

In [7]:

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, accuracy_score
from sklearn.model_selection import GridSearchCV
import pydot
from io import StringIO
from sklearn.tree import export_graphviz
#from dm_tools import data_prep
from sklearn.preprocessing import StandardScaler
from sklearn.neural_network import MLPClassifier
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import classification_report, accuracy_score
from collections import Counter
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import classification_report, accuracy_score
from sklearn.feature_extraction.text import CountVectorizer
from imblearn.over_sampling import SMOTE
from imblearn.under_sampling import ClusterCentroids
from sklearn.neural_network import MLPClassifier
```

In [16]:

```
# Preprocessing
#listing the nominal and numerical values
nominal_cols = ['Auction', 'Make', 'TopThreeAmericanName', 'Color', 'Transmission', 'Nationality', 'Size', 'VNST']

num_cols = ['WheelTypeID', 'VehYear', 'VehBCost', 'VehOdo', 'IsOnlineSale', 'WarrantyCost'] + ['MMRAcquisitionAuctionAveragePrice', 'MMRAcquisitionAuctionCleanPrice',
                                                'MMRAcquisitionRetailAveragePrice', 'MMRAcquisitionRetailCleanPrice',
                                                'MMRCurrentAuctionAveragePrice', 'MMRCurrentAuctionCleanPrice',
                                                'MMRCurrentRetailAveragePrice', 'MMRCurrentRetailCleanPrice']

## Defining missing values
def fill_missing_values(df):

    for col in num_cols:
        df[col] = df[col].fillna(df[col].median())

    for col in nominal_cols:
        mode = df[col].mode()[0]
        df[col] = df[col].fillna(mode)

    return df

#deleting the unwanted features
def feature_engineering(df):

    del df['WheelType'] #Wheeltype ID is used. Wheeltype and WheelTypeID is derived
    del df['PurchaseID']
    del df['ForSale']
    del df['PurchaseDate'] #We are using Vehicle year to measure the time series and hence,
    deleting it.
    del df['MMRCurrentRetailRatio'] #Derived type
    del df['PRIMEUNIT'] # more than 80% of the values are empty
    del df['AUCGUART'] # more than 80% of the values are empty
    del df['PurchaseTimestamp'] # We are using year as a way of time measure. Hence, Time stamp is
    not needed.

    return df

def data_type_change(df):
    # change Transmission into binary 0/1 variable
    Transmission_map = {"AUTO":0, "MANUAL": 1, 'Manual':1}
    df['Transmission'] = df['Transmission'].map(Transmission_map)
    df['Transmission'].fillna(df['Transmission'].mode(), inplace=True)

    # WheelType_map = {"Alloy":1, "Covers": 2, "Special": 3}
    # df['WheelType'] = df['WheelType'].map(WheelType_map)
```

```

Auction_map={'ADESA':0, 'MANHEIM':1, 'OTHER':2}
df['Auction']=df['Auction'].map(Auction_map)

Make_map={'ACURA':0, 'BUICK':1, 'CADILLAC':3, 'CHEVROLET':4, 'CHRYSLER':5, 'DODGE':6, 'FORD':7, 'GMC':
8, 'HONDA':9, 'HYUNDAI':10, 'INFINITI':11, 'ISUZU':12, 'JEEP':13, 'KIA':14, 'LEXUS':15, 'LINCOLN':16, 'MAZDA':
17, 'MERCURY':18, 'MINI':19, 'MITSUBISHI':20, 'NISSAN':21, 'OLDSMOBILE':22, 'PONTIAC':23, 'SATURN':24, 'S
CION':25, 'SUBARU':26, 'SUZUKI':27, 'TOYOTA':2, 'VOLKSWAGEN':28, 'VOLVO':29}
df['Make']=df['Make'].map(Make_map)
#df['Make'].fillna(df['Make'].mode(), inplace=True)

american_name_map={'CHRYSLER':0, 'FORD':1, 'GM':2, 'OTHER':3}
df['TopThreeAmericanName']=df['TopThreeAmericanName'].map(american_name_map)

Color_map={'BEIGE':0, 'BLACK':1, 'BLUE':2, 'BROWN':3, 'GOLD':4, 'GREEN':5, 'GREY':6, 'MAROON':7, 'NOT A
VAIL':8, 'ORANGE':9, 'OTHER':10, 'PURPLE':11, 'RED':12, 'SILVER':13, 'WHITE':14, 'YELLOW':15}
df['Color']=df['Color'].map(Color_map)

Nationality_map={'AMERICAN':0, 'OTHER':1, 'OTHER ASIAN':2, 'TOP LINE ASIAN':3, 'USA':4}
df['Nationality']=df['Nationality'].map(Nationality_map)

Size_map={'COMPACT':0, 'CROSSOVER':1, 'LARGE':2, 'LARGE SUV':3, 'LARGE TRUCK':4, 'MEDIUM':5, 'MEDIUM
SUV':6, 'SMALL SUV':7, 'SMALL TRUCK':8, 'SPECIALTY':9, 'SPORTS':10, 'VAN':11}
df['Size']=df['Size'].map(Size_map)

vnst_map = {'TX':0, 'FL':1, 'CO':2, 'NC':3, 'AZ':4, 'CA':5, 'OK':6, 'SC':7, 'TN':8, 'GA':9, 'VA':10, 'MO':
11, 'PA':12, 'NV':13, 'IN':14, 'MS':15, 'LA':16, 'NJ':17, 'NM':18, 'KY':19, 'AL':20, 'IL':21, 'UT':22, 'WV':2
3, 'WA':24, 'OR':25, 'NH':26, 'NE':27, 'OH':28, 'ID':29, 'NY':30}
df['VNST'] = df['VNST'].map(vnst_map)

df['Transmission'] = df['Transmission'].astype(float)
df['Auction'] = df['Auction'].astype(float)
df['Make'] = df['Make'].astype(float)
df['TopThreeAmericanName'] = df['TopThreeAmericanName'].astype(float)
df['Nationality'] = df['Nationality'].astype(float)
df['Size'] = df['Size'].astype(float)
df['VNST'] = df['VNST'].astype(float)
df['WheelTypeID'] = df['WheelTypeID'].astype(float)
df['VehBCost'] = df['VehBCost'].astype(float)
df['IsOnlineSale'] = df['IsOnlineSale'].astype(float)
df['MMRAcquisitionAuctionAveragePrice'] = df['MMRAcquisitionAuctionAveragePrice'].astype(float)
df['MMRAcquisitionAuctionCleanPrice'] = df['MMRAcquisitionAuctionCleanPrice'].astype(float)
df['MMRAcquisitionRetailAveragePrice'] = df['MMRAcquisitionRetailAveragePrice'].astype(float)
df['MMRAcquisitionRetailCleanPrice'] = df['MMRAcquisitionRetailCleanPrice'].astype(float)
df['MMRCurrentAuctionAveragePrice'] = df['MMRCurrentAuctionAveragePrice'].astype(float)
df['MMRCurrentAuctionCleanPrice'] = df['MMRCurrentAuctionCleanPrice'].astype(float)
df['MMRCurrentRetailAveragePrice'] = df['MMRCurrentRetailAveragePrice'].astype(float)
df['MMRCurrentRetailCleanPrice'] = df['MMRCurrentRetailCleanPrice'].astype(float)

return df

def error_replacing_For_IsOnlineSale(df):
    mask = df['IsOnlineSale'] == -1
    df.loc[mask, 'IsOnlineSale'] = np.nan
    mask = df['IsOnlineSale'] == 2
    df.loc[mask, 'IsOnlineSale'] = np.nan
    mask = df['IsOnlineSale'] == 4
    df.loc[mask, 'IsOnlineSale'] = np.nan

    return df

# Converting Nominal into Numerica type
def convert_nominal_cols(df):
    '''
    This function converts nominal cols to one-hot vectors
    '''
    global nominal_cols
    df_with_dummies = pd.get_dummies(df, columns = nominal_cols)

    return df_with_dummies

def analyse_feature_importance(dm_model, feature_names, n_to_display=20):
    # grab feature importances from the model

```

```

importances = dm_model.feature_importances_

# sort them out in descending order
indices = np.argsort(importances)
indices = np.flip(indices, axis=0)

# limit to 20 features, you can leave this out to print out everything
indices = indices[:n_to_display]

for i in indices:
    print(feature_names[i], ':', importances[i])

def visualize_decision_tree(dm_model, feature_names, save_name):
    import pydot
    from io import StringIO
    from sklearn.tree import export_graphviz

    dotfile = StringIO()
    export_graphviz(dm_model, out_file=dotfile, feature_names=feature_names)
    graph = pydot.graph_from_dot_data(dotfile.getvalue())
    graph[0].write_png(save_name) # saved in the following file

def error_replacing_For_IsOnlineSale(df):
    mask = df['IsOnlineSale'] == -1
    df.loc[mask, 'IsOnlineSale'] = np.nan
    mask = df['IsOnlineSale'] == 2
    df.loc[mask, 'IsOnlineSale'] = np.nan
    mask = df['IsOnlineSale'] == 4
    df.loc[mask, 'IsOnlineSale'] = np.nan

    return df

```

In [20]:

```

# Preprocssing

def preprocessing():

    df = pd.read_csv('Home.csv')

    #dropping the columns with continuously 10 null values
    new1_df=df.dropna(axis=0,thresh=10)

    #Rerplacing ? with null values
    new_df = new1_df.replace(['?'], np.nan, inplace=False)

    fill_missing_values(new_df)

    error_replacing_For_IsOnlineSale(df)

    feature_engineering(new_df)

    convert_nominal_cols(new_df)

    return new_df

```

In [21]:

```

df = pd.read_csv('Home.csv')
df = preprocessing()
df2 = data_type_change(df)
df2.to_csv("CaseStudyData1.csv")
print(df2.info())
y = df2['IsBadBuy']
X = df2.drop(['IsBadBuy'], axis=1)
rs=10
#Split the data based on training and testing with 70 and 30%
X_res, X_test, y_res, y_test = train_test_split(X, y, test_size=0.3, stratify=y, random_state=rs)

sm = SMOTE(random_state=42)
X_train, y_train = sm.fit_sample(X_res, y_res)

```

```
C:\ProgramData\Anaconda3\lib\site-packages\IPython\core\interactiveshell.py:3185: DtypeWarning:
Columns (27) have mixed types. Specify dtype option on import or set low_memory=False.
if (yield from self.run_code(code, result)):
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 41432 entries, 0 to 41475
Data columns (total 23 columns):
Auction                41432 non-null float64
VehYear                41432 non-null float64
Make                   41432 non-null float64
Color                  41432 non-null int64
Transmission           41432 non-null float64
WheelTypeID            41432 non-null float64
VehOdo                 41432 non-null float64
Nationality            41432 non-null float64
Size                   41432 non-null float64
TopThreeAmericanName   41432 non-null float64
MMRAcquisitionAuctionAveragePrice 41432 non-null float64
MMRAcquisitionAuctionCleanPrice    41432 non-null float64
MMRAcquisitionRetailAveragePrice    41432 non-null float64
MMRAcquisitionRetailCleanPrice      41432 non-null float64
MMRCurrentAuctionAveragePrice        41432 non-null float64
MMRCurrentAuctionCleanPrice          41432 non-null float64
MMRCurrentRetailAveragePrice         41432 non-null float64
MMRCurrentRetailCleanPrice           41432 non-null float64
VNST                                 41432 non-null float64
VehBCost                            41432 non-null float64
IsOnlineSale                        41432 non-null float64
WarrantyCost                        41432 non-null float64
IsBadBuy                            41432 non-null int64
dtypes: float64(21), int64(2)
memory usage: 7.6 MB
None
```

Task 4: Predictive Modeling Using Neural Networks

1. Build a Neural Network model using the default setting. Answer the following: a. What is the network architecture? We are just training our MPL Classifier with the random state as parameter. It will retain the data set in same state. We are using the number of iterations in order to reach the convergence. Convergence will help to provide the optimal solution for the loss function. It will help to reach the Global minima or maxima. Hence, convergence will be reach. Convergence will give the best or optimal solution.

b. How many iterations are needed to train this network? 360 iterations are needed to converge the network. there are many local minima and maxima present in the loss function. iterations helps to find the Global Minima or Global Maxima in the Loss function.

c. Do you see any sign of over-fitting? As the data is imbalanced, We balanced the dataset and we are checking the precision and f1score along with the train and testing accuracy of oversampled data. Here we got good accuracy with good precision and f1score. But difference between training and testing accuracies is a little bit more. however, both training and testing accuracy is performing well.

d. Did the training process converge and resulted in the best model? Training process get converged @iteration 360. Accuracy increased . But this is not the best model. beacuse we didn't mentioned the number of Hidden layers and the we didn't applied the reguralisation concept. There is chance that the model will get improve if we add those parameters.

Default model with no iterations(not converged): Train accuracy: 0.8685982492969462 Test accuracy: 0.7944489139179405

Default model with 360 iterations(converged): Train accuracy: 0.8787578722224423 Test accuracy: 0.7936444086886565

After 360 iterations, we didn't get the error message of Training accuracy is increased little bit and resulted the better model.

e. What is classification accuracy on training and test datasets? Default model with no iterations(not converged): Train accuracy: 0.8685982492969462 Test accuracy: 0.7944489139179405

Default model with 360 iterations(converged): Train accuracy: 0.8787578722224423 Test accuracy: 0.7936444086886565

In [22]:

```
#Scaling the entire Data set, which will make all the points in single scale.
```

```

scaler = StandardScaler()
X_train = scaler.fit_transform(X_train, y_train)
X_test = scaler.transform(X_test)

```

C:\ProgramData\Anaconda3\lib\site-packages\ipykernel_launcher.py:5: DataConversionWarning: Data with input dtype int64, float64 were all converted to float64 by StandardScaler.

a. What is the network architecture? Ans: Multi Layer Perceptron(MLP) with default setting is shown below.

In [23]:

```

#Default MLP without any iterations

model = MLPClassifier(random_state=rs)
model.fit(X_train, y_train)

print("Train accuracy:", model.score(X_train, y_train))
print("Test accuracy:", model.score(X_test, y_test))

y_pred = model.predict(X_test)
print(classification_report(y_test, y_pred))

print(model)

```

```

Train accuracy: 0.8685982492969462
Test accuracy: 0.7944489139179405

```

	precision	recall	f1-score	support
0	0.89	0.87	0.88	10820
1	0.24	0.27	0.25	1610
micro avg	0.79	0.79	0.79	12430
macro avg	0.56	0.57	0.57	12430
weighted avg	0.80	0.79	0.80	12430

```

MLPClassifier(activation='relu', alpha=0.0001, batch_size='auto', beta_1=0.9,
              beta_2=0.999, early_stopping=False, epsilon=1e-08,
              hidden_layer_sizes=(100,), learning_rate='constant',
              learning_rate_init=0.001, max_iter=200, momentum=0.9,
              n_iter_no_change=10, nesterovs_momentum=True, power_t=0.5,
              random_state=10, shuffle=True, solver='adam', tol=0.0001,
              validation_fraction=0.1, verbose=False, warm_start=False)

```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\neural_network\multilayer_perceptron.py:562: ConvergenceWarning: Stochastic Optimizer: Maximum iterations (200) reached and the optimization hasn't converged yet.

```

% self.max_iter, ConvergenceWarning)

```

In [8]:

```

model = MLPClassifier(max_iter=360, random_state=rs)
model.fit(X_train, y_train)

print("Train accuracy:", model.score(X_train, y_train))
print("Test accuracy:", model.score(X_test, y_test))

y_pred = model.predict(X_test)
print(classification_report(y_test, y_pred))

print(model)

```

```

Train accuracy: 0.878757872224423
Test accuracy: 0.7936444086886565

```

	precision	recall	f1-score	support
0	0.89	0.87	0.88	10820
1	0.25	0.30	0.28	1610
micro avg	0.79	0.79	0.79	12430
macro avg	0.57	0.59	0.58	12430
weighted avg	0.81	0.79	0.80	12430

```
MLPClassifier(activation='relu', alpha=0.0001, batch_size='auto', beta_1=0.9,
              beta_2=0.999, early_stopping=False, epsilon=1e-08,
              hidden_layer_sizes=(100,), learning_rate='constant',
              learning_rate_init=0.001, max_iter=360, momentum=0.9,
              n_iter_no_change=10, nesterovs_momentum=True, power_t=0.5,
              random_state=10, shuffle=True, solver='adam', tol=0.0001,
              validation_fraction=0.1, verbose=False, warm_start=False)
```

In [9]:

```
print("X train shape :", X_train.shape)
```

X train shape : (50494, 22)

2. Refine this network by tuning it with GridSearchCV.

a. What is the network architecture?

We have 90 hidden layers for this network when we have just done one iteration. It became 93 when we run it for 360 iterations. We are also using numberOfJobs(n_jobs) as -1. As it is doing Gridsearch search with 360 iterations it will take more time to give the good model.

b. How many iterations are needed to train this network? Here, 500 iterations are needed to iterate the model. Here, we are using the number of hidden layers and alpha value. If we give iterations to MLP classifier then it will take more time to converge. Due to computational complexity. we didn't added the number of iterations.

c. Sign of overfitting? Precision and f1score are increased with slight difference between training and testing accuracy. Overfitting has been improved when compare to the previous case. testing accuracy is more and hence, we are not seeing an sign of overfitting.

d. Did the training process converge and resulted in the best model?

We are tuning the second hyperparameter, Alpha as L2 regularization parameter used in each neuron's activation function. Training process is not converged for the default 1 iterations. We got warning message for the default MLP classifier. We need to introduce the number of iterations in order to attain the convergence. But, due to lack of computational power, we didn't added the iterations in MLP classifier.

below MLP classifier will attain the convergence params = {'hidden_layer_sizes': [(x,) for x in range(70, 95, 5)]}

```
cv = GridSearchCV(param_grid=params, estimator=MLPClassifier(max_iter=500, random_state=rs), cv=10, n_jobs=-1) cv.fit(X_train,
y_train)
```

```
print("Train accuracy:", cv.score(X_train, y_train)) print("Test accuracy:", cv.score(X_test, y_test))
```

```
y_pred = cv.predict(X_test) print(classification_report(y_test, y_pred))
```

```
print(cv.bestparams)
```

default MLP classifier with no iterations will not converge. If we use iterations in the MLP classifier then it will converge. It will attain the Global Maxima or Minima.

e. What is classification accuracy on training and test datasets? Is there any improvement in the outcome? Classification Accuracy:

Optimal Hidden layer size: 93 optimal Hyper-parameter: 1e-05 Train accuracy: 0.8852933021745157 Test accuracy: 0.8083668543845535

Accuray is improved a little bit when compare to the default model and precision & f1score also increased. Hence, GridSearch CV method helps to improve the accuracy of the model.

In [37]:

```
params = {'hidden_layer_sizes': [(x,) for x in range(70, 95, 5)]}
```

```
cv = GridSearchCV(param_grid=params, estimator=MLPClassifier(random_state=rs), cv=10, n_jobs=-1)
cv.fit(X_train, y_train)
```

```
print("Train accuracy:", cv.score(X_train, y_train))
print("Test accuracy:", cv.score(X_test, y_test))
```

```
y_pred = cv.predict(X_test)
```

```
print(classification_report(y_test, y_pred))

print(cv.best_params_)
```

```
Train accuracy: 0.8623994930090704
Test accuracy: 0.8120675784392598
      precision    recall  f1-score   support

      0         0.89      0.90      0.89      10820
      1         0.26      0.25      0.26       1610

   micro avg       0.81      0.81      0.81      12430
   macro avg       0.58      0.57      0.57      12430
weighted avg       0.81      0.81      0.81      12430

{'hidden_layer_sizes': (90,)}
```

```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\neural_network\multilayer_perceptron.py:562: Co
nvergenceWarning: Stochastic Optimizer: Maximum iterations (200) reached and the optimization
hasn't converged yet.
  % self.max_iter, ConvergenceWarning)
```

In [38]:

```
params = {'hidden_layer_sizes': [(90), (91,), (92,), (93,)], 'alpha': [0.01, 0.001, 0.0001, 0.00001]}

cv = GridSearchCV(param_grid=params, estimator=MLPClassifier(max_iter=360, random_state=rs), cv=10,
n_jobs=-1)
cv.fit(X_train, y_train)

print("Train accuracy:", cv.score(X_train, y_train))
print("Test accuracy:", cv.score(X_test, y_test))

y_pred = cv.predict(X_test)
print(classification_report(y_test, y_pred))

print(cv.best_params_)
```

```
Train accuracy: 0.8852933021745157
Test accuracy: 0.8083668543845535
      precision    recall  f1-score   support

      0         0.89      0.88      0.89      10820
      1         0.28      0.29      0.28       1610

   micro avg       0.81      0.81      0.81      12430
   macro avg       0.58      0.59      0.59      12430
weighted avg       0.81      0.81      0.81      12430

{'alpha': 1e-05, 'hidden_layer_sizes': (93,)}
```

```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\neural_network\multilayer_perceptron.py:562: Co
nvergenceWarning: Stochastic Optimizer: Maximum iterations (360) reached and the optimization
hasn't converged yet.
  % self.max_iter, ConvergenceWarning)
```

3. Build another Neural Network model with inputs selected from RFE with regression (use the best model generated in Task 3) and selection with decision tree (use the best model from Task 2).

a. Did feature selection help here? Any change in the network architecture? What inputs are being used as the network input?

We found many features are skewed. Hence, we are using the log transformation to make them into normal/gaussian distribution.

Yes, Feature selection helped to increase the accuracy of neural network.

We are using Recursive feature elimination to reduce the dimensional size. So unimportant features are removed in this model. Neural Network is not having this option. Hence we are using Decision Tree and Logistic regression model as the base elimination model. Log transformation improved the model performance. Hence we will transform the original dataset.

RFECV is used to find the best number of features for the model. It is used to find the best number of features for the model.

REF Logistic regression has selected 15 features. REF Decision Tree has reduced the dimension to 5.

b. What is classification accuracy on training and test datasets? Is there any improvement in the outcome?

Feature selection helps to reduce the number of dimensions. REF will help to give the weights to the features and it will remove the unimportant features. Neural network does not assign weight/feature importance to each feature. Hence we are using the Logistic Regression and Decision Tree.

Log transformation has proven to improve model performance. That we already performed this in Logistic Regression. We will use the Log transformed dataset to perform the RFE (Recursive Feature Elimination)

The RFE with logistic regression has selected 15 features as the best set of features. With these selected features, we are tuning an MLPClassifier with the transformed data set as training data.

a. RFE using Logistic Regression.

Classification accuracy for Logistic regression: Train accuracy: 0.8645070024048663 Test accuracy: 0.8495610271305036

SelectFromModel reduces the log transformed dataset into only 5 variables. With these selected features, we are tuning an MLPClassifier with the transformed data set as training data.

b. Feature selection using Decision Tree Classification accuracy for Decision Tree: Train accuracy: 0.8543589337347012 Test accuracy: 0.8409493161705551

Transforming features into log-transformation, to remove the skewness.

c. How many iterations are now needed to train this network?

360 iterations are needed to achieve the converge. below MLP classifier will attain the convergence params = {'hidden_layer_sizes': [(x,) for x in range(70, 95, 5)]}

```
cv = GridSearchCV(param_grid=params, estimator=MLPClassifier(max_iter=500, random_state=rs), cv=10, n_jobs=-1) cv.fit(X_train, y_train)
```

```
print("Train accuracy:", cv.score(X_train, y_train)) print("Test accuracy:", cv.score(X_test, y_test))
```

```
y_pred = cv.predict(X_test) print(classification_report(y_test, y_pred))
```

```
print(cv.best_params)
```

default MLP classifier with no iterations will not converge. If we use iterations in the MLP classifier then it will converge. It will attain the Global Maxima or Minima.

d. Do you see any sign of over-fitting? There is no sign of overfitting. Bcz, There is no huge difference between the train and Test accuracy. Testing accuracy is good. There is no sign of bias as well.

c. Did the training process converge and resulted in the best model? Nope, The result is not converged. Because, we didn't use the iterations in the MLP classifier. If we use the iterations in the MLP classifier then we can achieve the Convergence.

Default Accuracy: Train accuracy: 0.8685982492969462 Test accuracy: 0.7944489139179405

Accuracy using RFE Train accuracy: 0.8645070024048663 Test accuracy: 0.8495610271305036

Accuracy improved a little bit. when compare to the default MLP classifier. Hence, it can be the best or optimal model.

In [24]:

```
import numpy as np

# list columns to be transformed
columns_to_transform = ['Make', 'TopThreeAmericanName', 'Color', 'Transmission', 'Nationality', 'Size',
                        'VNST', 'MMRAcquisitionRetailAveragePrice', 'MMRCurrentAuctionAveragePrice', 'MMRAcquisitionAuctionAveragePrice']

df12 = pd.DataFrame(X_train, columns=['Auction', 'VehYear', 'Make', 'Color', 'Transmission', 'WheelType', 'VehOdo', 'Nationality', 'Size', 'TopThreeAmericanName', 'MMRAcquisitionAuctionAveragePrice', 'MMRAcquisitionAuctionCleanPrice', 'MMRAcquisitionRetailAveragePrice', 'MMRAcquisitionRetailCleanPrice', 'MMRCurrentAuctionAveragePrice', 'MMRCurrentAuctionCleanPrice', 'MMRCurrentRetailAveragePrice', 'MMRCurrentRetailCleanPrice', 'VNST', 'VehBCost', 'IsOnlineSale', 'WarrantyCost'])
df13 = pd.DataFrame(y_train, columns=['IsBadBuy'])
df14 = pd.concat([df12, df13], axis=1, sort=False)

df_log = df14.copy()
```



```

X_log = df_log.copy()
# transform the columns with np.log
for col in columns_to_transform:
    df_log[col] = df_log[col].apply(lambda x: x+1)
    df_log[col] = df_log[col].apply(np.log)

df_log = df_log.apply(lambda x:x.fillna(x.value_counts().index[0]))

# create X, y and train test data partitions
y_log = df_log['IsBadBuy']
X_log = df_log.drop(['IsBadBuy'], axis=1)
#X_mat_log = X_log.as_matrix()
X_train_log, X_test_log, y_train_log, y_test_log = train_test_split(X_log, y_log, test_size=0.3, stratify=y_log,
                                                                    random_state=rs)

# standardise them again
scaler_log = StandardScaler()
X_train_log = scaler_log.fit_transform(X_train_log, y_train_log)
X_test_log = scaler_log.transform(X_test_log)

```

In [10]:

```

params = {'hidden_layer_sizes': [(91,), (92,), (93,)], 'alpha': [0.01, 0.001, 0.0001, 0.00001]}

cv = GridSearchCV(param_grid=params, estimator=MLPClassifier(max_iter=360, random_state=rs), cv=10,
n_jobs=-1)
cv.fit(X_train_log, y_train_log)

print("Train accuracy:", cv.score(X_train_log, y_train_log))
print("Test accuracy:", cv.score(X_test_log, y_test_log))

y_pred = cv.predict(X_test_log)
print(classification_report(y_test_log, y_pred))

print(cv.best_params_)

```

Train accuracy: 0.8874522563304569

Test accuracy: 0.8571522872796884

	precision	recall	f1-score	support
0	0.82	0.91	0.86	7574
1	0.90	0.80	0.85	7575
micro avg	0.86	0.86	0.86	15149
macro avg	0.86	0.86	0.86	15149
weighted avg	0.86	0.86	0.86	15149

```
{'alpha': 0.0001, 'hidden_layer_sizes': (93,)}
```

In [25]:

```

#RFE with Logistic Regression

from sklearn.feature_selection import RFECV
from sklearn.linear_model import LogisticRegression

rfe = RFECV(estimator = LogisticRegression(random_state=rs), cv=10)
rfe.fit(X_train_log, y_train_log)

print(rfe.n_features_)

```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:433: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.

FutureWarning)

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:433: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.

FutureWarning)

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:433: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.

FutureWarning)

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:433: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.

FutureWarning)

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]


```
print("Test accuracy:", cv.score(X_test_rfe, y_test_log))

y_pred = cv.predict(X_test_rfe)
print(classification_report(y_test_log, y_pred))

print(cv.best_params_)
```

```
Train accuracy: 0.8645070024048663
Test accuracy: 0.8495610271305036
```

	precision	recall	f1-score	support
0	0.82	0.90	0.86	7574
1	0.89	0.80	0.84	7575
micro avg	0.85	0.85	0.85	15149
macro avg	0.85	0.85	0.85	15149
weighted avg	0.85	0.85	0.85	15149

```
{'alpha': 0.001, 'hidden_layer_sizes': (91,)}
```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\neural_network\multilayer_perceptron.py:562: ConvergenceWarning: Stochastic Optimizer: Maximum iterations (200) reached and the optimization hasn't converged yet.
% self.max_iter, ConvergenceWarning)

In [27]:

```
from sklearn.tree import DecisionTreeClassifier

params = {'criterion': ['gini', 'entropy'],
          'max_depth': range(3, 8),
          'min_samples_leaf': range(20, 61, 10)}

cv = GridSearchCV(param_grid=params, estimator=DecisionTreeClassifier(random_state=rs), cv=10)
cv.fit(X_train_log, y_train_log)
```

Out[27]:

```
GridSearchCV(cv=10, error_score='raise-deprecating',
             estimator=DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
                                              max_features=None, max_leaf_nodes=None,
                                              min_impurity_decrease=0.0, min_impurity_split=None,
                                              min_samples_leaf=1, min_samples_split=2,
                                              min_weight_fraction_leaf=0.0, presort=False, random_state=10,
                                              splitter='best'),
             fit_params=None, iid='warn', n_jobs=None,
             param_grid={'criterion': ['gini', 'entropy'], 'max_depth': range(3, 8), 'min_samples_leaf':
range(20, 61, 10)},
             pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
             scoring=None, verbose=0)
```

In [28]:

```
analyse_feature_importance(cv.best_estimator_, X_log.columns)
```

```
WheelType : 0.44757157420636856
Color : 0.1974940571093593
Auction : 0.11880470950360118
TopThreeAmericanName : 0.10297043468792931
Make : 0.05501741952038854
VehYear : 0.04141094755975841
Size : 0.020177359341664633
VehBCost : 0.007047476675064481
MMRCurrentRetailAveragePrice : 0.003085605662151227
VehOdo : 0.002922476001337584
VNST : 0.0016630976202473575
Nationality : 0.000880365785788005
MMRAcquisitionAuctionCleanPrice : 0.0008594987787911562
WarrantyCost : 9.497754755005371e-05
IsOnlineSale : 0.0
MMRAcquisitionRetailAveragePrice : 0.0
MMRAcquisitonRetailCleanPrice : 0.0
Transmission : 0.0
```

MMRCurrentAuctionAveragePrice : 0.0
MMRCurrentAuctionCleanPrice : 0.0

In [24]:

```
from sklearn.feature_selection import SelectFromModel

selectmodel = SelectFromModel(cv.best_estimator_, prefit=True)
X_train_sel_model = selectmodel.transform(X_train)
X_test_sel_model = selectmodel.transform(X_test)

print(X_train_sel_model.shape)

params = {'hidden_layer_sizes': [(91,), (92,), (93,)], 'alpha': [0.01, 0.001, 0.0001, 0.00001]}

cv = GridSearchCV(param_grid=params, estimator=MLPClassifier(random_state=rs), cv=10, n_jobs=-1)
cv.fit(X_train_sel_model, y_train)

print("Train accuracy:", cv.score(X_train_sel_model, y_train))
print("Test accuracy:", cv.score(X_test_sel_model, y_test))

y_pred = cv.predict(X_test_sel_model)
print(classification_report(y_test, y_pred))

print(cv.best_params_)
```

```
(50494, 5)
Train accuracy: 0.8543589337347012
Test accuracy: 0.8409493161705551
```

	precision	recall	f1-score	support
0	0.88	0.95	0.91	10820
1	0.21	0.09	0.12	1610
micro avg	0.84	0.84	0.84	12430
macro avg	0.54	0.52	0.52	12430
weighted avg	0.79	0.84	0.81	12430

```
['alpha': 0.0001, 'hidden_layer_sizes': (93,)]
```

```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\neural_network\multilayer_perceptron.py:562: ConvergenceWarning: Stochastic Optimizer: Maximum iterations (200) reached and the optimization hasn't converged yet.
  % self.max_iter, ConvergenceWarning)
```

4. Using the comparison methods, which of the models (i.e one with selected

Accuracy os all models Accuracy score on test for DT: 0.8194690265486726 Accuracy score on test for logistic regression: 0.6262268704746581 Accuracy score on test for NN: 0.8133547868061143

Among the above models, Decision Tree and Neural Network is roughly gave the same accuracy. Both can be considered as the best model.

1. Using the comparison methods, which of the models (i.e one with selected variables and another with all variables) appears to be better? From the better model, can you identify cars those could potential be "kicks"? Can you provide some descriptive summary of those cars? Is it easy to comprehend the performance of the best neural network model for decision making? Ans: We found Decision Tree is performing better.

Accuracy score on test for DT: 0.8194690265486726 Accuracy score on test for logistic regression: 0.6262268704746581 Accuracy score on test for NN: 0.8133547868061143 From the better model, can you identify cars those could potential be "kicks"? Can you provide some descriptive summary of those cars? Is it easy to comprehend the performance of the best neural network model for

In [29]:

```
# grid search CV for decision tree
params_dt = {'criterion': ['gini'],
             'max_depth': range(2, 10),
             'min_samples_leaf': range(15, 25, 5)}
```

```

cv = GridSearchCV(param_grid=params_dt, estimator=DecisionTreeClassifier(random_state=rs), cv=10)
cv.fit(X_train, y_train)

dt_model = cv.best_estimator_
print(dt_model)

# grid search CV for logistic regression
params_log_reg = {'C': [pow(10, x) for x in range(-6, 4)]}

cv = GridSearchCV(param_grid=params_log_reg, estimator=LogisticRegression(random_state=rs), cv=10,
n_jobs=-1)
cv.fit(X_train, y_train)

log_reg_model = cv.best_estimator_
print(log_reg_model)

# grid search CV for NN
params_nn = {'hidden_layer_sizes': [(91,), (92,), (93,)], 'alpha': [0.01, 0.001, 0.0001, 0.00001]}

cv = GridSearchCV(param_grid=params_nn, estimator=MLPClassifier(random_state=rs), cv=10, n_jobs=-1)
cv.fit(X_train, y_train)

nn_model = cv.best_estimator_
print(nn_model)

```

```

DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=9,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=15, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False, random_state=10,
splitter='best')

```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:433: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
FutureWarning)

```

LogisticRegression(C=1000, class_weight=None, dual=False, fit_intercept=True,
intercept_scaling=1, max_iter=100, multi_class='warn',
n_jobs=None, penalty='l2', random_state=10, solver='warn',
tol=0.0001, verbose=0, warm_start=False)
MLPClassifier(activation='relu', alpha=0.0001, batch_size='auto', beta_1=0.9,
beta_2=0.999, early_stopping=False, epsilon=1e-08,
hidden_layer_sizes=(92,), learning_rate='constant',
learning_rate_init=0.001, max_iter=200, momentum=0.9,
n_iter_no_change=10, nesterovs_momentum=True, power_t=0.5,
random_state=10, shuffle=True, solver='adam', tol=0.0001,
validation_fraction=0.1, verbose=False, warm_start=False)

```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\neural_network\multilayer_perceptron.py:562: ConvergenceWarning: Stochastic Optimizer: Maximum iterations (200) reached and the optimization hasn't converged yet.
% self.max_iter, ConvergenceWarning)

5. Generating an Ensemble Model and Comparing Models

Task5:

1. Generate an ensemble model to include the best regression model, best decision tree model, and best neural network model.

a. Does the Ensemble model outperform the underlying models? Resonate your answer.

Yes, Ensemble model performs well when compare to all the other model. Accuracuy of all the base models such as Decision Tree, Logistic Regression, and Neural Network is less than the Ensemble model. Which is mentioned below.

1. Use the comparison methods (or the comparison node) to compare the best decision tree model, the best regression model, the best neural network model and the ensemble model. a. Discuss the findings led by (a) ROC Chart (and Index); (b) Score Ranking (or Accuracy Score); (c) Fit Statistics; (or Classification report) and (4) Output.

ROC Index:

ROC index on test for DT: 0.6638120400454645 ROC index on test for logistic regression: 0.6855684779738465 ROC index on test for NN: 0.6651260605503956 ROC index on test for Ensemble: 0.7089515045751483

Accuracy Score: Accuracy score on test for DT: 0.8194690265486726 Accuracy score on test for logistic regression: 0.6262268704746581 Accuracy score on test for NN: 0.8133547868061143

Ensemble train accuracy: 0.8962847070939122 Ensemble test accuracy: 0.8332260659694288 precision recall f1-score support

0	0.89	0.92	0.91	10820
1	0.31	0.24	0.27	1610

micro avg 0.83 0.83 0.83 12430 macro avg 0.60 0.58 0.59 12430 weighted avg 0.82 0.83 0.82 12430

Output: Probability produced by decision tree for each class vs actual prediction on Kick (0 = Good Buy, 1 = Bad Buy). You should be able to see the default threshold of 0.5. (Probs on zero) (probs on one) (prediction made) 0.8437146092865232

0.15628539071347677 0 1.0 0.0 0 0.8711538461538462 0.12884615384615383 0 0.7387466902030009 0.2612533097969991 0
0.8992568125516103 0.10074318744838975 0 0.7874015748031497 0.2125984251968504 0 0.9087301587301587
0.09126984126984126 0 0.8037225042301185 0.19627749576988154 0 0.9681818181818181 0.0318181818181815 0 0.9125
0.0875 0 0.7479674796747967 0.25203252032520324 0 1.0 0.0 0 0.9875968992248062 0.012403100775193798 0 1.0 0.0 0
0.5797720797720798 0.4202279202279202 0 0.8708399366085579 0.12916006339144215 0 0.7547169811320755
0.24528301886792453 0 0.9087301587301587 0.09126984126984126 0 0.731610337972167 0.268389662027833 0
0.9087301587301587 0.09126984126984126 0

b. Do all the models agree on the cars characteristics? How do they vary? Yes, because the accuracy of the ensemble model improved. So, most of the models agreed the cars characteristics. Prediction of Decision Tree, Neural Network and Logistic regression on Kicked data is same. This helped to increase the accuracy of the bagged aggregation

In [31]:

```
# import the model
from sklearn.ensemble import VotingClassifier

# initialise the classifier with 3 different estimators
voting = VotingClassifier(estimators=[('dt', dt_model), ('lr', log_reg_model), ('nn', nn_model)], voting='soft')
```

In [38]:

```
# Test Accuracy
y_pred_dt = dt_model.predict(X_test)
y_pred_log_reg = log_reg_model.predict(X_test)
y_pred_nn = nn_model.predict(X_test)
y_pred_proba_ensemble = voting.predict_proba(X_test)

print("Accuracy score on test for DT:", accuracy_score(y_test, y_pred_dt))
print("Accuracy score on test for logistic regression:", accuracy_score(y_test, y_pred_log_reg))
print("Accuracy score on test for NN:", accuracy_score(y_test, y_pred_nn))
```

Accuracy score on test for DT: 0.8194690265486726
Accuracy score on test for logistic regression: 0.6262268704746581
Accuracy score on test for NN: 0.8133547868061143

In [38]:

```
# typical prediction
y_pred = dt_model.predict(X_test)

# probability prediction from decision tree
y_pred_proba_dt = dt_model.predict_proba(X_test)

print("Probability produced by decision tree for each class vs actual prediction on Kick (0 = Good Buy, 1 = Bad Buy). You should be able to see the default threshold of 0.5.")
print("(Probs on zero)\t(probs on one)\t(prediction made)")
# print top 10
for i in range(20):
    print(y_pred_proba_dt[i][0], '\t', y_pred_proba_dt[i][1], '\t', y_pred[i])
```

Probability produced by decision tree for each class vs actual prediction on Kick (0 = Good Buy, 1 = Bad Buy). You should be able to see the default threshold of 0.5.

```

(probs on zero) (probs on one) (prediction made)
0.8437146092865232 0.15628539071347677 0
1.0 0.0 0
0.8711538461538462 0.12884615384615383 0
0.7387466902030009 0.2612533097969991 0
0.8992568125516103 0.10074318744838975 0
0.7874015748031497 0.2125984251968504 0
0.9087301587301587 0.09126984126984126 0
0.8037225042301185 0.19627749576988154 0
0.9681818181818181 0.031818181818181815 0
0.9125 0.0875 0
0.7479674796747967 0.25203252032520324 0
1.0 0.0 0
0.9875968992248062 0.012403100775193798 0
1.0 0.0 0
0.5797720797720798 0.4202279202279202 0
0.8708399366085579 0.12916006339144215 0
0.7547169811320755 0.24528301886792453 0
0.9087301587301587 0.09126984126984126 0
0.731610337972167 0.268389662027833 0
0.9087301587301587 0.09126984126984126 0

```

In [40]:

```

voting.fit(X_train, y_train)

# evaluate train and test accuracy
print("Ensemble train accuracy:", voting.score(X_train, y_train))
print("Ensemble test accuracy:", voting.score(X_test, y_test))
y_perd_ens=voting.predict(X_test)
print(classification_report(y_test,y_perd_ens))
# evaluate ROC auc score
y_pred_proba_ensemble = voting.predict_proba(X_test)
print()
roc_index_ensemble = roc_auc_score(y_test, y_pred_proba_ensemble[:, 1])
print("ROC score of voting classifier:", roc_index_ensemble)

```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:433: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
FutureWarning)

```

Ensemble train accuracy: 0.8962847070939122
Ensemble test accuracy: 0.8332260659694288

```

	precision	recall	f1-score	support
0	0.89	0.92	0.91	10820
1	0.31	0.24	0.27	1610
micro avg	0.83	0.83	0.83	12430
macro avg	0.60	0.58	0.59	12430
weighted avg	0.82	0.83	0.82	12430

ROC score of voting classifier: 0.7089515045751483

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\neural_network\multilayer_perceptron.py:562: ConvergenceWarning: Stochastic Optimizer: Maximum iterations (200) reached and the optimization hasn't converged yet.
% self.max_iter, ConvergenceWarning)

b. Use the comparison methods (or the comparison node) to compare the best decision tree model, the best regression model, the best neural network model and the ensemble model.

In [36]:

```

from sklearn.metrics import roc_auc_score

y_pred_proba_dt = dt_model.predict_proba(X_test)
y_pred_proba_log_reg = log_reg_model.predict_proba(X_test)
y_pred_proba_nn = nn_model.predict_proba(X_test)
y_pred_proba_ensemble = voting.predict_proba(X_test)

roc_index_dt = roc_auc_score(y_test, y_pred_proba_dt[:, 1])

```



```

roc_index_dt = roc_auc_score(y_test, y_pred_proba_dt[:, 1])
roc_index_log_reg = roc_auc_score(y_test, y_pred_proba_log_reg[:, 1])
roc_index_nn = roc_auc_score(y_test, y_pred_proba_nn[:, 1])
roc_index_ensemble = roc_auc_score(y_test, y_pred_proba_ensemble[:, 1])

print("ROC index on test for DT:", roc_index_dt)
print("ROC index on test for logistic regression:", roc_index_log_reg)
print("ROC index on test for NN:", roc_index_nn)
print("ROC index on test for Ensemble:", roc_index_ensemble)

from sklearn.metrics import roc_curve

fpr_dt, tpr_dt, thresholds_dt = roc_curve(y_test, y_pred_proba_dt[:,1])
fpr_log_reg, tpr_log_reg, thresholds_log_reg = roc_curve(y_test, y_pred_proba_log_reg[:,1])
fpr_nn, tpr_nn, thresholds_nn = roc_curve(y_test, y_pred_proba_nn[:,1])
fpr_ensemble, tpr_ensemble, thresholds_ensemble = roc_curve(y_test, y_pred_proba_ensemble[:,1])

import matplotlib.pyplot as plt

plt.plot(fpr_dt, tpr_dt, label='ROC Curve for DT {:.3f}'.format(roc_index_dt), color='red', lw=0.5)
plt.plot(fpr_log_reg, tpr_log_reg, label='ROC Curve for Log reg {:.3f}'.format(roc_index_log_reg),
color='green', lw=0.5)
plt.plot(fpr_nn, tpr_nn, label='ROC Curve for NN {:.3f}'.format(roc_index_nn), color='darkorange',
lw=0.5)
plt.plot(fpr_ensemble, tpr_ensemble, label='ROC Curve for Ensemble model {:.3f}'.format(roc_index_ensemble), color='blue', lw=0.5)

plt.plot([0, 1], [0, 1], color='navy', lw=0.5, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.0])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic example')
plt.legend(loc="lower right")
plt.show()

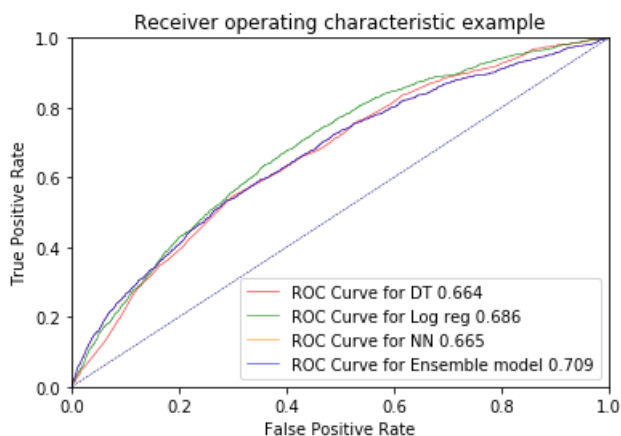
y_pred_proba_ensemble =

```

```

ROC index on test for DT: 0.6638120400454645
ROC index on test for logistic regression: 0.6855684779738465
ROC index on test for NN: 0.6651260605503956
ROC index on test for Ensemble: 0.7089515045751483

```



In [35]:

```

y_pred = dt_model.predict(X_test)

# probability prediction from decision tree
y_pred_proba_dt = dt_model.predict_proba(X_test)

print("Probability produced by decision tree for each class vs actual prediction on Kick (0 = Good Buy, 1 = Bad Buy). You should be able to see the default threshold of 0.5.")
print("(Probs on zero)\t(probs on one)\t(prediction made)")
# print top 10
for i in range(20):
    print(y_pred_proba_dt[i][0], '\t', y_pred_proba_dt[i][1], '\t', y_pred[i])

```


Probability produced by decision tree for each class vs actual prediction on Kick (0 = Good Buy, 1 = Bad Buy). You should be able to see the default threshold of 0.5.

(Probs on zero) (probs on one) (prediction made)

0.8437146092865232	0.15628539071347677	0
1.0	0.0	0
0.8711538461538462	0.12884615384615383	0
0.7387466902030009	0.2612533097969991	0
0.8992568125516103	0.10074318744838975	0
0.7874015748031497	0.2125984251968504	0
0.9087301587301587	0.09126984126984126	0
0.8037225042301185	0.19627749576988154	0
0.9681818181818181	0.031818181818181815	0
0.9125	0.0875	0
0.7479674796747967	0.25203252032520324	0
1.0	0.0	0
0.9875968992248062	0.012403100775193798	0
1.0	0.0	0
0.5797720797720798	0.4202279202279202	0
0.8708399366085579	0.12916006339144215	0
0.7547169811320755	0.24528301886792453	0
0.9087301587301587	0.09126984126984126	0
0.731610337972167	0.268389662027833	0
0.9087301587301587	0.09126984126984126	0

Task 6. Final Remarks: Decision Making

1. Finally, based on all models and analysis, is there a particular model you will use in decision making? Justify your choice.

Ensemble model is best suitable method for Decision making because a. We are using the majority vote to predict the Kicked cars. Base models such as Decision Tree, Logistic Regression, neural Network. We are aggregating all these models and based on the majority vote from these models we can get the best results. It can also be referred as Bagging.

1. Can you summarise positives and negatives of each predictive modelling method based on this analysis?

Decision Tree: Positive: a. Decision rules are intuitive. b. Can handle the non-linear features. c. It will take the variable interactions to predict the output. Negative: a. It is little bit biased towards the training set. b. Easy to get overfit. But we already used the Ensemble model to get rid of this problem.

Logistic Regression: Positive: a. Convenient probability scores for observations. b. We can regularize the features using L1 and L2. We don't need to worry much about the feature correlation. We applied the regularisation while doing the GridSearch CV. It increased the accuracy a little bit. c. Probabilistic interpretation is good, unlike decision trees or SVMs.

Negative: a. Not good to handle the non-linear features b. Doesn't handle large number of categorical features/variables well.

Neural Network: Positive: a. Can handle the non-linear features well. b. Negative: a. Difficult to interpret the model. b. Taking

1. How the outcome of this study can be used by decision makers?

Predictive models such as decision tree, regression model, neural network and ensemble model will help the car dealers in purchasing a used car at an auto auction. Kicked cars cause problem to the car dealers because it will be expensive if the car is not sold to the customer. We are using the features from the data set to predict the kicked cars using the classification approach. 0 Denotes that the car is not kicked and 1 denotes that the car is kicked. Accuracy from these models plays the vital role in prediction of kicked cars. Best model or more accurate model helps to save the money to car dealers. They will use this model to predict whether the car has potential to be "kicked" or not.