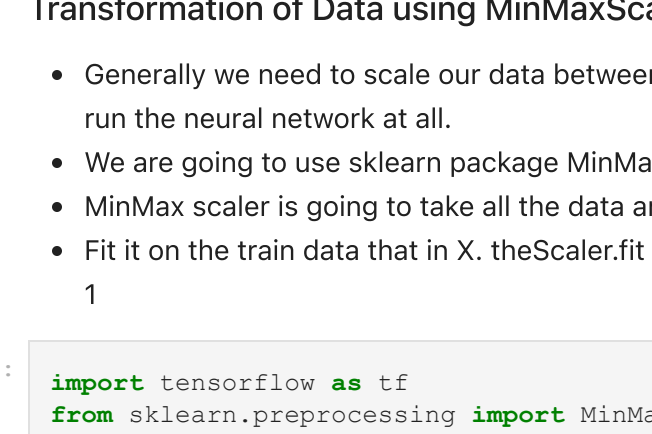


```
.....
argmax      (0, 1, 2, 3, 4)
avg_score   0.834263
Name: 5, dtype: object
('0', '1', '2', '3', '4')
TRUNC_LOAN  2.744553588137954
TRUNC_IMP_CLNO  2.046185940190754
TRUNC_IMP_DEBTINC  193.8455277024016
M_DEBTINC    5766.1336424739575
TRUNC_IMP_CLAGE  -0.0068396349908037115
REG_STEPISE_RMSE_ACCURACY
=====
REG_STEPISE_Train    = 0.885963087248322
REG_STEPISE         = 0.78842281979194631
=====
```



```
REG_STEPISE CLASSIFICATION ACCURACY
=====
REG_STEPISE_Train    = 0.885963087248322
REG_STEPISE         = 0.78842281979194631
=====
```

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

```
CRASH
-----
Total Variables: 10
INTERCEPT    = -3.357029791039792
TRUNC_IMP_DEBTINC  = 0.10930606571813604
M_VALUE        = 3.7092594452274983
M_DEROG        = -0.7853106477725912
W_DEBTINC      = 2.744553588137954
TRUNC_IMP_VALUE = -7.725928819978311e-07
TRUNC_IMP_DEROG = 0.6924412080601994
IMP_DELIQ      = 2.046185940190754
TRUNC_IMP_DELIQ = 0.6177869545230693
TRUNC_IMP_CLAGE = -0.0068396349908037115
REG_STEPISE_RMSE_ACCURACY
=====
REG_STEPISE_Train    = 4043.5513294053221
REG_STEPISE         = 4029.784013571014
=====
```

DAMAGES

```
-----
Total Variables: 6
INTERCEPT    = -12990.102898292584
TRUNC_LOAN     = 0.7811032153671928
TRUNC_IMP_CLNO = 296.7807595041021
TRUNC_IMP_DEBTINC = 193.8455277024016
M_DEBTINC      = 5766.1336424739575
TRUNC_IMP_CLAGE = -25.06345014576536
```

## Tensorflow Model

### Develop a model using Tensor Flow that will predict Loan Default.

#### Transformation of Data using MinMaxScaler

- Generally we need to scale our data between 0 and 1. Otherwise it takes longer to run or it wont run the neural network at all.
- We are going to use sklearn package MinMaxScaler over here
- MinMax scaler is going to take all the data and turn it between 0 and 1
- Fit it on the train data that in X. theScaler.fit transforms the new data and arrange it between 0 and 1

```
In [6]: import tensorflow as tf
from sklearn.preprocessing import MinMaxScaler
theScaler = MinMaxScaler()
theScaler.fit(X_train)
```

```
Out[6]: MinMaxScaler()
```

- First we are gonna do is build a model that predicts if the person has loans defaults or not
- We are going to create new variables U\_train and U\_test, so that theScaler transforms all the data between 0 and 1 for the variables X\_train and X\_test.
- We are going to convert numpy array U\_train and U\_test into pd.DataFrame
- Currently all the names for U\_train and U\_test are 0,1,2,3.. so that we are going to get the names from X\_train.columns.values and X\_test.columns.values put it into list and U\_train and U\_test

```
In [7]: WHO = "Tensor Flow"
U_train = theScaler.transform(X_train)
U_test = theScaler.transform(X_test)

U_train = pd.DataFrame(U_train)
U_test = pd.DataFrame(U_test)

U_train.columns = list(X_train.columns.values)
U_test.columns = list(X_test.columns.values)
```

#### Explore using a variable selection technique

- We are going select from the variables which Random Forest tree based model liked since it has performed well in terms overall accuracy for this particular dataset compare to other models. We are going to use those variables in the U\_train and U\_test

```
In [8]: U_train = U_train[RF_flag]
U_test = U_test[RF_flag]
U_test.head()
```

	0	1	2	3	4
TRUNC_IMP_DEBTINC	0.000000	0.000000	0.000000	0.000000	0.000000
TRUNC_IMP_DEBTINC	0.663835	0.755929	0.663555	0.523003	0.578634
TRUNC_IMP_CLAGE	0.385118	0.732522	0.271465	0.815440	0.271875
TRUNC_IMP_VALUE	0.464032	0.669815	0.197667	0.305699	0.509862
TRUNC_LOAN	0.400939	0.397027	0.608253	0.162331	1.000000
TRUNC_IMP_MORTDUE	0.769490	0.709614	0.068739	0.375530	0.327055
TRUNC_IMP_DELIQ	1.000000	0.000000	0.000000	0.000000	0.250000
TRUNC_IMP_CLNO	0.843137	0.333333	0.098039	0.901961	0.215686
TRUNC_IMP_YOJ	0.225806	0.225806	0.000000	0.774194	0.290323
TRUNC_IMP_DEROG	0.666667	0.000000	0.000000	0.000000	0.000000
TRUNC_IMP_NINQ	0.000000	0.000000	0.166667	0.000000	0.000000

## Tensorflow Model Accuracy Metrics Function

```
In [9]: def get_TP_ProbAccuracyScores( NAME, MODEL, X, Y ):
    probs = MODEL.predict( X ) #getting the probability scores
    pred_list = [] #printing an empty list
    for p in probs: #getting probability for both yes and no if they have loan default
        pred_list.append(np.argmax( p ))
    pred = np.array( pred_list )
    acc_score = metrics.accuracy_score(Y, pred)
    pl = probs[:,1] #we are getting only for people who have loans defaulted
    fpr, tpr, threshold = metrics.roc_curve( Y, pl)
    auc = metrics.auc(fpr,tpr)
    return (NAME, acc_score, fpr, tpr, auc)
```

## Try at least three different Activation Functions

### First Activation Function used RELU

- 1. Try one and two hidden layers
- 1. Try using a Dropout Layer

- F\_ is for the flag yes or no are if the person has default loans
- F\_ theShapeSize is we would like to know what's the size of the dataset.
- F\_ theLossMetric is SparseCategoricalCrossentropy() because this data as 0 or 1
- F\_ theOptimizer is Adam() function
- F\_ theEpochs go through the data 100 times
- F\_ theUnits is how many nodes should I have in my theUnits. I started off with 2 times the shape of the dataset lets do this as our starting off point which has better accuracy based on ROC curve and compare to when I started removing nodes. So, I decided to stick with 2 times the shape of dataset for the nodes used in the neural network.
- F\_ LAYER\_01 is a Dense layer
- F\_ LAYER\_DROP and F\_ LAYER\_02 without input his time. Adding a Dropout layer as well where we are simply saying everytime we run through the iteration throw away 20% of the nodes it will throw away 20% of the nodes from whatever nodes were in LAYER\_01 called before it this would prevent it from overfitting the model
- F\_ LAYER\_OUTPUT has 2 inputs one for YES and one for NO. The activation function used for categorical variables is softmax

- Create a model called CLM (CLAIM) thats a called Sequential()
- CLM.compile means (The compile simply say this is the loss function we are gonna use to train the neural network and this is the optimizer that is gonna adjust the weights)
- CLM.fit finally model.fit ( Fit the neural network with X(input variables), Y(output variables) and arguments we want)

```
In [10]: F_theShapeSize = U_train.shape[1]
F_theActivation = tf.keras.activations.relu
F_theLossMetric = tf.keras.losses.SparseCategoricalCrossentropy()
F_theOptimizer = tf.keras.optimizers.Adam()
F_theEpochs = 100

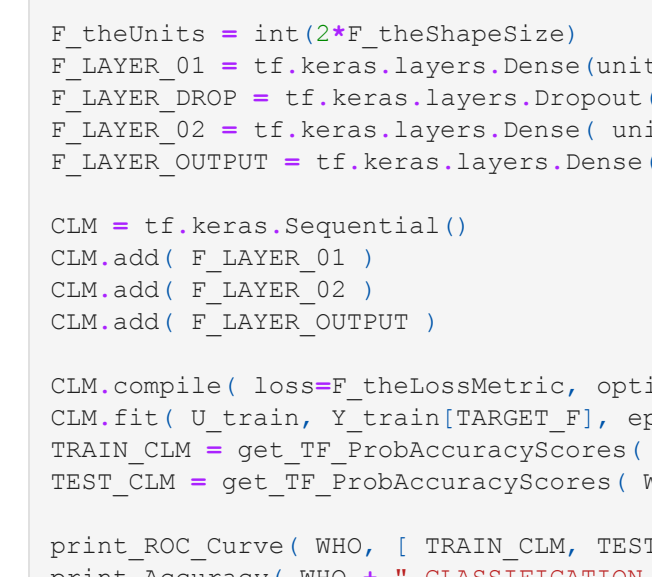
F_theUnits = int(2*F_theShapeSize)
F_LAYER_01 = tf.keras.layers.Dense(units=F_theUnits, activation = F_theActivation, input_shape=(F_LAYER_DROP,))
F_LAYER_DROP = tf.keras.layers.Dropout( 0.2 )
F_LAYER_02 = tf.keras.layers.Dense(units=F_theUnits, activation = F_theActivation)
F_LAYER_OUTPUT = tf.keras.layers.Dense(units=2, activation = tf.keras.activations.softmax)
```

```
CLM = tf.keras.Sequential()
CLM.add( F_LAYER_01 )
CLM.add( F_LAYER_02 )
CLM.add( F_LAYER_OUTPUT )

CLM.compile( loss=F_theLossMetric, optimizer=F_theOptimizer)
CLM.fit( U_train, Y_train[TARGET_F], epochs=F_theEpochs, verbose=False )
TRAIN_CLM = get_TP_ProbAccuracyScores( WHO + " Train", CLM, U_train, Y_train[ TARGET_F ])
TEST_CLM = get_TP_ProbAccuracyScores( WHO, CLM, U_test, Y_test[ TARGET_F ] )

print_ROC_Curve( WHO, [ TRAIN_CLM, TEST_CLM ] )
print_Accuracy( WHO + " CLASSIFICATION ACCURACY", [ TRAIN_CLM, TEST_CLM ] )

TF_CLM = TEST_CLM.copy()
TF_AMT = TEST_AMT.copy()
```



```
Tensor_Flow CLASSIFICATION ACCURACY
=====
Tensor_Flow_Train    = 0.902894295302134
Tensor_Flow         = 0.899261744966443
=====
```

## Try at least three different Activation Functions

### Second Activation Function used tanh (Hyperbolic function)

- 1. Try one and two hidden layers
- 1. Try using a Dropout Layer

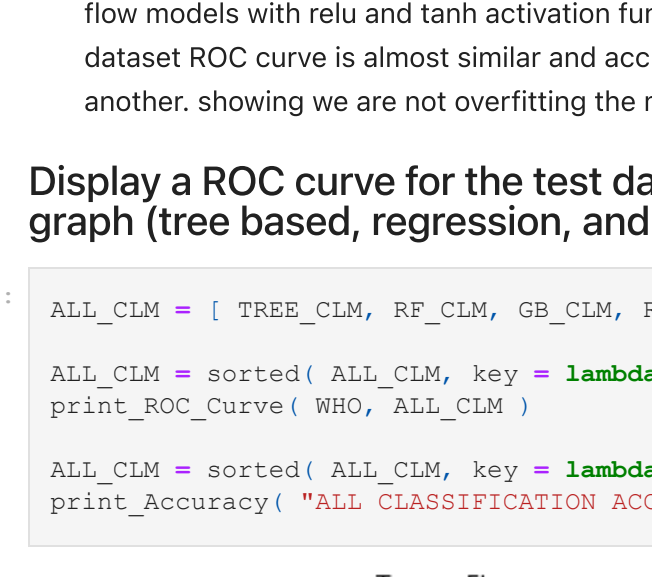
```
In [11]: F_theShapeSize = U_train.shape[1]
F_theActivation = tf.keras.activations.tanh
F_theLossMetric = tf.keras.losses.SparseCategoricalCrossentropy()
F_theOptimizer = tf.keras.optimizers.Adam()
F_theEpochs = 100

F_theUnits = int(2*F_theShapeSize)
F_LAYER_01 = tf.keras.layers.Dense(units=F_theUnits, activation = F_theActivation, input_shape=(F_LAYER_DROP,))
F_LAYER_DROP = tf.keras.layers.Dropout( 0.2 )
F_LAYER_02 = tf.keras.layers.Dense(units=F_theUnits, activation = F_theActivation)
F_LAYER_OUTPUT = tf.keras.layers.Dense(units=2, activation = tf.keras.activations.softmax)
```

```
CLM = tf.keras.Sequential()
CLM.add( F_LAYER_01 )
CLM.add( F_LAYER_02 )
CLM.add( F_LAYER_OUTPUT )

CLM.compile( loss=F_theLossMetric, optimizer=F_theOptimizer)
CLM.fit( U_train, Y_train[TARGET_F], epochs=F_theEpochs, verbose=False )
TRAIN_CLM = get_TP_ProbAccuracyScores( WHO + " Train", CLM, U_train, Y_train[ TARGET_F ])
TEST_CLM = get_TP_ProbAccuracyScores( WHO, CLM, U_test, Y_test[ TARGET_F ] )

print_ROC_Curve( WHO, [ TRAIN_CLM, TEST_CLM ] )
print_Accuracy( WHO + " CLASSIFICATION ACCURACY", [ TRAIN_CLM, TEST_CLM ] )
```



```
Tensor_Flow CLASSIFICATION ACCURACY
=====
Tensor_Flow_Train    = 0.8955336912751678
Tensor_Flow         = 0.8875838926174496
=====
```

## Try at least three different Activation Functions

### Third Activation Function used elu (Exponential Linear Unit)

- 1. Try one and two hidden layers
- 1. Try using a Dropout Layer

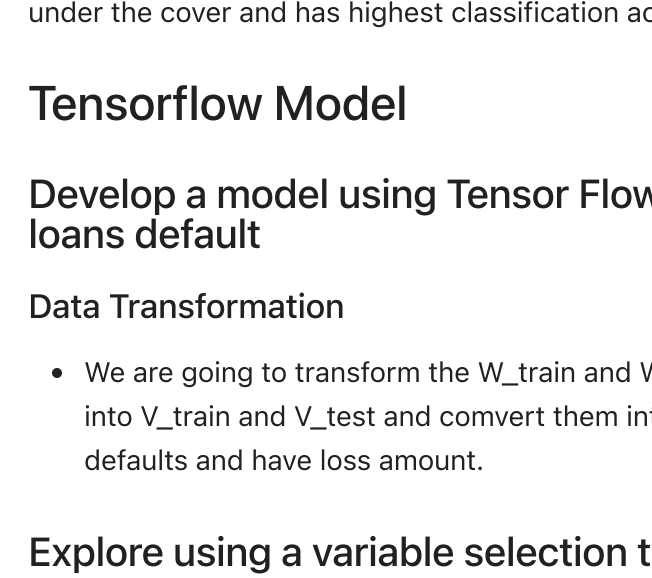
```
In [12]: F_theShapeSize = U_train.shape[1]
F_theActivation = tf.keras.activations.elu
F_theLossMetric = tf.keras.losses.SparseCategoricalCrossentropy()
F_theOptimizer = tf.keras.optimizers.Adam()
F_theEpochs = 100

F_theUnits = int(2*F_theShapeSize)
F_LAYER_01 = tf.keras.layers.Dense(units=F_theUnits, activation = F_theActivation, input_shape=(F_LAYER_DROP,))
F_LAYER_DROP = tf.keras.layers.Dropout( 0.2 )
F_LAYER_02 = tf.keras.layers.Dense(units=F_theUnits, activation = F_theActivation)
F_LAYER_OUTPUT = tf.keras.layers.Dense(units=2, activation = tf.keras.activations.softmax)
```

```
CLM = tf.keras.Sequential()
CLM.add( F_LAYER_01 )
CLM.add( F_LAYER_02 )
CLM.add( F_LAYER_OUTPUT )

CLM.compile( loss=F_theLossMetric, optimizer=F_theOptimizer)
CLM.fit( U_train, Y_train[TARGET_F], epochs=F_theEpochs, verbose=False )
TRAIN_CLM = get_TP_ProbAccuracyScores( WHO + " Train", CLM, U_train, Y_train[ TARGET_F ])
TEST_CLM = get_TP_ProbAccuracyScores( WHO, CLM, U_test, Y_test[ TARGET_F ] )

print_ROC_Curve( WHO, [ TRAIN_CLM, TEST_CLM ] )
print_Accuracy( WHO + " CLASSIFICATION ACCURACY", [ TRAIN_CLM, TEST_CLM ] )
```



```
Tensor_Flow CLASSIFICATION ACCURACY
=====
Tensor_Flow_Train    = 0.8888422818791947
Tensor_Flow         = 0.8657718120805369
=====
```

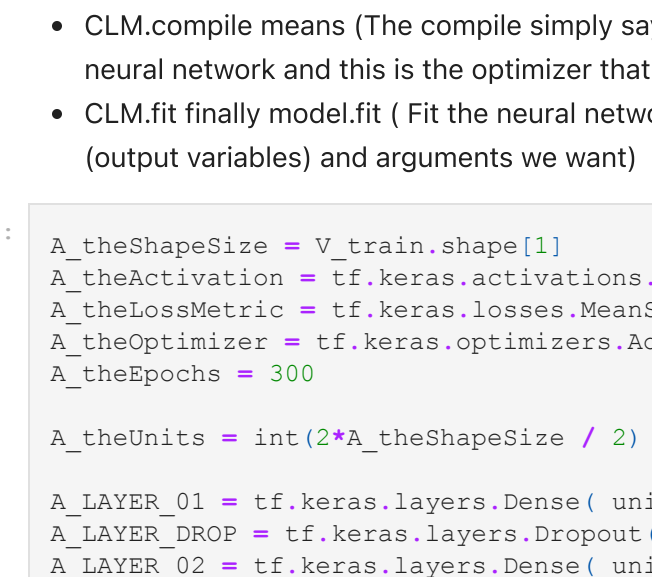
## Analysis for Tensor Flow ROC Curve (All 3 Activations Functions)

- I started off with 2 times the shape of the dataset for the nodes which I decided to stick with 2 times the shape of dataset and compared to when I started removing nodes. So, I better to stick with 2 times the shape of dataset for the nodes used in the neural network.
- The ROC Curve analysis for the Tensor Flow model showed that Area under the curve accuracy using relu activation was much better compare to other activation functions with 90% accuracy for and 89% AUC accuracy for predicting default loans. The ROC curve analysis using tanh (hyperbolic tangent) and elu (exponential linear unit) activation function also performed well but had a little bit lower accuracy at around 88% and 86% for both activation on this particular dataset. The Tensor flow models with relu and tanh activation function had much smoother curve the test and train dataset ROC curve is almost similar and accuracy for both train and test data is really close to one another. showing we are not overfitting the model.

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

```
In [13]: ALL_CLM = [ TREE_CLM, RF_CLM, GB_CLM, REG_ALL_CLM, REG_TREE_CLM, REG_RF_CLM, REG_GB_CLM ]
ALL_CLM = sorted( ALL_CLM, key = lambda x: x[4], reverse=True )
print_ROC_Curve( WHO, ALL_CLM )

ALL_CLM = sorted( ALL_CLM, key = lambda x: x[1], reverse=True )
print_Accuracy( "ALL CLASSIFICATION ACCURACY", ALL_CLM )
```



```
ALL CLASSIFICATION ACCURACY
=====
GB = 0.903254899328859
RF = 0.90184537832926
Tensor_Flow = 0.899261744966443
REG_ALL = 0.886744966442953
TREE = 0.886744966442953
REG_GB = 0.88874832214765
REG_RF = 0.8808724832214765
REG_TREE = 0.8674496644295302
=====
```

## Discuss which one is the most accurate. Which one would you recommend using?

The Tensorflow model using relu activation performed well based on analysis of ROC curve showing 90% accuracy after adding drop out and hidden layers to the neural network using relu activation function. Since we also selected variables that Random Forest tree model liked it had higher accuracy compare to other regression and tree based models in terms of variables selection techniques affected the results. I would recommend using random forest model since its covering 96% of area under the cover and has highest classification accuracy compared to other models.

## Tensorflow Model

### Develop a model using Tensor Flow that will predict loss amount given loans default

#### Data Transformation

- We are going to transform the W\_train and W\_test dataset so that its scales between 0 and 1 put it into V\_train and V\_test and convert them into dataframe. The data only has data if their loans are defaults and have loss amount.

#### Explore using a variable selection technique

- We are going select from the variables which Gradient Boosting tree based model liked since it has performed well in terms overall RMSE accuracy for this particular dataset compare to other models. We are going to use those variables in the V\_train and V\_test

```
In [14]: V_train = theScaler.transform( W_train )
V_test = theScaler.transform( W_test )

V_train = pd.DataFrame( V_train )
V_test = pd.DataFrame( V_test )

V_train.columns = list( W_train.columns.values )
V_test.columns = list( W_test.columns.values )

V_train = V_train[GB_amt]
V_test = V_test[GB_amt]
```

## Try at least three different Activation Functions

### First Activation Function used RELU

- 1. Try one and two hidden layers
- 1. Try using a Dropout Layer

- A\_ is for the flag loss amount given the person has default loans
- A\_ theShapeSize is we would like to know what's the size of the dataset.
- theActivation = tf.keras.activations.linear
- theLossMetric = tf.keras.losses.MeanSquaredError() the loss metric for linear regression purpose is that it minimizes the mean squared error.
- A\_ theOptimizer is Adam() function
- A\_ theEpochs go through the data 300 times
- A\_ theUnits is how many nodes should I have in my theUnits. I started off with 2 times the shape of the dataset which had the worse RMSE accuracy. To improve it we basically divided into half of the 2 times shape of dataset which is original number of variables that Gradient Boosting model like are used for the neural network.
- A\_ LAYER\_01 is a Dense layer
- A\_ LAYER\_DROP and A\_ LAYER\_02 without input his time. Adding a Dropout layer as well where we are simply saying everytime we run through the iteration throw away 20% of the nodes it will throw away 20% of the nodes from whatever nodes were in LAYER\_01 called before it this would prevent it from overfitting the model
- A\_ LAYER\_OUTPUT has 1 input for loss amount given loan defaults. The activation function used is linear

- Create a model called AMT (LOSS AMOUNT) thats a called Sequential()
- CLM.compile means (The compile simply say this is the loss function we are gonna use to train the neural network and this is the optimizer that is gonna adjust the weights)
- CLM.fit finally model.fit ( Fit the neural network with V\_train(input variables), Z\_train[TARGET\_A] (output variables) and arguments we want)

```
In [15]: A_theShapeSize = V_train.shape[1]
A_theActivation = tf.keras.activations.relu
A_theLossMetric = tf.keras.losses.MeanSquaredError()
A_theOptimizer = tf.keras.optimizers.Adam()
A_theEpochs = 300

A_theUnits = int(2*A_theShapeSize / 2)

A_LAYER_01 = tf.keras.layers.Dense(units=A_theUnits, activation=A_theActivation, input_shape=(A_LAYER_DROP,))
A_LAYER_DROP = tf.keras.layers.Dropout( 0.2 )
A_LAYER_02 = tf.keras.layers.Dense(units=A_theUnits, activation=A_theActivation)
A_LAYER_OUTPUT = tf.keras.layers.Dense(units=1, activation=tf.keras.activations.linear)
```

```
AMT = tf.keras.Sequential()
AMT.add( A_LAYER_01 )
AMT.add( A_LAYER_DROP )
AMT.add( A_LAYER_02 )
AMT.add( A_LAYER_OUTPUT )

AMT.compile( loss=A_theLossMetric, optimizer=A_theOptimizer)

AMT.fit( V_train, Z_train[TARGET_A], epochs=A_theEpochs, verbose=False )

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

TF_AMT = TEST_AMT.copy()
TF_AMT = TEST_AMT.copy()
```

```
Tensor_Flow RMSE ACCURACY
=====
Tensor_Flow_Train    = 8769.23894974666
Tensor_Flow         = 9019.736396725542
=====
```

## Try at least three different Activation Functions

### Second Activation Function used selu (Scaled Exponential Linear Unit)

- 1. Try one and two hidden layers
- 1. Try using a Dropout Layer

```
In [16]: A_theShapeSize = V_train.shape[1]
A_theActivation = tf.keras.activations.selu
A_theLossMetric = tf.keras.losses.MeanSquaredError()
A_theOptimizer = tf.keras.optimizers.Adam()
A_theEpochs = 300

A_theUnits = int(2*A_theShapeSize / 2)

A_LAYER_01 = tf.keras.layers.Dense(units=A_theUnits, activation=A_theActivation, input_shape=(A_LAYER_DROP,))
A_LAYER_DROP = tf.keras.layers.Dropout( 0.2 )
A_LAYER_02 = tf.keras.layers.Dense(units=A_theUnits, activation=A_theActivation)
A_LAYER_OUTPUT = tf.keras.layers.Dense(units=1, activation=tf.keras.activations.linear)
```

```
AMT = tf.keras.Sequential()
AMT.add( A_LAYER_01 )
AMT.add( A_LAYER_DROP )
AMT.add( A_LAYER_02 )
AMT.add( A_LAYER_OUTPUT )

AMT.compile( loss=A_theLossMetric, optimizer=A_theOptimizer)

AMT.fit( V_train, Z_train[TARGET_A], epochs=A_theEpochs, verbose=False )

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

Tensor_Flow RMSE ACCURACY
=====
Tensor_Flow_Train    = 8835.90533955421
Tensor_Flow         = 9076.76812743953
=====
```

## Try at least three different Activation Functions

### Third Activation Function used elu (Exponential Linear Unit)

- 1. Try one and two hidden layers
- 1. Try using a Dropout Layer

```
In [17]: A_theShapeSize = V_train.shape[1]
A_theActivation = tf.keras.activations.elu
A_theLossMetric = tf.keras.losses.MeanSquaredError()
A_theOptimizer = tf.keras.optimizers.Adam()
A_theEpochs = 300

A_theUnits = int(2*A_theShapeSize / 2)

A_LAYER_01 = tf.keras.layers.Dense(units=A_theUnits, activation=A_theActivation, input_shape=(A_LAYER_DROP,))
A_LAYER_DROP = tf.keras.layers.Dropout( 0.2 )
A_LAYER_02 = tf.keras.layers.Dense(units=A_theUnits, activation=A_theActivation)
A_LAYER_OUTPUT = tf.keras.layers.Dense(units=1, activation=tf.keras.activations.linear)
```

```
AMT = tf.keras.Sequential()
AMT.add( A_LAYER_01 )
AMT.add( A_LAYER_DROP )
AMT.add( A_LAYER_02 )
AMT.add( A_LAYER_OUTPUT )

AMT.compile( loss=A_theLossMetric, optimizer=A_theOptimizer)

AMT.fit( V_train, Z_train[TARGET_A], epochs=A_theEpochs, verbose=False )

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

Tensor_Flow RMSE ACCURACY
=====
Tensor_Flow_Train    = 8902.04803148962
Tensor_Flow         = 9050.460310570483
=====
```

## Analysis for Tensor Flow RMSE Accuracy (All 3 Activations Functions)

- I started off with 2 times the shape of the dataset which had the worse RMSE accuracy. To improve to we basically half of the 2 time shape of dataset which original variables used. Also, added dropout, hidden layer, and had to increase the number of iteration of data so made epochs to 300 iteration. It gave much better RMSE accuracy compared first testing phase of the model.
- The RMSE accuracy using Tensor flow for all 3 activation functions are really close to one another. Only the relu activation was little bit better the most used one which we used to compare with other models but we still worst than all the other regression and tree models.

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

```
In [18]: ALL_AMT = [ TREE_AMT, RF_AMT, GB_AMT, REG_ALL_AMT, REG_TREE_AMT, REG_RF_AMT, REG_GB_AMT ]
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 = 2201.8194437335387
RF = 2780.278461970596
REG_ALL = 3115.7180904557335
REG_GB = 4029.784013571014
REG_TREE = 4301.963358361031
REG_RF = 4360.610292020044
TREE = 5322.483680969033
Tensor_Flow = 9019.736396725542
=====
```

## Discuss which one is the most accurate. Which one would you recommend using?

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. Surprisingly, Tensorflow Model had worse RMSE accuracy for this dataset even after adding dropout, hidden layers, removing nodes, and increase iteration for running the data compared to all other regression and tree based models. Based on the All Loss Amount Model Accuracy we can recommend using Gradient Boosting and Random Forest model since they have the highest RMSE square error.

## 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 in predicting loans defaults.
- The Tensorflow model using relu activation performed well based on analysis of ROC curve showing 90% accuracy after adding drop out and hidden layers to the neural network. I would recommend using random forest model since its covering 96% of area under the cover and has highest classification accuracy compared to other models. But recommend trying out simple Decision Tree use and Regression ALL variables model to see how the accuracy and how model behaves when we use all variables for real world situation.

For any analysis of the coefficients, discuss whether or not they make sense. If any variable does not make sense, what would you recommend?

- The variable selection techniques for Tensor flow models based on the variables which Random Forest tree based model liked since it has performed well in terms overall ROC curve for this particular dataset for predicting default loans. And the variable which Gradient Boosting liked for predicting loss amount since it performed well from RMSE score perspective compare to other models. The variables made sense and to improve Tensor model accuracy for this dataset I would recommend doing variable selection technique based the variables that tree based model liked.

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 higher accuracy compare to other models. My second choice would have been Gradient Boosting because as it tends to build very shallow trees but we have to look at it is legal to use based business rules before deploying in production. I would recommend selecting Tensorflow model in production because it performed really bad for predicting loss amount not repaid given default loans in the RMSE analysis. Also its harder to deploy in production for Home Equity Loan area because of if something illegal happened such as discriminatory practices it would be harder to debug.

## Tensor Flow Model To Predict Loan Defaults:

### Discuss which one is the most accurate. Which one would you recommend using?

- The Tensorflow model using relu activation performed well based on analysis of ROC curve showing 90% accuracy after adding drop out and hidden layers to the neural network using relu activation function. Since we also selected variables that Random Forest tree model liked it had higher accuracy compare to other regression and tree based models in terms of variables selection techniques affected the results. I would recommend using random forest model since its covering 96% of area under the cover and has highest classification accuracy compared to other models.

## Tensor Flow Model to Predict Loss Given Defaults:

### Discuss which one is the most accurate. Which one would you recommend using?

- 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. Surprisingly, Tensorflow Model had worse RMSE accuracy for this dataset even after adding dropout, hidden layers, removing nodes, and increase iteration for running the data compared to all other regression and tree based models. Based on the All Loss Amount Model Accuracy we can recommend using Gradient Boosting and Random Forest model since they have the highest RMSE square error.

```
In [19]:
```

```
In [19]:
```