

In [1]:

```
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 PIL import Image
import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
import seaborn as sns
import matplotlib.pyplot as plt
```

In [2]:

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

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

In [3]:

```
## 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
```

In [4]:

```
#deleting the unwanted features
def feature_engineering(df):

    del df['WheelTypeID'] #Wheeltype is used. WheelTypeID is derived from WheelType
    del df['PurchaseID'] # Just a serial number.
    del df['ForSale'] #Just inclined to yes excluding 6 records-No contribution to the prediction
    del df['PurchaseDate'] #We are using Vehicle year to measure the time series and hence, deleting it.
    del df['MMRCurrentRetailRatio'] #Derived from MMRCurrentRetailAveragePrice and MMRCurrentRetailCleanPrice
    del df['PRIMEUNIT'] # more than 80% of the values are undefined(?)
    del df['AUCGUART'] # more than 80% of the values are undefined(?)
```

```
del df['PurchaseTimestamp'] # We are using year as a way of time measure, derived from
PurchaseDate. We are deleting both of them as we have year to measure.
```

```
return df
```

In [5]:

```
def data_type_change(df):

    #Assigning nominal values to either binary or numerical to support the datatype change
    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)

    # changing datatypes as required
    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['VehBCost'] = df['VehBCost'].astype(float)
    df['WheelType'] = df['WheelType'].astype(int)
    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['MMRAcquisitonRetailCleanPrice'] = df['MMRAcquisitonRetailCleanPrice'].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
```

In [6]:

```
# IsOnlineSale is a binary varibale which accepts wither zero or one. So we are replacing other va
```

lues wiht nan, which will be rpelaced by the median

```
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 [7]:

```
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
```

In [8]:

```
def preprocessing():

    df = pd.read_csv('CaseStudyData.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)
    error_replacing_For_IsOnlineSale(new_df)
    fill_missing_values(new_df)
    feature_engineering(new_df)

    return new_df
```

In [9]:

```
df = pd.read_csv('CaseStudyData.csv')
df = preprocessing()
df2 = data_type_change(df)

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:3020: DtypeWarning: Columns (27) have mixed types. Specify dtype option on import or set low_memory=False.
interactivity=interactivity, compiler=compiler, result=result)
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 not hasattr(df, 'numeric_types'):

```

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
WheelType             41432 non-null int32
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(20), int32(1), int64(2)
memory usage: 7.4 MB
None

```

Logistic Regression

In [10]:

```

from sklearn.preprocessing import StandardScaler

# initialise a standard scaler object
scaler = StandardScaler()

# visualise min, max, mean and standard dev of data before scaling
print("Before scaling\n-----")
for i in range(5):
    col = X_train[:,i]
    print("Variable #{}: min {}, max {}, mean {:.2f} and std dev {:.2f}".
          format(i, min(col), max(col), np.mean(col), np.std(col)))

# learn the mean and std.dev of variables from training data
# then use the learned values to transform training data
X_train = scaler.fit_transform(X_train, y_train)

print("After scaling\n-----")
for i in range(5):
    col = X_train[:,i]
    print("Variable #{}: min {}, max {}, mean {:.2f} and std dev {:.2f}".
          format(i, min(col), max(col), np.mean(col), np.std(col)))

# use the statistic that you learned from training to transform test data
# NEVER learn from test data, this is supposed to be a set of dataset
# that the model has never seen before
X_test = scaler.transform(X_test)

```

Before scaling

```

-----
Variable #0: min 0.0, max 2.0, mean 0.90 and std dev 0.64
Variable #1: min 2001.0, max 2010.0, mean 2005.07 and std dev 1.72
Variable #2: min 0.0, max 29.0, mean 9.23 and std dev 6.39
Variable #3: min 0.0, max 15.0, mean 7.99 and std dev 4.76
Variable #4: min 0.0, max 1.0, mean 0.03 and std dev 0.17
After scaling
-----

```

```

Variable #0: min -1.3952671651113548, max 1.7110295528962218, mean -0.00 and std dev 1.00

```

Variable #1: min -2.3629034833274973, max 2.8572877618480654, mean 0.00 and std dev 1.00
Variable #2: min -1.444416052719967, max 3.0925817631025105, mean 0.00 and std dev 1.00
Variable #3: min -1.6790550888103855, max 1.4718306939308479, mean 0.00 and std dev 1.00
Variable #4: min -0.19697746890836196, max 5.799330211492291, mean 0.00 and std dev 1.00

```
C:\ProgramData\Anaconda3\lib\site-packages\ipykernel_launcher.py:26: DataConversionWarning: Data with input dtype int32, int64, float64 were all converted to float64 by StandardScaler.
```

Task 3. Predictive Modeling Using Regression

1. In preparation for regression, is any imputation of missing values needed for this data set? List the variables that needed this.

Answer: Logistic regression is used to describe data and to explain the relationship between one dependent binary variable(IsBadBuy) and one or more nominal, ordinal, interval or ratio-level independent variables(remaining variables). We converted undefined values into nan values and then nan values into median and mode for Numerical and Nominal respectively. Since most of the features are having the null values and undefined values. Imputation is needed.

1. Apply transformation method(s) to the variable(s) that need it. List the variables that needed it. We used the Log transformation to remove the skewness in the variables. Regression models are sensitive to extreme or outlying values in the input space. Inputs with highly skewed or kurtotic distributions are often selected over inputs with better overall predictions. This approach will help to improve the performance of the model.

Since most of the variables are facing the skewness. we applied the Log transformation for all the input features. The skewness is determined using Boxplot. That is already mentioned in the data preprocessing part. Based on the analysis from the box plot. We found skewness in all input features and applied the log transformation on all the features.

Features used for the Log-transformation is: Auction VehYear Make Color Transmission WheelType VehOdo Nationality Size TopThreeAmericanName MMRAcquisitionAuctionAveragePrice MMRAcquisitionAuctionCleanPrice MMRAcquisitionRetailAveragePrice MMRAcquisitionRetailCleanPrice MMRCurrentAuctionAveragePrice MMRCurrentAuctionCleanPrice MMRCurrentRetailAveragePrice MMRCurrentRetailCleanPrice VNST VehBCost

1. Build a regression model using the default regression method with all inputs. Once you done it, build another one and tune it using GridSearchCV. Answer

h. Name the regression function used. We are using Sigmoid function which converts continuous values into binary values to predict the target.

i. How much was the difference in performance of two models build, default and optimal?

Default:

Train accuracy: 0.6340555313502594 Test accuracy: 0.6270313757039421

GridSearch CV:

Train accuracy: 0.6340951400166357 Test accuracy: 0.6270313757039421

{'C': 100}

Hyper parameter we used is C, which denotes the inverse of regularisation strength. Smaller C means stronger regularisation. It will help to balance the weight associated with each feature and minimize the overfitting.

j. Show the set parameters for the best model. What are the parameters used? Explain your decision. What are the optimal parameters?

The Hyper parameter used is C which is inverse of regularisation strength. Gridsearch will take the params as the input parameter and it will search iteratively from 10^{-6} to 10^4 . It will find the best C from the above range. Random state will help to keep the state consistent. n_jobs will enable the parallel computing. It will help to make use of all the resource of CPU to compute the Logistic regression model.

The Optimal parameters that we obtained is {'C': 100}

k. Report which variables are included in the regression model.

Answer: The important features which are used for the regression model is shown below.

Auction VehYear Make Color Transmission WheelType VehOdo Nationality Size TopThreeAmericanName MMRAcquisitionAuctionAveragePrice MMRAcquisitionAuctionCleanPrice MMRAcquisitionRetailAveragePrice MMRAcquisitionRetailCleanPrice MMRCurrentAuctionAveragePrice MMRCurrentAuctionCleanPrice MMRCurrentRetailAveragePrice MMRCurrentRetailCleanPrice VNST VehBCost

l. Report the top-5 important variables (in the order) in the model.

After using the Decision Tree model to select the MMRAcquisitionRetailAveragePrice MMRAcquisitionAuctionAveragePrice MMRCurrentAuctionAveragePrice MMRCurrentRetailAveragePrice MMRAcquisitionRetailCleanPrice

m. What is classification accuracy on training and test datasets? Train accuracy: 0.6770688923468666 Test accuracy: 0.6810350518186019

n. Report any sign of overfitting. Answer: We can see only slight difference between the train and test accuracy. So we can say that there is no overfitting as test data is having better accuracy than the train data.

In [97]:

```
#Default Logistic Regression
from sklearn.linear_model import LogisticRegression

model = LogisticRegression(random_state=rs)

# fit it to training data
model.fit(X_train, y_train)

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

# classification report on test data
y_pred = model.predict(X_test)
print(classification_report(y_test, y_pred))
```

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.

```
Train accuracy: 0.6340555313502594
Test accuracy: 0.6270313757039421
```

	precision	recall	f1-score	support
0	0.92	0.62	0.74	10820
1	0.21	0.65	0.31	1610
micro avg	0.63	0.63	0.63	12430
macro avg	0.56	0.64	0.53	12430
weighted avg	0.83	0.63	0.69	12430

In [98]:

```
print(model.coef_)
```

```
[[-0.05863466 -0.38087898  0.05053854  0.00854311 -0.0865566  -0.36904942
  0.05398771  0.11876354 -0.01123924 -0.14643453 -1.38591341  0.34128111
  1.63067541 -0.46056358  1.37398981 -0.44912909 -0.78586113 -0.12636073
 -0.00648164 -0.31814296  0.00518234 -0.02038268]]
```

In [99]:

```
feature_names = X.columns
coef = model.coef_[0]

# limit to 20 features, you can comment the following line to print out everything
coef = coef[:20]

for i in range(len(coef)):
    print(feature_names[i], ': ', coef[i])
```

```
Auction : -0.0586346628926053
VehYear : -0.38087897708425467
Make : 0.050538538184015554
Color : 0.008543110459752938
Transmission : -0.08655659948246425
WheelType : -0.36904942460927653
```

```
--
VehOdo : 0.05398771382540754
Nationality : 0.11876354104189467
Size : -0.011239239234750809
TopThreeAmericanName : -0.14643453244524487
MMRAcquisitionAuctionAveragePrice : -1.385913414644366
MMRAcquisitionAuctionCleanPrice : 0.34128110737186523
MMRAcquisitionRetailAveragePrice : 1.6306754144144062
MMRAcquisitionRetailCleanPrice : -0.46056358195057817
MMRCurrentAuctionAveragePrice : 1.3739898116071951
MMRCurrentAuctionCleanPrice : -0.44912908997325846
MMRCurrentRetailAveragePrice : -0.7858611333621506
MMRCurrentRetailCleanPrice : -0.12636072580169702
VNST : -0.0064816412846934015
VehBCost : -0.3181429621599786
```

In [100]:

```
# grab feature importances from the model and feature name from the original X
coef = model.coef_[0]
feature_names = X.columns

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

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

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

```
MMRAcquisitionRetailAveragePrice : 1.6306754144144062
MMRAcquisitionAuctionAveragePrice : -1.385913414644366
MMRCurrentAuctionAveragePrice : 1.3739898116071951
MMRCurrentRetailAveragePrice : -0.7858611333621506
MMRAcquisitionRetailCleanPrice : -0.46056358195057817
MMRCurrentAuctionCleanPrice : -0.44912908997325846
VehYear : -0.38087897708425467
WheelType : -0.36904942460927653
MMRAcquisitionAuctionCleanPrice : 0.34128110737186523
VehBCost : -0.3181429621599786
TopThreeAmericanName : -0.14643453244524487
MMRCurrentRetailCleanPrice : -0.12636072580169702
Nationality : 0.11876354104189467
Transmission : -0.08655659948246425
Auction : -0.0586346628926053
VehOdo : 0.05398771382540754
Make : 0.050538538184015554
WarrantyCost : -0.02038267664509494
Size : -0.011239239234750809
Color : 0.008543110459752938
```

In [101]:

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

# use all cores to tune logistic regression with C parameter
cv = GridSearchCV(param_grid=params, estimator=LogisticRegression(random_state=rs), cv=10,
n_jobs=-1)
cv.fit(X_train, y_train)

# test the best model
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 parameters of the best model
print(cv.best_params_)
```

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.

```
Train accuracy: 0.6340951400166357
Test accuracy: 0.6270313757039421
      precision    recall  f1-score   support

0         0.92        0.62        0.74       10820
1         0.21        0.65        0.31        1610

 micro avg       0.63        0.63        0.63       12430
 macro avg       0.56        0.64        0.53       12430
weighted avg       0.83        0.63        0.69       12430

{'C': 100}
```

1. Build another regression model using the subset of inputs selected by RFE and selection by model method. Answer the followings:

a. Report which variables are included in the regression model. Answer: Below are the variables included in the regression model. WheelType Color Auction TopThreeAmericanName VehYear Make Size VehBCost VehOdo VNST Nationality MMRCurrentRetailAveragePrice MMRCurrentRetailCleanPrice MMRAcquisitionAuctionCleanPrice MMRAcquisitionRetailAveragePrice MMRAcquisitionAuctionAveragePrice MMRCurrentAuctionAveragePrice WarrantyCost MMRAcquisitionRetailCleanPrice

b. Report the top-5 important variables (in the order) in the model.

Answer: Below are the top five important variables in the model. We are analysin feature importance from the tuned decision tree against log transformed features.

WheelType Color Auction TopThreeAmericanName VehYear

c. What are the parameters used? Explain your choices. What are the optimal parameters? Which regression function is being used?

Answer: We are passing the hyperparameter C, logistic regression model and numberOfJobs as parameters to the gridsearch model for the best model. We can see that in the below cells. Here C is the inverse of regularisation, which will help to avoid the overfitting. Typical values for C range from 10^{-6} to 10^4 . Smaller C means stronger regularisation. To speedup the searching process in GridsearchCV, we are enabling the parallel processing using numberOfJobs = -1.

Optimal parameters for this model are: Hyperparamter C=100.

d. Report any sign of overfitting. Answer: We can see only slight difference between the train and test accuracy. So we can say that there is no overfitting as test data is having better than the train data. Train accuracy: 0.6773235252510963 Test accuracy: 0.6821572381015248

Testing accuracy is Good. Hence, there is no sign of overfitting.

e. What is classification accuracy on training and test datasets?

The accuracy we obtained after Log transformation and Reccursive Feature elimination is shown below.

Train accuracy: 0.6770688923468666 Test accuracy: 0.6810350518186019

f. Did it improve/worsen the performance? Explain why those changes may have happened. When compare to default Accuracy, The accuracy we obtained after using Reccursive Feature elimination is improved a little bit. Before REF, we eliminated the skewness from all the features. REF will choose the important features by assigning weights to individual features. based on that it will eliminate the unwanted features. This will also reduce the space and time complexity. hence, it will help to improve the accuracy of the model.

In [102]:

```
def plot_skewed_columns(df):
    # setting up subplots for easier visualisation
    f, axes = plt.subplots(2,4, figsize=(10,10), sharex=False)

    # gift avg plots
    sns.distplot(df['VehYear'].dropna(), hist=False, ax=axes[0,0])
    sns.distplot(df['WheelType'].dropna(), hist=False, ax=axes[0,1])
    sns.distplot(df['TopThreeAmericanName'].dropna(), hist=False, ax=axes[1,0])
    sns.distplot(df['Nationality'].dropna(), hist=False, ax=axes[1,1])

    # gift avg plots
```



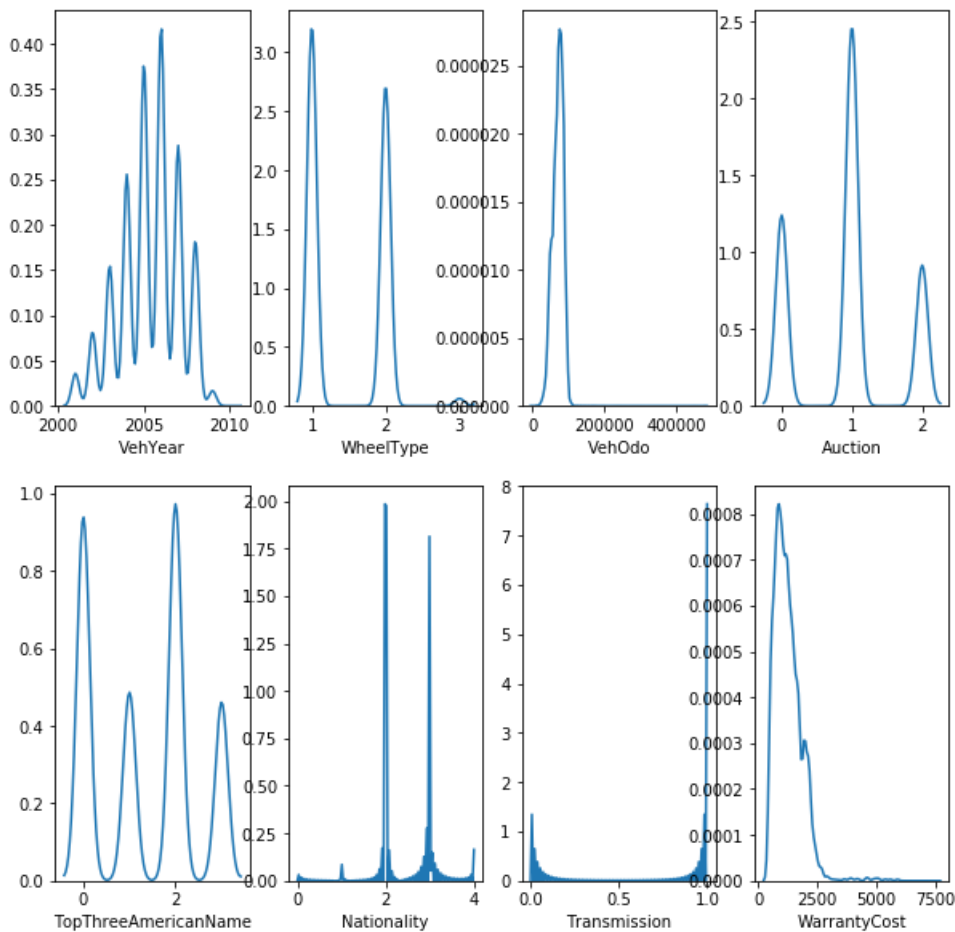
```
# give the plots
sns.distplot(df['VehOdo'].dropna(), hist=False, ax=axes[0,2])
sns.distplot(df['Auction'].dropna(), hist=False, ax=axes[0,3])
sns.distplot(df['Transmission'].dropna(), hist=False, ax=axes[1,2])
sns.distplot(df['WarrantyCost'].dropna(), hist=False, ax=axes[1,3])

plt.show()
```

```
plot_skewed_columns(df2)
```

C:\ProgramData\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning:

Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.



In [103]:

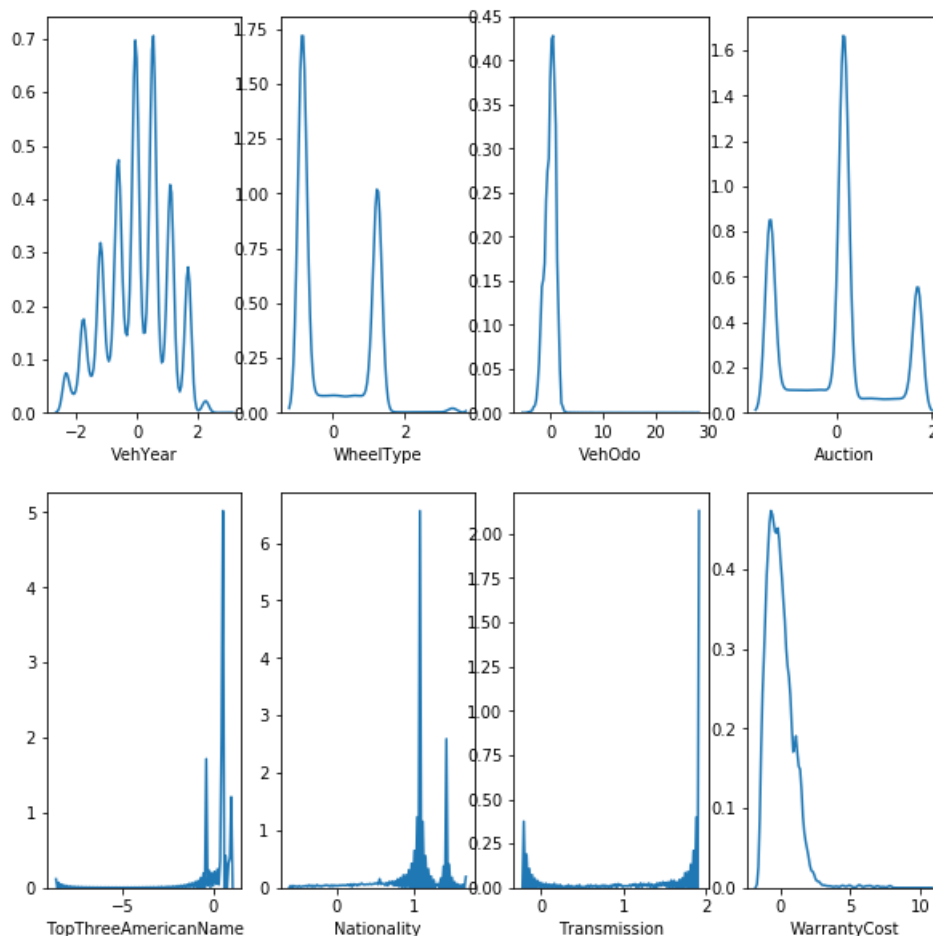
```
import numpy as np

# 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()
# 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]))
# plot them again to show the distribution
plot_skewed_columns(df_log)
```



In [104]:

```
# 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 [105]:

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

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

# test the best model
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 parameters of the best model
print(cv.best_params_)
```

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.

```
Train accuracy: 0.6770688923468666
Test accuracy: 0.6810350518186019
      precision    recall  f1-score   support

     0         0.67       0.70       0.69         7574
     1         0.69       0.66       0.68         7575

 micro avg       0.68       0.68       0.68        15149
 macro avg       0.68       0.68       0.68        15149
 weighted avg    0.68       0.68       0.68        15149

{'C': 100}
```

1. Using the best regression model, which cars could potential be “kicks”? Can you provide some descriptive summary of those cars?

Wheeltype of the car, majorly contributing to predict the kick cars. Then Based on the color also we can predict the kick cars. Majorly Wheeltype,color, Auction company name, TopThreeAmericanName and Vehyear are contributing to predict the kick cars.As Vehyear is increasing the cars not a kick car.

In [106]:

```
from sklearn.feature_selection import RFECV

rfe = RFECV(estimator = LogisticRegression(random_state=rs), cv=10)
rfe.fit(X_train, y_train) # run the RFECV

# comparing how many variables before and after
print("Original feature set", X_train.shape[1])
print("Number of features after elimination", rfe.n_features )
```

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

```
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.

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.

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.

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.

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.

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.

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.
```

Original feature set 22
Number of features after elimination 14

In [107]:

```
X_train_sel = rfe.transform(X_train)
X_test_sel = rfe.transform(X_test)
```

In [108]:

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

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

# test the best model
print("Train accuracy:", cv.score(X_train_sel, y_train))
print("Test accuracy:", cv.score(X_test_sel, y_test))

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

# print parameters of the best model
print(cv.best_params_)
```

```
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.
```

```
Train accuracy: 0.6367687249970293
Test accuracy: 0.6283990345937248
      precision    recall  f1-score   support
```

In [109]:

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

```
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.

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.

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.

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.

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.

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.

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.

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.

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.
```

Original feature set 22
Number of features after elimination 15

```
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.
```

```
Train accuracy: 0.6773235252510963
Test accuracy: 0.6821572381015248
```

	precision	recall	f1-score	support
0	0.68	0.70	0.69	7574
1	0.69	0.66	0.68	7575
micro avg	0.68	0.68	0.68	15149
macro avg	0.68	0.68	0.68	15149
weighted avg	0.68	0.68	0.68	15149

```
{'C': 10}
```

```
-----
AttributeError                                Traceback (most recent call last)
<ipython-input-109-c6260fd2812d> in <module>
    26 print(cv.best_params_)
    27
--> 28 analyse_feature_importance(cv.best_estimator_, X_log.columns)

<ipython-input-93-4b6f5664d395> in analyse_feature_importance(dm_model, feature_names,
n_to_display)
     1 def analyse_feature_importance(dm_model, feature_names, n_to_display=20):
     2     # grab feature importances from the model
--> 3     importances = dm_model.feature_importances_
     4
     5     # sort them out in descending order

AttributeError: 'LogisticRegression' object has no attribute 'feature_importances_'
```

Feature selection using Decision Tree model

In [112]:

```
from sklearn.tree import DecisionTreeClassifier

#Setting the Parameters for Decision Tree like Maximum depth, and minimim number of leaf per sample
params = {'criterion': ['gini', 'entropy'],
          'max_depth': range(4, 10),
          'min_samples_leaf': range(10,15,2)}

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

print(cv.best_params_)
```

```
{'criterion': 'gini', 'max_depth': 9, 'min_samples_leaf': 12}
```

In [113]:

```
# analyse feature importance from the tuned decision tree against log transformed X
analyse_feature_importance(cv.best_estimator_, X_log.columns)
```

```
WheelType : 0.4034058826987275
Color : 0.2285520490774161
Auction : 0.10993086355875598
TopThreeAmericanName : 0.10372067635790408
VehYear : 0.05350481335247447
Make : 0.0478307481097652
Size : 0.015864243347691584
VehBCost : 0.012702916836217463
VehOdo : 0.006029399698476508
VNST : 0.0036936869238576213
Nationality : 0.003142391353828301
MMRCurrentRetailAveragePrice : 0.0030955310222805415
MMRCurrentRetailCleanPrice : 0.003012583437579162
MMRAcquisitionAuctionCleanPrice : 0.002538609293902061
MMRAcquisitionRetailAveragePrice : 0.0008225093954484635
MMRAcquisitionAuctionAveragePrice : 0.0007690042164613255
MMRCurrentAuctionAveragePrice : 0.0007011664052984512
WarrantyCost : 0.0004108694624065102
MMRAcquisitonRetailCleanPrice : 0.00027205545150884084
IsOnlineSale : 0.0
```

In [114]:

```
from sklearn.feature_selection import SelectFromModel

# use the trained best decision tree from GridSearchCV to select features
# supply the prefit=True parameter to stop SelectFromModel to re-train the model
selectmodel = SelectFromModel(cv.best_estimator_, prefit=True)
X_train_sel_model = selectmodel.transform(X_train_log)
X_test_sel_model = selectmodel.transform(X_test_log)

print(X_train_sel_model.shape)
```

```
(35345, 6)
```

In [115]:

```
params = {'C': [pow(10, x) for x in range(-6, 4)]}

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

print("Train accuracy:", cv.score(X_train_sel_model, y_train_log))
print("Test accuracy:", cv.score(X_test_sel_model, y_test_log))
```

```
# test the best model
y_pred = cv.predict(X_test_sel_model)
print(classification_report(y_test_log, y_pred))

# print parameters of the best model
print(cv.best_params_)
```

Train accuracy: 0.6339510538972981

Test accuracy: 0.6383919730675293

	precision	recall	f1-score	support
0	0.64	0.64	0.64	7574
1	0.64	0.63	0.64	7575
micro avg	0.64	0.64	0.64	15149
macro avg	0.64	0.64	0.64	15149
weighted avg	0.64	0.64	0.64	15149

{'C': 1e-06}

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.