.columns.values )

# REG_TREE_AMT_COEF = getCoefLinear( AMT, X_train[vars_tree_amt] )
REG_TREE_AMT = TEST_AMT.copy()

=====
REG_TREE_RMSE ACCURACY
=====
REG_TREE_Train    = 4307.479817092925
REG_TREE    = 4301.963358361646
-----
```

### Analysis:

List the RMSE for the test data set for all of the models created (tree based and regression).

- The average RMSE score for loss amount not repaid with Linear Regression Decision Tree (REG\_TREE) is 4301 dollars of for test dataset is much better than RMSE score for loss amount loan not repaid of 16,000 dollars there is almost 2000 dollars difference between RMSE score of each model. They both are using same variables for build models which trees liked. I would recommend going with Regression Tree instead of tree based models because its smaller loan amount not repaid to the bank assume loans defaulted for this particular dataset.

Develop a linear regression model to determine the expected loss if the loan defaults. Use the variables that were selected by a RANDOM FOREST.

```
In [24]: """
LINEAR REGRESSION RANDOM FOREST
"""

WHO = "REG_RF"

print("\n\n")
RF_amt = {}
for i in vars_RF_amt :
    # print(i)
    theVar = i[0]
    RF_amt.append( theVar )

AMT = LinearRegression()
AMT = AMT.fit( X_train[RF_amt], Z_train[TARGET_A] )

TRAIN_AMT = getAmtAccuracyScores( WHO + " Train", AMT, X_train[RF_amt], Z_train[TARGET_A] )
TEST_AMT = getAmtAccuracyScores( WHO, AMT, X_test[RF_amt], Z_test[TARGET_A] )
print_Accuracy( WHO + " RMSE ACCURACY", { TRAIN_AMT, TEST_AMT } )

# REG_RF_AMT_COEF = getCoefLinear( AMT, X_train[RF_amt] )
REG_RF_AMT = TEST_AMT.copy()

=====
REG_RF_RMSE ACCURACY
=====
REG_RF_Train    = 4384.862092423268
REG_RF    = 4360.610292020044
-----
```

### Analysis:

List the RMSE for the test data set for all of the models created (tree based and regression).

- The average RMSE score for loss amount not repaid with Linear Regression Gradient Boosting (REG\_TREE) is 4301 dollars of for test dataset is much better than RMSE score for loss amount loan not repaid of 16,000 dollars there is almost 2000 dollars difference between RMSE score of each model. They both are using same variables for build models which Gradient Boosting model liked. I would recommend going with Gradient Boosting tree based model instead of Regression Random Forest model because its smaller loan amount not repaid to the bank for this particular dataset.

Develop a linear regression model to determine the expected loss if the loan defaults. Use the variables that were selected by a GRADIENT BOOSTING model.

```
In [25]: """
LINEAR REGRESSION GRADIENT BOOSTING
"""

WHO = "REG_GB"

print("\n\n")
GB_amt = {}
for i in vars_GB_amt :
    # print(i)
    theVar = i[0]
    GB_amt.append( theVar )

AMT = LinearRegression()
AMT = AMT.fit( X_train[GB_amt], Z_train[TARGET_A] )

TRAIN_AMT = getAmtAccuracyScores( WHO + " Train", AMT, X_train[GB_amt], Z_train[TARGET_A] )
TEST_AMT = getAmtAccuracyScores( WHO, AMT, X_test[GB_amt], Z_test[TARGET_A] )
print_Accuracy( WHO + " RMSE ACCURACY", { TRAIN_AMT, TEST_AMT } )

# REG_GB_AMT_COEF = getCoefLinear( AMT, X_train[GB_amt] )
REG_GB_AMT = TEST_AMT.copy()

=====
REG_GB_RMSE ACCURACY
=====
REG_GB_Train    = 4043.551329405321
REG_GB    = 4029.7840135710135
-----
```

### Analysis:

List the RMSE for the test data set for all of the models created (tree based and regression).

- The average RMSE score for loss amount not repaid with Linear Regression Gradient Boosting (REG\_TREE) is 4029 dollars of for test dataset is much higher than average RMSE score for Random Forest tree based models (RF) is 2201 dollars there is almost 2000 dollars difference between RMSE score of each model. They both are using same variables for build models which Gradient Boosting model liked. I would recommend going with Gradient Boosting tree based model instead of Regression Random Forest model because its smaller loan amount not repaid to the bank for this particular dataset. Also, Regression Stepwise would be larger and less efficient in terms of variables selection technique than Gradient Boosting model for this dataset.

Develop a linear regression model to determine the expected loss if the loan defaults. Use the variables that were selected by STEPWISE SELECTION.

```
In [26]: # Copy train data into new dataset called V_train based on the variables what Gradient Boosting tree based model using 1 to as many variable needed
# This will be a little bit faster since we are using linear regression

V_train = X_train[ GB_amt ]
stepVarNames = list( V_train.columns.values )
maxCols = V_train.shape[1]

sfs = SFS( LinearRegression(),
           k_features=(1, maxCols),
           forward=True,
           floating=False,
           scoring='r2',
           cv=5
           )

sfs.fit(V_train.values, Z_train[ TARGET_A ].values)

theFigure = plot_sfs(sfs.get_metric_dict(), kind=None )
plt.title('Loss Amount Sequential Forward Selection (w. StdErr)')
plt.grid()
plt.show()

dfm = pd.DataFrame.from_dict( sfs.get_metric_dict() ).T
dfm = dfm[ ['feature_names', 'avg_score' ] ]
dfm.avg_score = dfm.avg_score.astype(float)

print(" .....")
maxIndex = dfm.avg_score.argmax()
print("argmax")
print( dfm.iloc[ maxIndex, ] )
print(" .....")

stepVars = dfm.iloc[ maxIndex, ]
stepVars = stepVars.feature_names
print( stepVars )

finalStepVars = []
for i in stepVars :
    try :
        theName = stepVarNames[ index ]
        finalStepVars.append( theName )
    except :
        pass

for i in finalStepVars :
    print( i )

V_train = X_train[ finalStepVars ]
V_test = X_test[ finalStepVars ]

=====
Loss Amount Sequential Forward Selection (w. StdErr)
=====
Performance    0.650    0.700    0.725    0.750    0.775    0.800    0.825
Number of Features    1        2        3        4        5

.....
argmax
feature_names    (0, 1, 2, 3, 4)
avg_score    0.834263
Name: 5, dtype: object
.....
('0', '1', '2', '3', '4')
TRUNC_LOAN
TRUNC_IMP_CLNO
TRUNC_IMP_DEBTINC
M_DEBTINC
TRUNC_IMP_CLAGE
```

```
In [27]: """
LINEAR REGRESSION STEPWISE
"""

WHO = "REG_STEPWISE"

AMT = LinearRegression()
AMT = AMT.fit( V_train, Z_train[TARGET_A] )

TRAIN_AMT = getAmtAccuracyScores( WHO + " Train", AMT, V_train, Z_train[TARGET_A] )
TEST_AMT = getAmtAccuracyScores( WHO, AMT, V_test, Z_test[TARGET_A] )
print_Accuracy( WHO + " RMSE ACCURACY", { TRAIN_AMT, TEST_AMT } )

REG_STEP_AMT_COEF = getCoefLinear( AMT, V_train )
REG_STEP_AMT = TEST_AMT.copy()

=====
REG_STEPWISE_RMSE ACCURACY
=====
REG_STEPWISE_Train    = 4043.551329405321
REG_STEPWISE    = 4029.7840135710135
-----
```

LOSS AMOUNT

```
-----
Total Variables: 6
INTERCEPT    = -12990.102998292584
TRUNC_LOAN    = 0.7811032153671928
TRUNC_IMP_CLNO    = 296.780759504102
TRUNC_IMP_DEBTINC    = 193.84552770224016
W_DEBTINC    = 5766.133648473953
TRUNC_IMP_CLAGE    = -25.06345014576534
```

### Analysis:

List the RMSE for the test data set for all of the models created (tree based and regression).

- The average RMSE score for loss amount not repaid with Linear Regression Stepwise Function (REG\_STEPWISE) is 4029 dollars using test dataset is much higher than average RMSE score for Gradient Boosting tree based models (GB) is 2201 dollars there is almost 2000 dollars difference between RMSE score of each model. They both are using same variables for build models which Gradient Boosting model liked. I would recommend going with Gradient Boosting tree based model instead of Regression Stepwise model because its smaller loan amount not repaid to the bank for this particular dataset. Also, Regression Stepwise would be larger and less efficient in terms of variables selection technique than Gradient Boosting model for this dataset.

For one of the Regression Models, print the coefficients. Do the variables make sense? If not, what would you recommend?

- I printed the coefficients for Linear Regression variables Step wise function. All the variables mentioned below do make sense.
- TRUNC\_LOAN = 0.7811 (The positive value of LOAN amounts shows more people have higher amount of loan not repaid if they have bigger amount more risky the person which makes sense based on the criteria.)
- TRUNC\_IMP\_CLNO = 296.78 (Showing a really positive value means a lot of people have higher number of credit lines, and this variables could have gone either scenario either person is more trustworthy. But in our case based on the model a too many credit lines means that they have potential to run up a lot of debt causing higher loss amount not repaid to the bank)
- TRUNC\_IMP\_DEBTINC = 193.84 (The higher Debt Income Ratio means that people are loss amount increase if person has higher Debt Income Ratio)
- M\_DEBTINC = 5766 (Similarly if the debt income ratio value is missing which was fixed have higher loans amounts because they are trustworthy compare to imputed debtcinc ratio kind of doesn't make sense.)
- TRUNC\_IMP\_CLAGE = -25.06 (In our dataset looks like most of people had credit for a long time, considered less risky showing smaller amount of loan not repaid )

## Linear Regression RMSE Accuracy for ALL Models

```
In [28]: #Linear Regression Models
#Creating a list of all the models we have created
ALL_AMT = [ TREE_AMT, RF_AMT, GB_AMT, REG_ALL_AMT, REG_TREE_AMT, REG_RF_AMT, REG_GB_AMT ]

#Sort by Accuracy of the Damages
ALL_AMT = sorted( ALL_AMT, key = lambda x: x[1] )

print_Accuracy( "ALL LOSS AMOUNT MODEL ACCURACY", ALL_AMT )

=====
ALL LOSS AMOUNT MODEL ACCURACY
=====
GB    = 2201.819449335387
REG_GB    = 2780.278461970596
REG_GB    = 4029.7840135710135
REG_STEPWISE    = 4029.7840135710135
W_DEBTINC    = 4301.963358361646
REG_RF    = 4360.610292020044
TREE    = 5522.483680969053
REG_ALL    = 16998.678121282483
-----
```

### Analysis:

List the RMSE for the test data set for all of the models created (tree based and regression).

- Based on the RMSE accuracy curve results measuring the difference between values predicted by a model and their actual values for loss amount loan Gradient Boosting tree model (GB) as lowest average RMSE score of 2201 dollars loss amount not repaid compare to all other models. Second is Random Forest model (RF) with Root Mean Square error of 2780 dollars loss amount. I have seen these model doing fine but was surprise they are better than regression based models. The Regression of All (REG\_ALL) variables is larger difference in terms of RMSE accuracy of 16000 dollars performed really bad was surprising even though used all the variables in the dataset for predicting target variables which not efficient and time consuming. Based on the All Loss Model Accuracy we can that Gradient Boosting and Random Forest model are on the top in terms of accuracy as they have closer range of RMSE scores.

## Summary Report

Include a discussion of the which models were most accurate, and which ones would you recommend using in a real world situation.

- The random forest tree model (RF) considerably more accurate with 96% covering area under the curve compare to all other models. Second is Gradient boosting model (GB) is at 93% accuracy is a little bit closer to Gradient Boosting showing ensemble approach has worked surprisingly better than regression based models in predicting loans defaults. I would recommend using Tree Based to start off simple understanding how the input variables is having an effect on predicting target variables. And the REG\_ALL variables model for regression model as a starting point as well to see the accuracy and how model behaves when we use all variables for real world situation.

If you were to select one of these models to put into production, which would it be? Why would you select this model?

- I would Random Forest tree based model to put into production for this dataset to predict both default loans and loss amount if loan was not repaid. The reason for selecting this model based on ROC curve and RMSE accuracy results for Random Forest and Gradient was almost similar high accuracy compare to other models. My second choice would have been Gradient Boosting because as int tends to build very shallow trees but we have to look at is it legal to use based business rules before deploying in production

## Logistics Regression

Display a ROC curve for the test data with all your models on the same graph (tree based and regression).

I have compared ROC curve accuracy results between both Logistics Regression and tree based after each models under analysis section.

- Display a ROC curve for the test data with all your models on the same graph (tree based and regression). Based on the ROC accuracy curve results random forest tree model (RF) as highest accuracy compare to all other models. Second is Gradient boosting model (GB) at 93% I have seen these model doing fine but was surprise they are better than regression tree based models. The Regression of All (REG\_ALL) variables is 90% accuracy even though it has higher accuracy its uses all the variables in the dataset for predicting target variables which not efficient and time consuming. I would recommend using Regression Tree Based Models (REG\_TREE) which is 88% accuracy because its faster and simpler model using less variables its fine to loss some accuracy but more efficient than REG\_ALL model.

- Based on the All Classification Accuracy we can see that Random Forest and Gradient Boosting model are on the top in terms of accuracy similar to ROC results. But Regression Tree has very small margin of higher accuracy then regression all variables models which makes sense based on the reasoning mentioned above it uses less variables, not time consuming and simpler model compare to REG\_ALL model. Its surprising to see ROC curve results for simple Decision Tree Based model is lowest at 83% accuracy compare all other models but has higher classification accuracy rate of 88% compare to regression (REG\_GB) gradient boosting model and (REG\_RF) random forest model

For one of the Regression Models, print the coefficients. Do the variables make sense? If not, what would you recommend?

- I printed the coefficients variables for Stepwise selection logistics regression model. All the variables mentioned below do make sense
- TRUNC\_LOAN = 0.7811 (The positive value of LOAN amounts shows more people have higher amount of loan not repaid if they have bigger amount more risky the person which makes sense based on the criteria.)
- TRUNC\_IMP\_CLNO = 296.78 (Showing a really positive value means a lot of people have higher number of credit lines, and this variables could have gone either scenario either person is more trustworthy. But in our case based on the model a too many credit lines means that they have potential to run up a lot of debt causing higher loss amount not repaid to the bank)
- TRUNC\_IMP\_DEBTINC = 193.84 (The higher Debt Income Ratio means that people are loss amount increase if person has higher Debt Income Ratio)
- M\_DEBTINC = 5766 (Similarly if the debt income ratio value is missing which was fixed have higher loans amounts because they are trustworthy compare to imputed debtcinc ratio kind of doesn't make sense.)
- TRUNC\_IMP\_CLAGE = -25.06 (In our dataset looks like most of people had credit for a long time, considered less risky showing smaller amount of loan not repaid )

## BINGO BONUS ON NEXT PAGE

Suggestions: Explore different parameters for the Logistic and Linear Regression models. Briefly discuss whether or not they had an effect on results

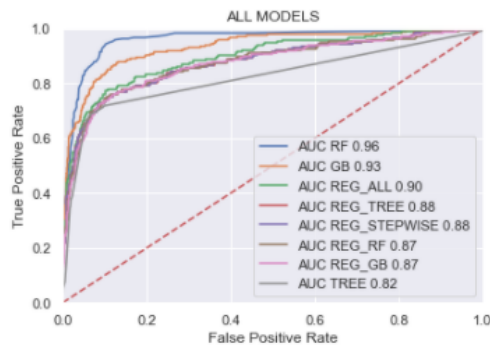
Logistics Regression, we will be trying different 'solver' parameters to see how they affect the results.

In [ ] :



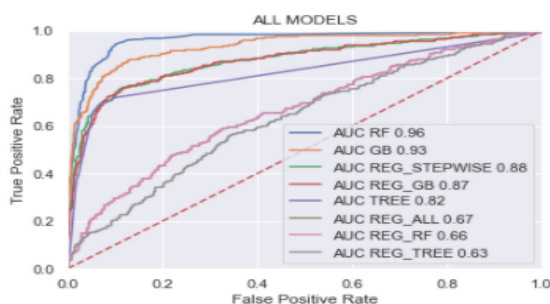
**BINGO BONUS: Suggestions: Explore different parameters for the Logistic and Linear Regression models. Briefly discuss whether they had an effect on results**

**Logistics Regression, we will be trying different 'solver' parameters to see how they affect the results.** We tried **solver= 'newton-cg'** for this Assignment it turned out to have real good accurate scores the model based on ROC curve for this dataset. The newton method uses an exact Hessian matrix. It slows for large datasets, its accurate for this dataset because its computing second derivative.



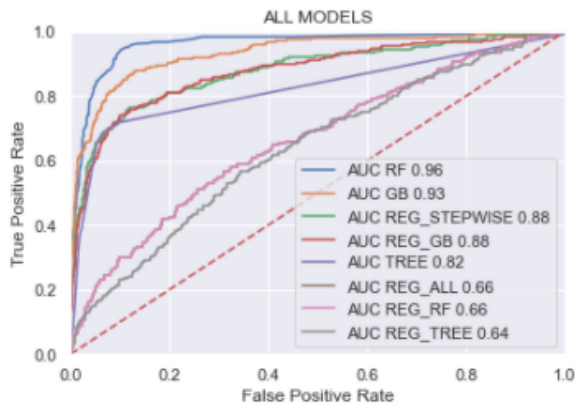
```
ALL MODELS CLASSIFICATION ACCURACY
=====
RF      = 0.9144295302013423
GB      = 0.9035234899328859
REG_TREE = 0.8875838926174496
REG_STEPWISE = 0.8875838926174496
REG_ALL  = 0.886744966442953
TREE     = 0.8859060402684564
REG_GB   = 0.8808724832214765
REG_RF   = 0.87751677852349
-----
```

Second **solver= 'lbfgs'** It approximates second derivative matrix updates gradient evaluations. Its stores in memory but works fast for large dataset. It affected the accuracy results for the Logistics Regression Random Forest (REG\_RF), All Variables (REG\_ALL) and Tree (REG\_TREE) models dropping accuracy by more than 20% compared to 'newton-cg' solver. It can only use one-vs.-rest to solve multi-class problems. It also penalizes the intercept, which isn't good for interpretation.



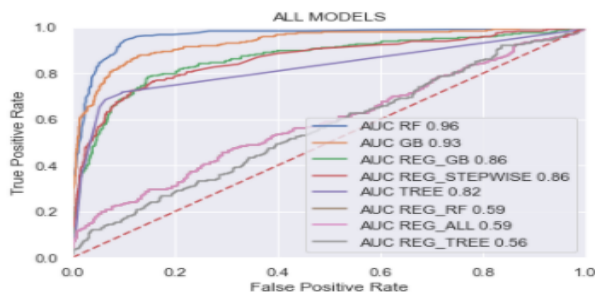
```
ALL MODELS CLASSIFICATION ACCURACY
=====
RF      = 0.9144295302013423
GB      = 0.9035234899328859
REG_STEPWISE = 0.8859060402684564
TREE     = 0.8859060402684564
REG_GB   = 0.8808724832214765
REG_ALL  = 0.7919463087248322
REG_RF   = 0.7919463087248322
REG_TREE = 0.7919463087248322
-----
```

Third **solver** = 'liblinear' is for large linear classification Uses a coordinate descent algorithm. Coordinate descent is based on minimizing a multivariate function by solving univariate optimization problems in a loop. In other words, it moves toward the minimum in one direction at a time. The ROC curve results using liblinear solver are almost similar to lbfgs solver. It did not work well for this dataset.



```
ALL MODELS CLASSIFICATION ACCURACY
=====
RF      = 0.9144295302013423
GB      = 0.9035234899328859
REG_STEPWISE = 0.8859060402684564
TREE    = 0.8859060402684564
REG_GB   = 0.8783557046979866
REG_ALL  = 0.7927852348993288
REG_RF   = 0.7927852348993288
REG_TREE = 0.7927852348993288
-----
```

Fourth **solver** = 'saga' stands for stochastic average gradient descent. A variation of gradient descent and incremental aggregated gradient approaches that uses a random sample of previous gradient values and also allows for L1 regularization trains models faster than 'sag' but it affected the results a lot it had the worst ROC curve accuracy for this dataset compare to other solver algorithm.



```
ALL MODELS CLASSIFICATION ACCURACY
=====
RF      = 0.9144295302013423
GB      = 0.9035234899328859
TREE    = 0.8859060402684564
REG_STEPWISE = 0.8682885906040269
REG_GB   = 0.8557046979865772
REG_RF   = 0.7919463087248322
REG_ALL  = 0.7919463087248322
REG_TREE = 0.7919463087248322
-----
```

The following snippet code is a good a way to see which solver performs better with our test dataset. As we can we **'newton-cg' solver** has good accuracy compared to other solvers which saw above in our ROC curve results as well.

```
In [67]: from sklearn.model_selection import GridSearchCV
solver_list = ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']
parameters = dict(solver=solver_list)
lr = LogisticRegression(random_state=1)
clf = GridSearchCV(lr, parameters, cv=5)
clf.fit(X_test, Y_test[TARGET_F ])
scores = clf.cv_results_['mean_test_score']
```

```
for score, solver, in zip(scores, solver_list):
    print(f"{solver}: {score:.3f}")
```

```
newton-cg: 0.884
lbfgs: 0.797
liblinear: 0.791
sag: 0.792
saga: 0.792
```