```
1   #Importing all the required libraries
2
3   import os
4   import pandas as pd
5   import numpy as np
6   import matplotlib.pyplot as plt
7   from scipy.stats import chi2_contingency
8   import seaborn as sns
9   from random import randrange, uniform
```

```
1   #Reading the train data
2
3   train_dataset = pd.read_csv("/content/train_Df64byy.csv")
4   train_dataset.head()
```

|   | ID | City_Code | Region_Code | Accomodation_Type | Reco_Insurance_Type | Upper_Age | Lower_Age | Is_Spouse | Health Indicator | Holc |
|---|----|-----------|-------------|-------------------|---------------------|-----------|-----------|-----------|------------------|------|
| 0 | 1  | C3        | 3213        | Rented            | Individual          | 36        | 36        | No        | X1               |      |
| 1 | 2  | C5        | 1117        | Owned             | Joint               | 75        | 22        | No        | X2               |      |
| 2 | 3  | C5        | 3732        | Owned             | Individual          | 32        | 32        | No        | NaN              |      |
| 3 | 4  | C24       | 4378        | Owned             | Joint               | 52        | 48        | No        | X1               |      |
| 4 | 5  | C8        | 2190        | Rented            | Individual          | 44        | 44        | No        | X2               |      |

```
1   #Describing the train data
2
3   train_dataset.describe()
```

| | ID | Region_Code | Upper_Age | Lower_Age | Holding_Policy_Type | Reco_Policy_Cat | Reco_Policy_Premium |
|---|---|---|---|---|---|---|---|
| **count** | 66720.000000 | 66720.000000 | 66720.000000 | 66720.000000 | 40123.000000 | 66720.000000 | 66720.000000 |
| **mean** | 29600.349460 | 1733.354332 | 44.846927 | 42.734937 | 2.442664 | 15.112125 | 14183.940833 |
| **std** | 15003.069766 | 1424.021443 | 17.299650 | 17.310907 | 1.025191 | 6.340442 | 6585.290312 |
| **min** | 1.000000 | 1.000000 | 18.000000 | 16.000000 | 1.000000 | 1.000000 | 2280.000000 |
| **25%** | 16679.750000 | 527.000000 | 28.000000 | 27.000000 | 1.000000 | 12.000000 | 9252.000000 |

```
1    #Returning the data types of train
2
3    train_dataset.dtypes
```

```
ID                        int64
City_Code                 object
Region_Code               float64
Accomodation_Type         object
Reco_Insurance_Type       object
Upper_Age                 float64
Lower_Age                 float64
Is_Spouse                 object
Health Indicator          object
Holding_Policy_Duration   object
Holding_Policy_Type       float64
Reco_Policy_Cat           float64
Reco_Policy_Premium       float64
Response                  float64
dtype: object
```

```
1    #Returning the no of classified and non classified
2
3    print("No of observations classified as 1 are : ", len(train_dataset[train_dataset['Response']==1]))
4
5    print("No of observations classified as 0 are : ", len(train_dataset[train_dataset['Response']==0]))
```

```
No of observations classified as 1 are :  16080
No of observations classified as 0 are :  50640
```

```
1   #Checking no of missing values in train data
2
3   missing_val = pd.DataFrame(train_dataset.isnull().sum())
4   missing_val
```

|  | 0 |
|---|---|
| ID | 0 |
| City_Code | 1 |
| Region_Code | 1 |
| Accomodation_Type | 1 |
| Reco_Insurance_Type | 1 |
| Upper_Age | 1 |
| Lower_Age | 1 |
| Is_Spouse | 1 |
| Health Indicator | 11692 |
| Holding_Policy_Duration | 20252 |
| Holding_Policy_Type | 20252 |
| Reco_Policy_Cat | 1 |
| Reco_Policy_Premium | 1 |
| Response | 1 |

```
1   #Returning the value counts for each unique value in Holding_Policy_Duration Column
2
3   train_dataset['Holding_Policy_Duration'].value_counts()

    1.0     5889
    14+     5647
    2.0     5582
    3.0     4724
    4.0     3612
```

```
5.0     3103
6.0     2495
7.0     2161
8.0     1726
9.0     1451
10.0    1070
11.0     709
13.0     688
12.0     659
14.0     607
Name: Holding_Policy_Duration, dtype: int64
```

```
1    #Replacing 14+ value in Holding_Policy_Duration with 15
2
3    # mapping = {'14+':'15'}
4    train_dataset['Holding_Policy_Duration'] = train_dataset['Holding_Policy_Duration'].replace("14+", "15")
```

```
1    #Returning the value counts for each unique value in Health_Indicator Column
2
3    train_dataset["Health Indicator"].value_counts()
```

```
X1    17087
X2    13522
X3     8890
X4     7550
X5     2242
X6     1677
X7      258
X8      105
X9       86
Name: Health Indicator, dtype: int64
```

```
1    #Returning the value counts for each unique value in Holding_Policy_Type Column
2
3    train_dataset["Holding_Policy_Type"].value_counts()
```

```
3.0    17432
1.0    10647
2.0     6556
```

```
       4.0     5488
       Name: Holding_Policy_Type, dtype: int64
```

```
1   #Checking the percentage of missing values for train data
2
3   missing_val = missing_val.reset_index()
4   missing_val = missing_val.rename(columns = {'index':'Variables', 0: 'Missing_percentage'})
5   missing_val['Missing_percentage'] = (missing_val['Missing_percentage']/len(train_dataset))*100
6   missing_val = missing_val.sort_values('Missing_percentage', ascending = False).reset_index(drop = True)
7   missing_val
```

|    | Variables | Variables | Missing_percentage |
|----|-----------|-----------|--------------------|
| 0  | 1         | City_Code | 0.000003 |
| 1  | 2         | Region_Code | 0.000003 |
| 2  | 3         | Accomodation_Type | 0.000003 |
| 3  | 4         | Reco_Insurance_Type | 0.000003 |
| 4  | 5         | Upper_Age | 0.000003 |
| 5  | 6         | Lower_Age | 0.000003 |
| 6  | 7         | Is_Spouse | 0.000003 |
| 7  | 11        | Reco_Policy_Cat | 0.000003 |
| 8  | 12        | Reco_Policy_Premium | 0.000003 |
| 9  | 13        | Response | 0.000003 |
| 10 | 0         | ID | 0.000000 |
| 11 | 8         | Health Indicator | 0.000000 |
| 12 | 9         | Holding_Policy_Duration | 0.000000 |
| 13 | 10        | Holding_Policy_Type | 0.000000 |

```
1   #Imputing the missing values for required columns
2
```

```
2
3    from sklearn.impute import SimpleImputer
4
5    impute_size = SimpleImputer(strategy = "most_frequent")
6    train_dataset['Health Indicator'] = impute_size.fit_transform(train_dataset[['Health Indicator']])
7    train_dataset['Holding_Policy_Duration'] = impute_size.fit_transform(train_dataset[['Holding_Policy_Duration']])
8    train_dataset['Holding_Policy_Type'] = train_dataset['Holding_Policy_Type'].fillna(train_dataset['Holding_Policy_Type']
```

```
1    #Checking the percentage of missing values for train data after imputing the missing values
2
3    missing_val = pd.DataFrame(train_dataset.isnull().sum())
4    missing_val = missing_val.reset_index()
5    missing_val = missing_val.rename(columns = {'index':'Variables', 0: 'Missing_percentage'})
6    missing_val['Missing_percentage'] = (missing_val['Missing_percentage']/len(train_dataset))*100
7    missing_val = missing_val.sort_values('Missing_percentage', ascending = False).reset_index(drop = True)
8    missing_val
```

| Variables | Missing_percentage |
|---|---|

```
1  #Still we can some missing values percentage, so we have to drop columns which have nan value
2
3  train_dataset.dropna(how='any')
```

| | ID | City_Code | Region_Code | Accomodation_Type | Reco_Insurance_Type | Upper_Age | Lower_Age | Is_Spouse | Health Indicato |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | C3 | 3213.0 | Rented | Individual | 36.0 | 36.0 | No | ⟩ |
| 1 | 2 | C5 | 1117.0 | Owned | Joint | 75.0 | 22.0 | No | ⟩ |
| 2 | 3 | C5 | 3732.0 | Owned | Individual | 32.0 | 32.0 | No | ⟩ |
| 3 | 4 | C24 | 4378.0 | Owned | Joint | 52.0 | 48.0 | No | ⟩ |
| 4 | 5 | C8 | 2190.0 | Rented | Individual | 44.0 | 44.0 | No | ⟩ |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 50877 | 35040 | C1 | 102.0 | Rented | Joint | 48.0 | 47.0 | Yes | ⟩ |
| 50878 | 35041 | C1 | 332.0 | Owned | Individual | 70.0 | 70.0 | No | ⟩ |
| 50879 | 35042 | C6 | 1165.0 | Owned | Joint | 57.0 | 56.0 | Yes | ⟩ |
| 50880 | 35043 | C11 | 1032.0 | Owned | Joint | 72.0 | 69.0 | Yes | ⟩ |
| 50881 | 35044 | C1 | 729.0 | Rented | Individual | 25.0 | 25.0 | No | ⟩ |

50882 rows × 14 columns

```
1  #Checking the percentage of missing values for train data after dropping the missing values
2
3  missing_val = pd.DataFrame(train_dataset.isnull().sum())
4  missing_val = missing_val.reset_index()
5  missing_val = missing_val.rename(columns = {'index':'Variables', 0: 'Missing_percentage'})
6  missing_val['Missing_percentage'] = (missing_val['Missing_percentage']/len(train_dataset))*100
7  missing_val1 = missing_val.sort_values('Missing_percentage', ascending = False).reset_index(drop = True)
8  missing_val1
```

| | Variables | Missing_percentage |
|---|---|---|
| **0** | ID | 0.0 |
| **1** | City_Code | 0.0 |
| **2** | Region_Code | 0.0 |
| **3** | Accomodation_Type | 0.0 |
| **4** | Reco_Insurance_Type | 0.0 |
| **5** | Upper_Age | 0.0 |
| **6** | Lower_Age | 0.0 |
| **7** | Is_Spouse | 0.0 |
| **8** | Health Indicator | 0.0 |
| **9** | Holding_Policy_Duration | 0.0 |
| **10** | Holding_Policy_Type | 0.0 |
| **11** | Reco_Policy_Cat | 0.0 |
| **12** | Reco_Policy_Premium | 0.0 |
| **13** | Response | 0.0 |

**I'm trying to consider only 'Reco_Policy_Premium' column for outliers because it's the only column which deals with money and all other columns belongs to ID, Age Duration Etc**

```
1  #Process to detect outliers with the help of quartile
2
3  max_thresold = train_dataset['Reco_Policy_Premium'].quantile(0.95)
4  max_thresold
```

```
26852.0
```

```
1  #Checking the data which are outliers with max_thresold
```

```
2
3  train_dataset[train_dataset['Reco_Policy_Premium'] > max_thresold]
```

| | ID | City_Code | Region_Code | Accomodation_Type | Reco_Insurance_Type | Upper_Age | Lower_Age | Is_Spouse | Healt Indicato |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | C5 | 1117.0 | Owned | Joint | 75.0 | 22.0 | No | ⟩ |
| 7 | 8 | C1 | 3175.0 | Owned | Joint | 75.0 | 73.0 | Yes | ⟩ |
| 8 | 9 | C15 | 3497.0 | Owned | Joint | 52.0 | 43.0 | No | ⟩ |
| 48 | 49 | C2 | 2858.0 | Owned | Joint | 57.0 | 55.0 | Yes | ⟩ |
| 49 | 50 | C1 | 85.0 | Owned | Joint | 73.0 | 68.0 | Yes | ⟩ |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 50793 | 34956 | C1 | 3014.0 | Owned | Joint | 47.0 | 46.0 | Yes | ⟩ |
| 50798 | 34961 | C2 | 484.0 | Owned | Joint | 69.0 | 68.0 | Yes | ⟩ |
| 50812 | 34975 | C1 | 907.0 | Owned | Joint | 75.0 | 72.0 | Yes | ⟩ |
| 50842 | 35005 | C1 | 4370.0 | Owned | Joint | 71.0 | 68.0 | No | ⟩ |
| 50852 | 35015 | C1 | 1383.0 | Owned | Joint | 59.0 | 58.0 | Yes | ⟩ |

2544 rows × 14 columns

```
1  #Process to detect outliers with the help of quartile
2
3  min_thresold = train_dataset['Reco_Policy_Premium'].quantile(0.05)
4  min_thresold
```

5224.0

```
1  #Checking the data which are outliers with min_thresold
2
3  train_dataset[train_dataset['Reco_Policy_Premium'] < min_thresold]
```

| | ID | City_Code | Region_Code | Accomodation_Type | Reco_Insurance_Type | Upper_Age | Lower_Age | Is_Spouse | Health Indicato |
|---|---|---|---|---|---|---|---|---|---|
| 10 | 11 | C28 | 600.0 | Owned | Individual | 21.0 | 21.0 | No | › |
| 15 | 16 | C3 | 1484.0 | Rented | Individual | 20.0 | 20.0 | No | › |
| 22 | 23 | C25 | 787.0 | Rented | Individual | 18.0 | 18.0 | No | › |
| 27 | 28 | C9 | 855.0 | Rented | Individual | 21.0 | 21.0 | No | › |
| 46 | 47 | C3 | 1475.0 | Rented | Individual | 21.0 | 21.0 | No | › |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 50775 | 34938 | C3 | 1246.0 | Owned | Individual | 19.0 | 19.0 | No | › |
| 50814 | 34977 | C6 | 155.0 | Rented | Individual | 28.0 | 28.0 | No | › |
| 50822 | 34985 | C1 | 1984.0 | Rented | Individual | 19.0 | 19.0 | No | › |
| 50836 | 34999 | C3 | 1418.0 | Rented | Individual | 20.0 | 20.0 | No | › |
| 50874 | 35037 | C27 | 3277.0 | Owned | Individual | 29.0 | 29.0 | No | › |

2537 rows × 14 columns

```
1   train_dataset = train_dataset[(train_dataset['Reco_Policy_Premium'] < max_thresold) & (train_dataset['Reco_Policy_Premi
```

```
1   train_dataset
```

| | ID | City_Code | Region_Code | Accomodation_Type | Reco_Insurance_Type | Upper_Age | Lower_Age | Is_Spouse | Health Indicato |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | C3 | 3213.0 | Rented | Individual | 36.0 | 36.0 | No | ⟩ |
| 2 | 3 | C5 | 3732.0 | Owned | Individual | 32.0 | 32.0 | No | ⟩ |
| 3 | 4 | C24 | 4378.0 | Owned | Joint | 52.0 | 48.0 | No | ⟩ |
| 4 | 5 | C8 | 2190.0 | Rented | Individual | 44.0 | 44.0 | No | ⟩ |
| 5 | 6 | C9 | 1785.0 | Rented | Individual | 52.0 | 52.0 | No | ⟩ |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |

```
1  #Returning the no of classified and non classified after imputing missing values and removing outliers from train data
2
3  print("No of observations classified as 1 are : ", len(train_dataset[train_dataset['Response'] == 1]))
4
5  print("No of observations classified as 0 are : ", len(train_dataset[train_dataset['Response'] == 0]))
```

```
No of observations classified as 1 are :  16080
No of observations classified as 0 are :  50640
```

45770 rows × 14 columns

```
1  #Describing the test_data
2  test_dataset =pd.read_csv('/content/test_YCcRUnU.csv')
3  test_dataset.describe()
```

| ID | Region_Code | Upper_Age | Lower_Age | Holding_Policy_Type | Reco_Policy_Cat | Reco_Policy_Premium |
|----|-------------|-----------|-----------|---------------------|-----------------|---------------------|

```
1    #Checking the data types for test data
2    test_dataset = pd.read_csv("/content/test_YCcRUnU.csv")
3
4    test_dataset.dtypes
```

```
ID                         int64
City_Code                  object
Region_Code                int64
Accomodation_Type          object
Reco_Insurance_Type        object
Upper_Age                  int64
Lower_Age                  int64
Is_Spouse                  object
Health Indicator           object
Holding_Policy_Duration    object
Holding_Policy_Type        float64
Reco_Policy_Cat            int64
Reco_Policy_Premium        float64
dtype: object
```

```
1    #Checking the percentage of missing values for test data
2
3    missing_val = pd.DataFrame(test_dataset.isnull().sum())
4    missing_val = missing_val.reset_index()
5    missing_val = missing_val.rename(columns = {'index':'Variables', 0: 'Missing_percentage'})
6    missing_val['Missing_percentage'] = (missing_val['Missing_percentage']/len(test_dataset))*100
7    missing_val = missing_val.sort_values('Missing_percentage', ascending = False).reset_index(drop = True)
8    missing_val
```

|   | Variables | Missing_percentage |
|---|---|---|
| **0** | Holding_Policy_Duration | 39.454254 |
| **1** | Holding_Policy_Type | 39.454254 |
| **2** | Health Indicator | 23.054345 |
| **3** | ID | 0.000000 |
| **4** | City_Code | 0.000000 |
| **5** | Region_Code | 0.000000 |
| **6** | Accomodation_Type | 0.000000 |
| **7** | Reco_Insurance_Type | 0.000000 |

```
1   #Returning the value counts for each unique value in Holding_Policy_Duration Column
2
3   test_dataset['Holding_Policy_Duration'].value_counts()
```

```
14+     1892
1.0     1891
2.0     1772
3.0     1606
4.0     1205
5.0      992
6.0      903
7.0      664
8.0      569
9.0      493
10.0     333
11.0     254
13.0     221
14.0     211
12.0     196
Name: Holding_Policy_Duration, dtype: int64
```

```
1   #Replacing 14+ value with 15 for 'Holding_Policy_Duration' column
2
3   test_dataset['Holding_Policy_Duration'] = test_dataset['Holding_Policy_Duration'].replace("14+", "15")
```

```
1    #Returning the value counts for each unique value in Health Indicator Column
2
3    test_dataset['Health Indicator'].value_counts()
```

```
X1    5614
X2    4516
X3    2846
X4    2442
X5     681
X6     514
X7      96
X8      41
X9      28
Name: Health Indicator, dtype: int64
```

```
1    #Returning the value counts for each unique value in Holding_Policy_Type Column
2
3    test_dataset['Holding_Policy_Type'].value_counts()
```

```
3.0    5572
1.0    3574
2.0    2150
4.0    1906
Name: Holding_Policy_Type, dtype: int64
```

```
1    #Imputing missing values for test data
2
3    impute_size = SimpleImputer(strategy = "most_frequent")
4    test_dataset['Health Indicator'] = impute_size.fit_transform(test_dataset[['Health Indicator']])
5    test_dataset['Holding_Policy_Duration'] = impute_size.fit_transform(test_dataset[['Holding_Policy_Duration']])
6    test_dataset['Holding_Policy_Type'] = test_dataset['Holding_Policy_Type'].fillna(test_dataset['Holding_Policy_Type'].me
```

```
1    #Checking the percentage of missing values for test data after imputing the missing values
2
3    missing_val = pd.DataFrame(test_dataset.isnull().sum())
4    missing_val = missing_val.reset_index()
5    missing_val = missing_val.rename(columns = {'index':'Variables', 0: 'Missing_percentage'})
```

```
6  missing_val['Missing_percentage'] = (missing_val['Missing_percentage']/len(test_dataset))*100
7  missing_val = missing_val.sort_values('Missing_percentage', ascending = False).reset_index(drop = True)
8  missing_val
```

| | Variables | Missing_percentage |
|---|---|---|
| **0** | ID | 0.0 |
| **1** | City_Code | 0.0 |
| **2** | Region_Code | 0.0 |
| **3** | Accomodation_Type | 0.0 |
| **4** | Reco_Insurance_Type | 0.0 |
| **5** | Upper_Age | 0.0 |
| **6** | Lower_Age | 0.0 |
| **7** | Is_Spouse | 0.0 |
| **8** | Health Indicator | 0.0 |
| **9** | Holding_Policy_Duration | 0.0 |
| **10** | Holding_Policy_Type | 0.0 |
| **11** | Reco_Policy_Cat | 0.0 |
| **12** | Reco_Policy_Premium | 0.0 |

```
1  from sklearn.preprocessing import LabelEncoder, OneHotEncoder
2
3  le = LabelEncoder()
4  objList = train_dataset.select_dtypes(include="object").columns
5
6  for feat in objList:
7      train_dataset[feat] = le.fit_transform(train_dataset[feat].astype(str))
```

```
1  train_dataset
```

| | ID | City_Code | Region_Code | Accomodation_Type | Reco_Insurance_Type | Upper_Age | Lower_Age | Is_Spouse | Health Indicato |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 22 | 3213 | 1 | 0 | 36 | 36 | 0 | |
| **1** | 2 | 31 | 1117 | 0 | 1 | 75 | 22 | 0 | |
| **2** | 3 | 31 | 3732 | 0 | 0 | 32 | 32 | 0 | |
| **3** | 4 | 16 | 4378 | 0 | 1 | 52 | 48 | 0 | |
| **4** | 5 | 34 | 2190 | 1 | 0 | 44 | 44 | 0 | |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| **66715** | 50878 | 30 | 845 | 1 | 0 | 22 | 22 | 0 | |
| **66716** | 50879 | 31 | 4188 | 1 | 0 | 27 | 27 | 0 | |
| **66717** | 50880 | 0 | 442 | 1 | 0 | 63 | 63 | 0 | |
| **66718** | 50881 | 0 | 4 | 0 | 1 | 71 | 49 | 0 | |
| **66719** | 50882 | 22 | 3866 | 1 | 0 | 24 | 24 | 0 | |

66720 rows × 14 columns

```
1   from sklearn.preprocessing import LabelEncoder, OneHotEncoder
2
3   le = LabelEncoder()
4   objList = test_dataset.select_dtypes(include="object").columns
5
6   for feat in objList:
7       test_dataset[feat] = le.fit_transform(test_dataset[feat].astype(str))
```

```
1   test_dataset.head()
```

| | ID | City_Code | Region_Code | Accomodation_Type | Reco_Insurance_Type | Upper_Age | Lower_Age | Is_Spouse | Health Indicator | H |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 50883 | 0 | 156 | 0 | 0 | 30 | 30 | 0 | 0 | |
| **1** | 50884 | 30 | 7 | 0 | 1 | 69 | 68 | 1 | 0 | |

```
1  X = train_dataset.drop("Response", axis=1)    #Feature Matrix
2  y = train_dataset["Response"]
```

| **4** | 50887 | 0 | 951 | 0 | 0 | 75 | 75 | 0 | 2 | |

```
1  X.head()
```

| | ID | City_Code | Region_Code | Accomodation_Type | Reco_Insurance_Type | Upper_Age | Lower_Age | Is_Spouse | Health Indicator | Hold |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 22 | 3213 | 1 | 0 | 36 | 36 | 0 | 0 | |
| **2** | 3 | 31 | 3732 | 0 | 0 | 32 | 32 | 0 | 0 | |
| **3** | 4 | 16 | 4378 | 0 | 1 | 52 | 48 | 0 | 0 | |
| **4** | 5 | 34 | 2190 | 1 | 0 | 44 | 44 | 0 | 1 | |
| **5** | 6 | 35 | 1785 | 1 | 0 | 52 | 52 | 0 | 1 | |

```
1  # separate dataset into train and test
2  from sklearn.model_selection import train_test_split
3  X_train, X_val, y_train, y_val = train_test_split(
4      X,
5      y,
6      test_size=0.2,
7      random_state=39)
8
9  X_train.shape, X_val.shape

((53376, 13), (13344, 13))
```

## ▾ Handling Imablance

```
1   from sklearn.datasets import make_classification
2   from sklearn.linear_model import LogisticRegression
3   from sklearn.dummy import DummyClassifier
4   from sklearn.model_selection import train_test_split
5   from sklearn.metrics import roc_curve
6   from sklearn.metrics import roc_auc_score
7   from matplotlib import pyplot
8
9   # plot no skill and model roc curves
10  def plot_roc_curve(test_y, naive_probs, model_probs):
11    # plot naive skill roc curve
12    fpr, tpr, _ = roc_curve(test_y, naive_probs)
13    pyplot.plot(fpr, tpr, linestyle='--', label='No Skill')
14    # plot model roc curve
15    fpr, tpr, _ = roc_curve(test_y, model_probs)
16    pyplot.plot(fpr, tpr, marker='.', label='Logistic')
17    # axis labels
18    pyplot.xlabel('False Positive Rate')
19    pyplot.ylabel('True Positive Rate')
20    # show the legend
21    pyplot.legend()
22    # show the plot
23    pyplot.show()
24  from sklearn.metrics import precision_recall_curve
```

```
1   # no skill model, stratified random class predictions
2   model = DummyClassifier(strategy='stratified')
3   model.fit(X_train,y_train)
4   yhat = model.predict_proba(X_val)
5   pos_probs = yhat[:, 1]
6   # calculate the precision-recall auc
7   precision, recall, _ = precision_recall_curve(y_val, pos_probs)
8   auc_score = auc(recall, precision)
9   print('No Skill PR AUC: %.3f' % auc_score)
```

```
No Skill PR AUC: 0.333
```

```
1   # skilled model
2   model = LogisticRegression(solver='lbfgs')
3   model.fit(X_train, y_train)
4   yhat = model.predict_proba(X_val)
5   model_probs = yhat[:, 1]
6   # calculate roc auc
7   roc_auc = roc_auc_score(y_val, model_probs)
8   print('Logistic ROC AUC %.3f' % roc_auc)
9
```

```
Logistic ROC AUC 0.501
```

```
1   # example of evaluating a decision tree with random oversampling
2   from numpy import mean
3   from sklearn.datasets import make_classification
4   from sklearn.model_selection import cross_val_score
5   from sklearn.model_selection import RepeatedStratifiedKFold
6   from sklearn.tree import DecisionTreeClassifier
7   from imblearn.pipeline import Pipeline
8   from imblearn.over_sampling import RandomOverSampler
```

```
/usr/local/lib/python3.7/dist-packages/sklearn/externals/six.py:31: FutureWarning: The module is deprecated in version
  "(https://pypi.org/project/six/).", FutureWarning)
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:144: FutureWarning: The sklearn.neighbors.base modu
  warnings.warn(message, FutureWarning)
```

```
1   # define pipeline
2   steps = [('over', RandomOverSampler()), ('model', DecisionTreeClassifier())]
3   pipeline = Pipeline(steps=steps)
4   # evaluate pipeline
5   cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
6   scores = cross_val_score(pipeline, X, y, scoring='f1_micro', cv=cv, n_jobs=-1)
7   score = mean(scores)
```

```
  8    print('F-measure: %.3f' % score)
```

F-measure: 0.668

```
  1    from numpy import where
  2    from collections import Counter
  3    from imblearn.over_sampling import SMOTE
  4    # summarize class distribution
  5    counter = Counter(y)
  6    print(counter)
  7    # scatter plot of examples by class label
  8    # transform the dataset
  9    oversample = SMOTE()
 10    X, y = oversample.fit_resample(X, y)
```

Counter({0: 34779, 1: 10991})
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function safe_indexing is deprec
  warnings.warn(msg, category=FutureWarning)

```
  1    # summarize the new class distribution
  2    counter = Counter(y)
  3    print(counter)
```

Counter({0: 34779, 1: 34779})

```
  1    # define model
  2    model = DecisionTreeClassifier()
  3    # evaluate pipeline
  4    cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
  5    scores = cross_val_score(model, X, y, scoring='roc_auc', cv=cv, n_jobs=-1)
  6    print('Mean ROC AUC: %.3f' % mean(scores))
```

Mean ROC AUC: 0.758

```
  1    # define pipeline
  2    steps = [('over', SMOTE()), ('model', DecisionTreeClassifier())]
```

```
3    pipeline = Pipeline(steps=steps)
4    # evaluate pipeline
5    cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
6    scores = cross_val_score(pipeline, X, y, scoring='roc_auc', cv=cv, n_jobs=-1)
7    print('Mean ROC AUC: %.3f' % mean(scores))
```

```
Mean ROC AUC: 0.758
```

```
1    from imblearn.over_sampling import BorderlineSMOTE
2    # summarize class distribution
3    counter = Counter(y)
4    print(counter)
5    # transform the dataset
6    oversample = BorderlineSMOTE()
7    X, y = oversample.fit_resample(X_train, y_train)
8    # summarize the new class distribution
9    counter = Counter(y)
10   # summarize the new class distribution
11   counter = Counter(y)
12   print(counter)
```

```
Counter({0: 9900, 1: 100})
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function safe_indexing is deprec
  warnings.warn(msg, category=FutureWarning)
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function safe_indexing is deprec
  warnings.warn(msg, category=FutureWarning)
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function safe_indexing is deprec
  warnings.warn(msg, category=FutureWarning)
Counter({1: 27835, 0: 27835})
```

◀                                                                                                                    ▶

```
1    # combined SMOTE and Tomek Links sampling for imbalanced classification
2    from numpy import mean
3    from sklearn.datasets import make_classification
4    from sklearn.model_selection import cross_val_score
5    from sklearn.model_selection import RepeatedStratifiedKFold
6    from imblearn.pipeline import Pipeline
7    from sklearn.tree import DecisionTreeClassifier
8    from imblearn.combine import SMOTETomek
```

```
 9   from imblearn.under_sampling import TomekLinks
10
11   # define model
12   model = DecisionTreeClassifier()
13   # define sampling
14   resample = SMOTETomek(tomek=TomekLinks(sampling_strategy='majority'))
15   # define pipeline
16   pipeline = Pipeline(steps=[('r', resample), ('m', model)])
17   # define evaluation procedure
18   cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
19   # evaluate model
20   scores = cross_val_score(pipeline, X_train, y_train, scoring='roc_auc', cv=cv, n_jobs=-1)
21   # summarize performance
22   print('Mean ROC AUC: %.3f' % mean(scores))
```

    Mean ROC AUC: 0.529

```
 1   # fit a logistic regression model on an imbalanced classification dataset
 2   from numpy import mean
 3   from sklearn.datasets import make_classification
 4   from sklearn.model_selection import cross_val_score
 5   from sklearn.model_selection import RepeatedStratifiedKFold
 6   from sklearn.linear_model import LogisticRegression
 7
 8   # define model
 9   weights = {0:0.01, 1:1.0}
10   model = LogisticRegression(solver='lbfgs', class_weight=weights)
11   model = LogisticRegression(solver='lbfgs')
12   # define evaluation procedure
13   cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
14   # evaluate model
15   scores = cross_val_score(model, X_train, y_train, scoring='roc_auc', cv=cv, n_jobs=-1)
16   # summarize performance
17   print('Mean ROC AUC: %.3f' % mean(scores))
```

    Mean ROC AUC: 0.507

```
 1   # calculate heuristic class weighting
```

```
2    from sklearn.utils.class_weight import compute_class_weight
3    # calculate class weighting
4    weighting = compute_class_weight('balanced', [0,1], y_train)
5    print(weighting)
```

```
[0.65773307 2.08495616]
```

```
1    # define model
2    model = LogisticRegression(solver='lbfgs', class_weight='balanced')
3    # define evaluation procedure
4    cv = RepeatedStratifiedKFold(n_splits=25, n_repeats=10, random_state=1)
5    # evaluate model
6    scores = cross_val_score(model, X_train, y_train, scoring='roc_auc', cv=cv, n_jobs=-1)
7    # summarize performance
8    print('Mean ROC AUC: %.3f' % mean(scores))
```

```
Mean ROC AUC: 0.552
```

```
1    # define model
2    model = LogisticRegression(solver='lbfgs')
3    # define grid
4    balance = [{0:100,1:1}, {0:10,1:1}, {0:1,1:1}, {0:1,1:10}, {0:1,1:100}]
5    param_grid = dict(class_weight=balance)
6    # define evaluation procedure
7    cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
8    # define grid search
9    grid = GridSearchCV(estimator=model, param_grid=param_grid, n_jobs=-1, cv=cv,
10   scoring='roc_auc')
11   # execute the grid search
12   grid_result = grid.fit(X_train, y_train)
13   # report the best configuration
14   print('Best: %f using %s' % (grid_result.best_score_, grid_result.best_params_))
15   # report all configurations
16   means = grid_result.cv_results_['mean_test_score']
17   stds = grid_result.cv_results_['std_test_score']
18   params = grid_result.cv_results_['params']
19   for mean, stdev, param in zip(means, stds, params):
20       print('%f (%f) with: %r' % (mean, stdev, param))
```

```
          print( %i (%i) witm  %i   % (mcun, stucv, purum))
```

```
     Best: 0.561343 using {'class_weight': {0: 1, 1: 10}}
     0.478032 (0.009518) with: {'class_weight': {0: 100, 1: 1}}
     0.479372 (0.011027) with: {'class_weight': {0: 10, 1: 1}}
     0.506829 (0.009842) with: {'class_weight': {0: 1, 1: 1}}
     0.561343 (0.011775) with: {'class_weight': {0: 1, 1: 10}}
     0.555669 (0.011346) with: {'class_weight': {0: 1, 1: 100}}
```

```
 1    # fit a decision tree on an imbalanced classification dataset
 2    from numpy import mean
 3    # define model
 4    model = DecisionTreeClassifier()
 5    # define evaluation procedure
 6    cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
 7    # evaluate model
 8    scores = cross_val_score(model, X_train, y_train, scoring='roc_auc', cv=cv, n_jobs=-1)
 9    # summarize performance
10    print('Mean ROC AUC: %.3f' % mean(scores))
```

```
     Mean ROC AUC: 0.544
```

```
 1    # define model
 2    model = DecisionTreeClassifier(class_weight='balanced')
 3    # define evaluation procedure
 4    cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
 5    # evaluate model
 6    scores = cross_val_score(model, X_train, y_train, scoring='roc_auc', cv=cv, n_jobs=-1)
 7    # summarize performance
 8    print('Mean ROC AUC: %.3f' % mean(scores))
```

```
     Mean ROC AUC: 0.542
```

```
 1    #define model
 2    model = DecisionTreeClassifier()
 3    # define grid
 4    balance = [{0:100,1:1}, {0:10,1:1}, {0:1,1:1}, {0:1,1:10}, {0:1,1:100}]
 5    param_grid = dict(class_weight=balance)
 6    # define evaluation procedure
```

```
 6     # define evaluation procedure
 7     cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
 8     # define grid search
 9     grid = GridSearchCV(estimator=model, param_grid=param_grid, n_jobs=-1, cv=cv,
10     scoring='roc_auc')
11     # execute the grid search
12     grid_result = grid.fit(X_train, y_train)
13     # report the best configuration
14     print('Best: %f using %s' % (grid_result.best_score_, grid_result.best_params_))
15     # report all configurations
16     means = grid_result.cv_results_['mean_test_score']
17     stds = grid_result.cv_results_['std_test_score']
18     params = grid_result.cv_results_['params']
19     for mean, stdev, param in zip(means, stds, params):
20       print('%f (%f) with: %r' % (mean, stdev, param))
```

```
Best: 0.544592 using {'class_weight': {0: 10, 1: 1}}
0.542249 (0.008030) with: {'class_weight': {0: 100, 1: 1}}
0.544592 (0.011834) with: {'class_weight': {0: 10, 1: 1}}
0.543369 (0.009397) with: {'class_weight': {0: 1, 1: 1}}
0.536144 (0.008191) with: {'class_weight': {0: 1, 1: 10}}
0.536047 (0.006741) with: {'class_weight': {0: 1, 1: 100}}
```

```
1     # define model
2     model = SVC(gamma='scale', class_weight='balanced')
3     # define evaluation procedure
4     cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
5     # evaluate model
6     scores = cross_val_score(model, X, y, scoring='roc_auc', cv=cv, n_jobs=-1)
7     # summarize performance
8     print('Mean ROC AUC: %.3f' % mean(scores))
```

```
1     from sklearn.metrics import roc_auc_score
2     from keras.layers import Dense
3     from keras.models import Sequential
4     # define the neural network model
5     def define_model(n_input):
6     # define model
7       model = Sequential()
```

```
8    # define first hidden layer and visible layer
9      model.add(Dense(10, input_dim=n_input, activation='relu',
10     kernel_initializer='he_uniform'))
11   # define output layer
12     model.add(Dense(1, activation='sigmoid'))
13   # define loss and optimizer
14     model.compile(loss='binary_crossentropy', optimizer='sgd')
15     return model
```

```
1    # define the model
2    n_input = X_train.shape[1]
3    model = define_model(n_input)
```

```
1    # fit model
2    model.fit(X_train,y_train,epochs=100, verbose=1)
3    # make predictions on the test dataset
4    yhat = model.predict(y_val)
5    # evaluate the ROC AUC of the predictions
6    score = roc_auc_score(testy, yhat)
7    print('ROC AUC: %.3f' % score)
```

```
1    model = define_model(n_input)
2    # fit model
3    weights = {0:1, 1:100}
4    history = model.fit(X_train,y_train, class_weight=weights, epochs=100, verbose=1)
5    # evaluate model
6    yhat = model.predict(X_val)
7    score = roc_auc_score(y_val, yhat)
8    print('ROC AUC: %.3f' % score)

    Epoch 1/100
    1145/1145 [==============================] - 2s 2ms/step - loss: 12005825.8617
    Epoch 2/100
    1145/1145 [==============================] - 2s 2ms/step - loss: 3.3963
    Epoch 3/100
    1145/1145 [==============================] - 2s 2ms/step - loss: 3.3979
    Epoch 4/100
```

```
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3947
Epoch 5/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3975
Epoch 6/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3968
Epoch 7/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3958
Epoch 8/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3965
Epoch 9/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3964
Epoch 10/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3955
Epoch 11/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3959
Epoch 12/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3961
Epoch 13/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3973
Epoch 14/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3957
Epoch 15/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3970
Epoch 16/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3957
Epoch 17/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3955
Epoch 18/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3959
Epoch 19/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3939
Epoch 20/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3957
Epoch 21/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3974
Epoch 22/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3978
Epoch 23/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3959
Epoch 24/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3956
Epoch 25/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3973
```

```
Epoch 26/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3957
Epoch 27/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3961
Epoch 28/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3959
Epoch 29/100
1145/1145 [==============================] - 2s 2ms/step - loss: 3.3967
```

```
1    # fit xgboost on an imbalanced classification dataset
2    from numpy import mean
3    from sklearn.model_selection import cross_val_score
4    from sklearn.model_selection import RepeatedStratifiedKFold
5    from xgboost import XGBClassifier
```

```
1    # define model
2    model = XGBClassifier()
3    # define evaluation procedure
4    cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
5    # evaluate model
6    scores = cross_val_score(model, X_train, y_train, scoring='roc_auc', cv=cv, n_jobs=-1)
7    # summarize performance
8    print('Mean ROC AUC: %.5f' % mean(scores))
```

```
Mean ROC AUC: 0.62498
```

```
1    # define model
2    model = XGBClassifier(scale_pos_weight=99)
3    # define evaluation procedure
4    cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
5    # evaluate model
6    scores = cross_val_score(model, X_train, y_train, scoring='roc_auc', cv=cv, n_jobs=-1)
7    # summarize performance
8    print('Mean ROC AUC: %.5f' % mean(scores))
```

```
Mean ROC AUC: 0.62435
```

```
1    # define grid
```

```
 1    # define grid
 2    weights = [1, 10, 25, 50, 75, 99, 100, 1000]
 3    param_grid = dict(scale_pos_weight=weights)
 4    # define evaluation procedure
 5    cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
 6    # define grid search
 7    grid = GridSearchCV(estimator=model, param_grid=param_grid, n_jobs=-1, cv=cv,
 8    scoring='roc_auc')
 9    # execute the grid search
10    grid_result = grid.fit(X_train, y_train)
11    # report the best configuration
12    print("Best: %f using %s" % (grid_result.best_score_, grid_result.best_params_))
13    # report all configurations
14    means = grid_result.cv_results_['mean_test_score']
15    stds = grid_result.cv_results_['std_test_score']
16    params = grid_result.cv_results_['params']
17    for mean, stdev, param in zip(means, stds, params):
18      print("%f (%f) with: %r" % (mean, stdev, param))
```

```
 1    # class imbalnce with Gaussian
 2
 3    from numpy import sqrt
 4    from numpy import argmax
 5    from sklearn.datasets import make_classification
 6    from sklearn.linear_model import LogisticRegression
 7    from sklearn.model_selection import train_test_split
 8    from sklearn.metrics import roc_curve
 9    from matplotlib import pyplot
10
11
12    # fit a model
13    model = LogisticRegression(solver='lbfgs')
14    model.fit(X_train, y_train)
15    # predict probabilities
16    yhat = model.predict_proba(X_val)
17    # keep probabilities for the positive outcome only
18    yhat = yhat[:, 1]
19    # calculate roc curves
20    fpr, tpr, thresholds = roc_curve(y_val, yhat)
```
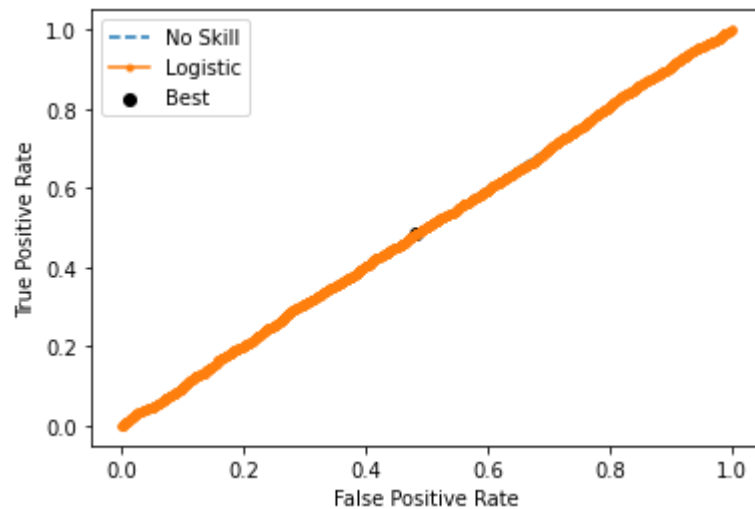
```
21    # calculate the g-mean for each threshold
22    gmeans = sqrt(tpr * (1-fpr))
23    # locate the index of the largest g-mean
24    ix = argmax(gmeans)
25    print('Best Threshold=%f, G-mean=%.3f' % (thresholds[ix], gmeans[ix]))
26    # plot the roc curve for the model
27    pyplot.plot([0,1], [0,1], linestyle='--', label='No Skill')
28    pyplot.plot(fpr, tpr, marker='.', label='Logistic')
29    pyplot.scatter(fpr[ix], tpr[ix], marker='o', color='black', label='Best')
30    # axis labels
31    pyplot.xlabel('False Positive Rate')
32    pyplot.ylabel('True Positive Rate')
33    pyplot.legend()
34    # show the plot
35    pyplot.show()
```

Best Threshold=0.252692, G-mean=0.501



```
1    # Class imbalance using SVM
2
3    from numpy import mean
4    from sklearn.datasets import make_classification
5    from sklearn.model_selection import cross_val_score
6    from sklearn.model_selection import RepeatedStratifiedKFold
7    from sklearn.calibration import CalibratedClassifierCV
```

```
7    from sklearn.calibration import CalibratedClassifierCV
8    from sklearn.svm import SVC
```

```
1    # define model
2    model = SVC(gamma='scale')
3    # wrap the model
4    calibrated = CalibratedClassifierCV(model, method='isotonic', cv=3)
5    # define evaluation procedure
6    cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
7    # evaluate model
8    scores = cross_val_score(calibrated, X_train, y_train, scoring='roc_auc', cv=cv, n_jobs=-1)
9    # summarize performance
10   print('Mean ROC AUC: %.3f' % mean(scores))
```

```
/usr/local/lib/python3.7/dist-packages/joblib/externals/loky/process_executor.py:691: UserWarning: A worker stopped whi
  "timeout or by a memory leak.", UserWarning
```

```
1    train_dataset
```

| | City_Code | Region_Code | Accomodation_Type | Reco_Insurance_Type | Upper_Age | Lower_Age | Is_Spouse | Health Indicator | Hol |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 22 | 3213.0 | 1 | 0 | 36.0 | 36.0 | 0 | 0 | |

```
1  y = train_dataset.values
2  #train_dataset.drop(['ID', 'Response'], inplace=True, axis=1)
3
4  x = train_dataset.values
5
```

```
1  from lightgbm import LGBMClassifier
2
3  lgbm = LGBMClassifier(random_state=5)
4
5  lgbm.fit(X_train, y_train)
6
7  y_pred = lgbm.predict(X_val)
```
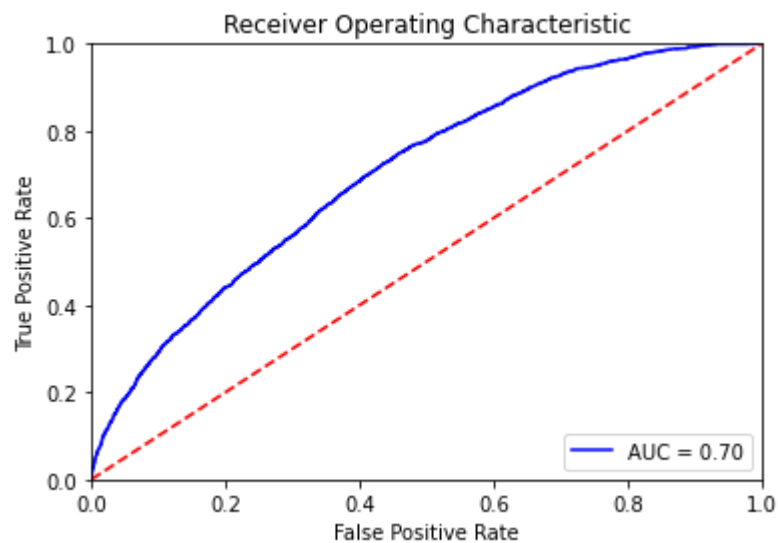
| | City_Code | Region_Code | Accomodation_Type | Reco_Insurance_Type | Upper_Age | Lower_Age | Is_Spouse | Health Indicator | Hol |
|---|---|---|---|---|---|---|---|---|---|
| **50881** | 0 | 729.0 | 1 | 0 | 25.0 | 25.0 | 0 | 0 | |

```
1  import sklearn.metrics as metrics
2  # calculate the fpr and tpr for all thresholds of the classification
3  probs = lgbm.predict_proba(X_val)
4  preds = probs[:,1]
5  fpr, tpr, threshold = metrics.roc_curve(y_val, preds)
6  roc_auc = metrics.auc(fpr, tpr)
7
8  # method I: plt
9  import matplotlib.pyplot as plt
10 plt.title('Receiver Operating Characteristic')
11 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
12 plt.legend(loc = 'lower right')
13 plt.plot([0, 1], [0, 1],'r--')
14 plt.xlim([0, 1])
15 plt.ylim([0, 1])
16 plt.ylabel('True Positive Rate')
17 plt.xlabel('False Positive Rate')
18 plt.show()
```

```
1   import sklearn.metrics as metrics
2   # calculate the fpr and tpr for all thresholds of the classification
3   probs = lgbm.predict_proba(test_dataset)
4   preds = probs[:,1]
5
```

```
1   ID = test_dataset['ID'].values
2
```

```
1   submission = pd.DataFrame({'ImageId':ID,
2                              'Response':preds,
3                             }).set_index('ImageId')
```

```
1   submission.head()
```

|  | Response |
|---|---|
| **ImageId** |  |
| **50883** | 0.221873 |
| **50884** | 0.368947 |

```
1   submission.to_csv('mnist_lgbm.csv', columns=['Response'])
```

These all observation.

1. Random Forest class_balanced will never impacted your model accuracy
2. class imbalnce with Gaussian never gives us right accuracy
3. Weights Paremeters doesn't work imbalance
4. Probabalistics methods also didn't solved the problem imbalance
5. Poor Pefroramace on STOME
6. Poor Perforomance on SMOTE +SMOTETomek
7. No performance impormace over SMOTE over sampling
8. U never think of Performance if your using Undersampling
7. Lots f Books and Online class room talks about imbalance data. but i may work on small sample.
8. XG Boost works some better
9. LightGBM Worked very well.its master piece for class imbalance.
10.Even trained Model on Deep learning No improvement performance with help of

Pycaret having paramete class_imbalance is True, But No effect on IT.even pycaret works well because it handles XGBOOST, LIGHTGBM perfe

AutoViml having class imbalance = SMOTE, I will never effect model ROC AUC score.

These all are my exprience with analystics vidya competition.

ned Model on Deep learning No improvement perforamance

Bottom line : Don't follow Bilindly on online platform posts. understand the core concept

Please find the experinments for the same.

## This is formatted as c **bold text**ode

https://neptune.ai/blog/lightgbm-parameters-guide

1