

Importing libraries and datasets

```
In [32]: import numpy as np
import pandas as pd
import datetime as dt
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split
from sklearn.cross_validation import RandomOverSampler
from imblearn.under_sampling import RandomOverSampler
from imblearn.over_sampling import SMOTE

from sklearn.linear_model import LogisticRegression
import xgboost as xgb
import tensorflow as tf

from sklearn.metrics import precision_score, recall_score, f1_score, auc, roc_auc_score, accuracy_score,
classification_report, confusion_matrix, roc_curve
from xgboost import plot_importance

In [33]: df_response = pd.read_csv('Retail_Data_Response.csv')
df_transactions = pd.read_csv('Retail_Data_Transactions.csv', parse_dates=['trans_date'])

In [4]: df_response.head()
```

customer_id	response
0	CS1112 0
1	CS1113 0
2	CS1114 1
3	CS1115 1
4	CS1116 1

```
In [5]: df_transactions.head()
```

customer_id	trans_date	trans_amount
0	CS5295	2013-02-11 35
1	CS2478	2015-03-15 39
2	CS2122	2013-02-26 52
3	CS1217	2011-11-16 99
4	CS1850	2013-11-20 78

```
In [6]: print(df_transactions['trans_date'].min())
print(df_transactions['trans_date'].max())

2011-05-16 00:00:00
2015-03-16 00:00:00
```

Data Preparation

```
In [7]: # since the last date of the data is 16 March 2015, the campaign date is assumed to be 17 March 2015
# RFM model will be used to predict campaign response. Recency is calculated

campaign_date = dt.datetime(2015,3,17)
df_transactions['recent'] = campaign_date - df_transactions['trans_date']
df_transactions['recent'].astype('timedelta64[D]')
df_transactions['recent'] = df_transactions['recent'] / np.timedelta64(1, 'D')
df_transactions.head()
```

customer_id	trans_date	trans_amount	recent
0	CS5295	2013-02-11	35
1	CS4768	2015-03-15	39
2	CS2122	2013-02-26	52
3	CS1217	2011-11-16	99
4	CS1850	2013-11-20	78

```
In [8]: # create data set with RFM variables

df_rfm = df_transactions.groupby('customer_id').agg({'recent': lambda x:x.min(), #
                                                    'recency': lambda x: len(x), #
                                                    'frequency': lambda x: x.sum(), #
                                                    'monetary_value': lambda x: x.sum(), #
                                                    'trans_date': lambda x: (x.max() - x.min()).days})

df_rfm.rename(columns={'recent': 'recency',
                       'customer_id': 'customer_id',
                       'trans_date': 'trans_date',
                       'trans_amount': 'monetary_value'}, inplace=True)

In [9]: df_rfm.head()
df_rfm.reset_index()
```

customer_id	recency	frequency	monetary_value
0	CS1112	620	15
1	CS1113	360	20
2	CS1114	330	19
3	CS1115	120	22
4	CS1116	2040	13

```
In [10]: # create data set with CVV variables

df_cvv = df_transactions.groupby('customer_id').agg({'recent': lambda x:x.min(), #
                                                    'recency': lambda x: len(x), #
                                                    'frequency': lambda x: x.sum(), #
                                                    'monetary_value': lambda x: x.sum(), #
                                                    'trans_date': lambda x: (x.max() - x.min()).days})

df_cvv.rename(columns={'recent': 'recency',
                       'customer_id': 'customer_id',
                       'trans_date': 'trans_date',
                       'trans_amount': 'monetary_value'}, inplace=True)

df_cvv['ticket_size'] = df_cvv['monetary_value'] / df_cvv['frequency']

In [11]: df_cvv.head()
df_cvv.reset_index()
```

customer_id	recency	frequency	monetary_value	AOU	ticket_size
0	CS1112	620	15	1012	1309
1	CS1113	360	20	1490	1354
2	CS1114	330	19	1432	1309
3	CS1115	120	22	1659	1303
4	CS1116	2040	13	857	1155

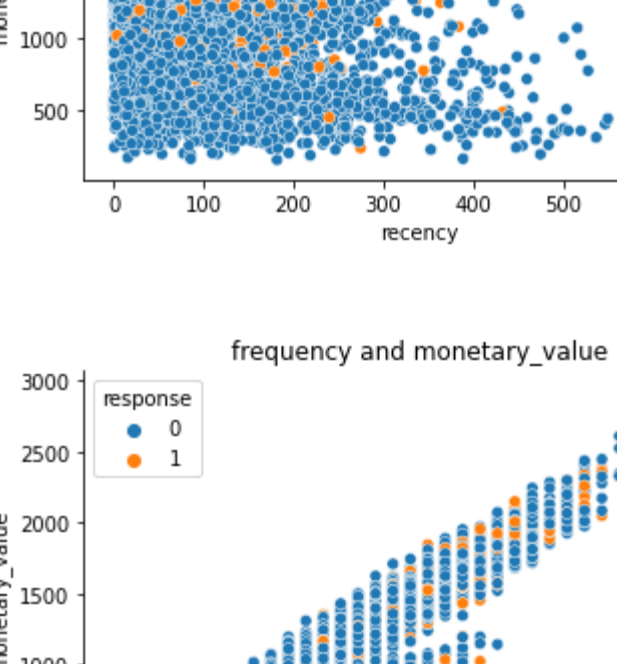
Calculating response rate

```
In [12]: response_rate = df_response.groupby('response').agg({'customer_id': lambda x: len(x)}).reset_index()
response_rate.head()
```

response	customer_id
0	6237
1	647

```
In [13]: plt.figure(figsize=(5,5))
x = range(2)
plt.bar(x, response_rate['customer_id'])
plt.xticks(response_rate.index)
plt.title('Response Distribution')
plt.xlabel('Convert or Not')
plt.ylabel('No. of Users')
plt.show()

# data is imbalanced
```



```
In [14]: # merging two data sets - RFM

df_modeling_rfm = pd.merge(df_response, df_rfm)
df_modeling_rfm.head()
```

customer_id	response	recency	frequency	monetary_value
0	CS1112	0	620	15
1	CS1113	0	360	20
2	CS1114	1	330	19
3	CS1115	1	120	22
4	CS1116	1	2040	13

```
In [15]: # merging two data sets - CVV

df_modeling_cvv = pd.merge(df_response, df_cvv)
df_modeling_cvv.head()
```

customer_id	response	recency	frequency	monetary_value	AOU	ticket_size
0	CS1112	0	620	15	1012	1309
1	CS1113	0	360	20	1490	1354
2	CS1114	1	330	19	1432	1309
3	CS1115	1	120	22	1659	1303
4	CS1116	1	2040	13	857	1155

Creating train and test dataset

```
In [16]: # splitting dataframe into X and y

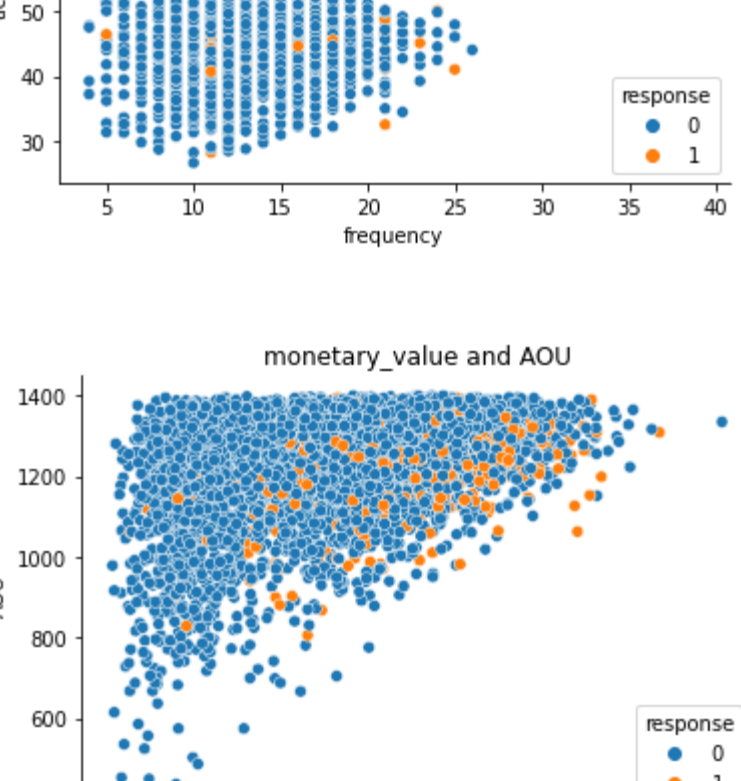
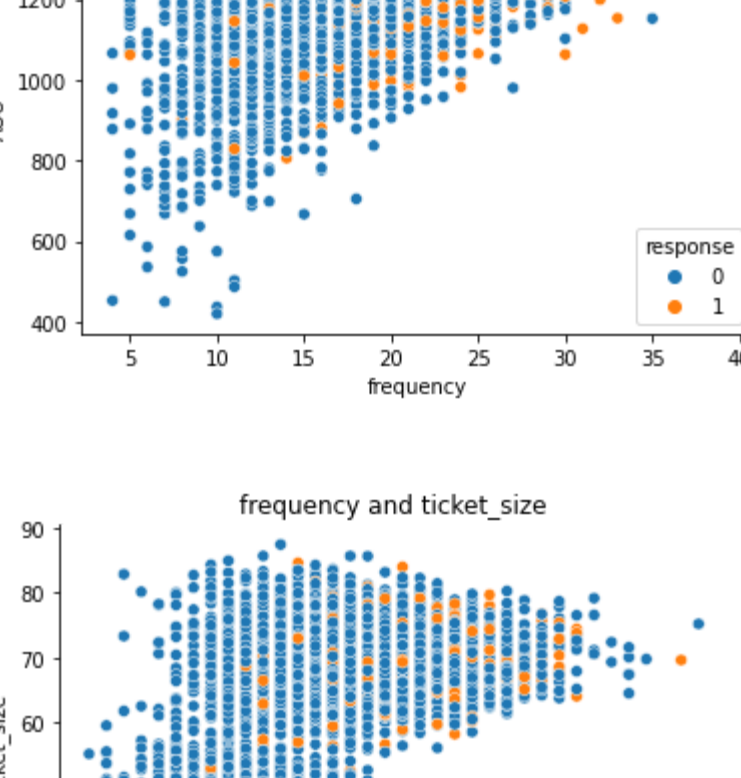
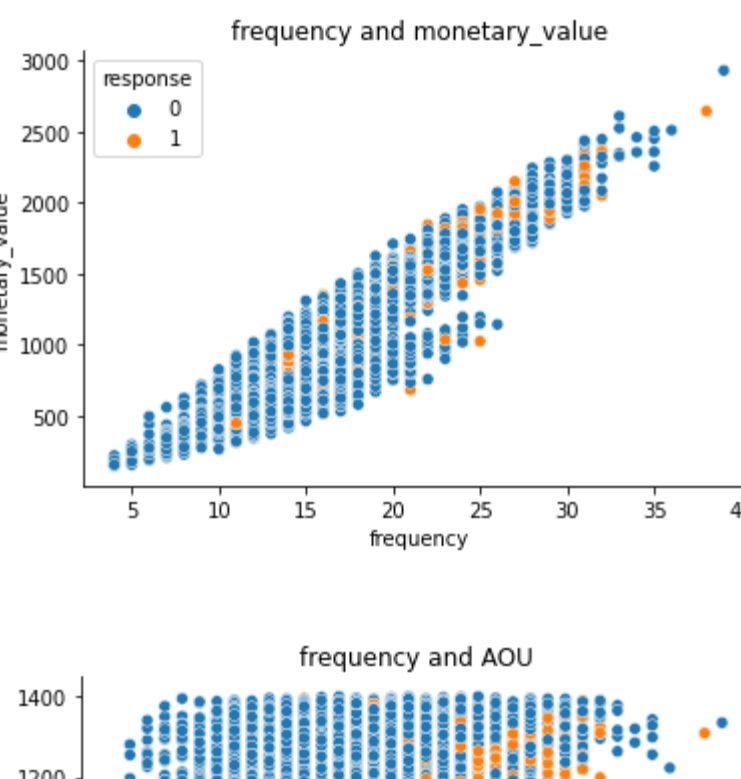
X_rfm = df_modeling_rfm.drop(columns=['response', 'customer_id'])
y_rfm = df_modeling_rfm['response']

X_cvv = df_modeling_cvv.drop(columns=['response', 'customer_id'])
y_cvv = df_modeling_cvv['response']

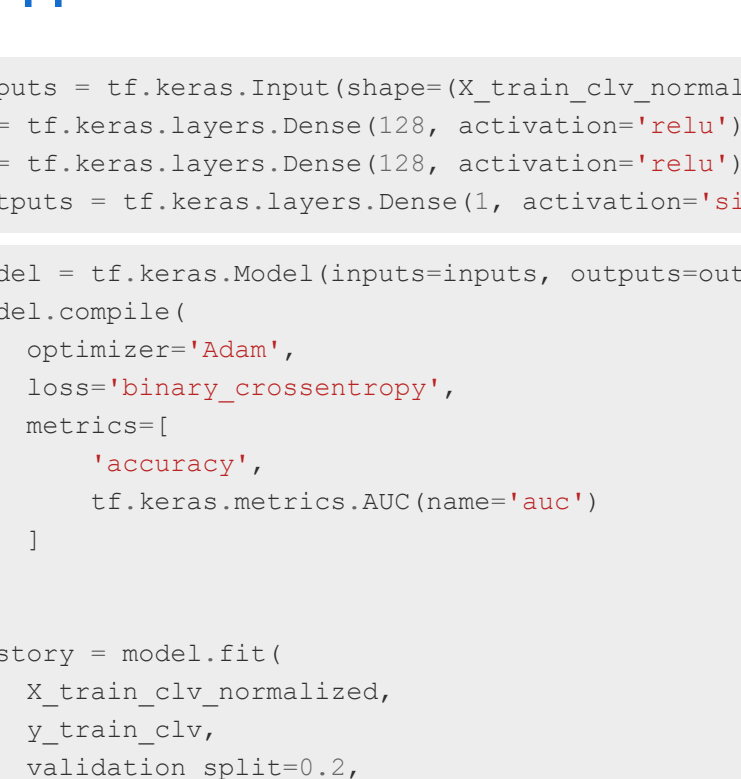
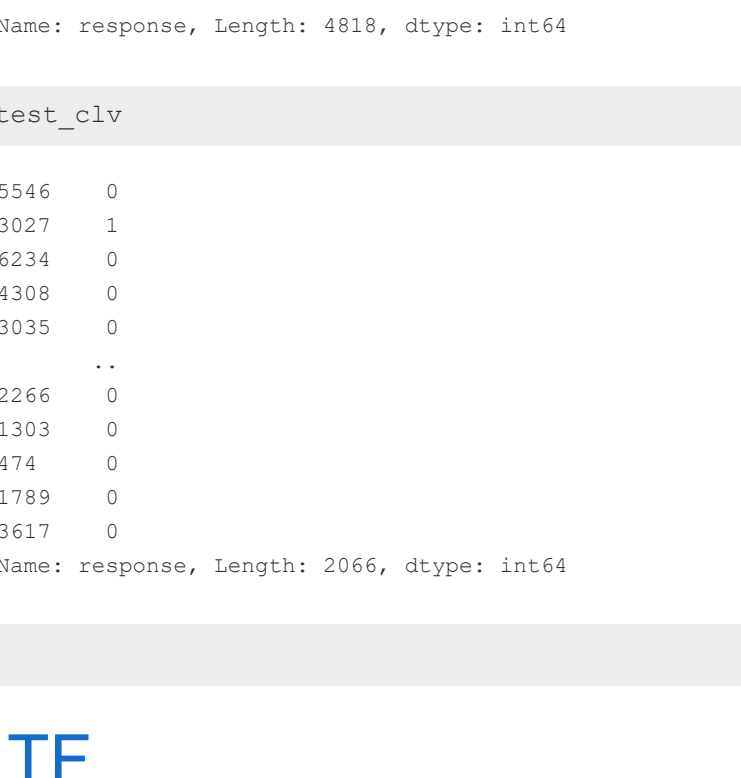
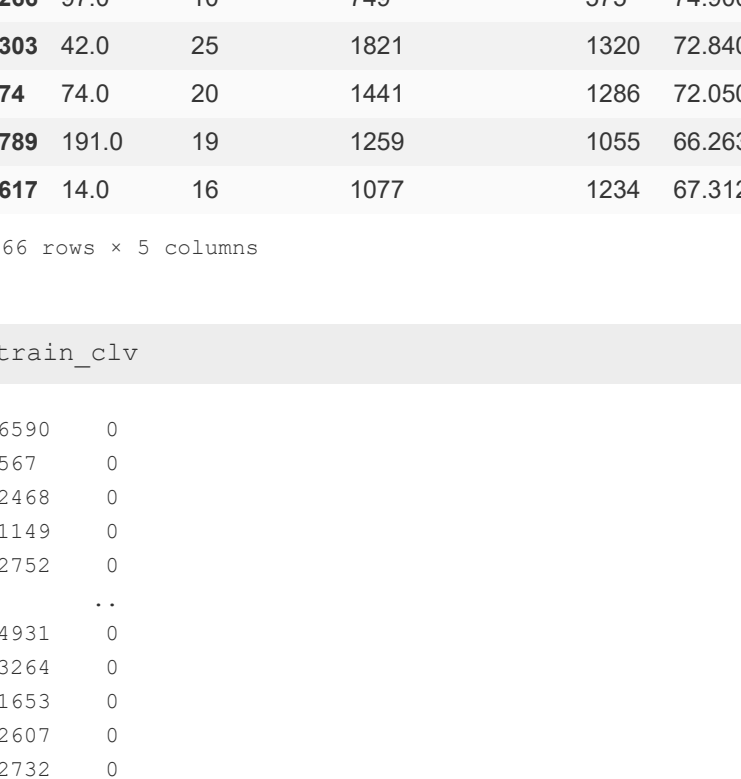
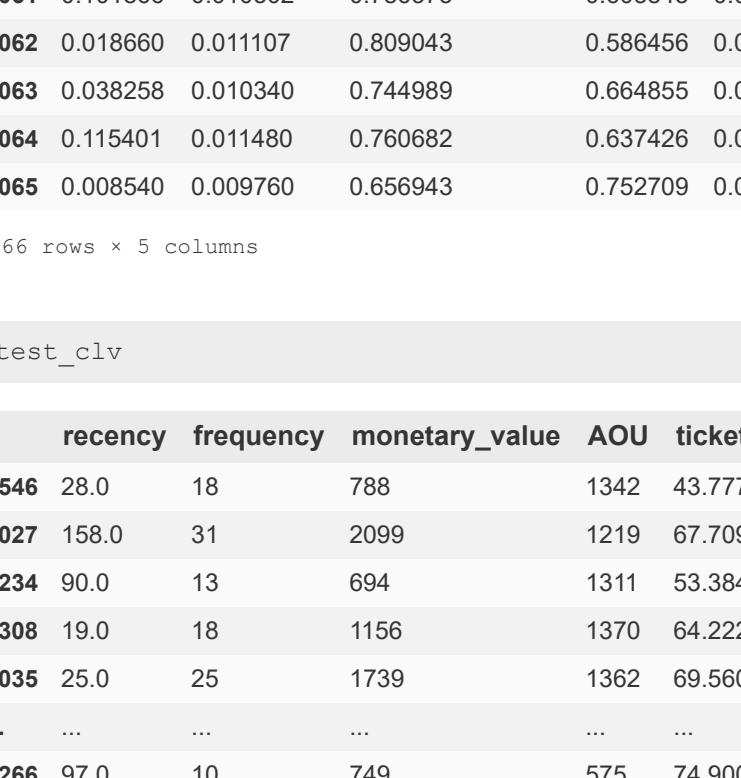
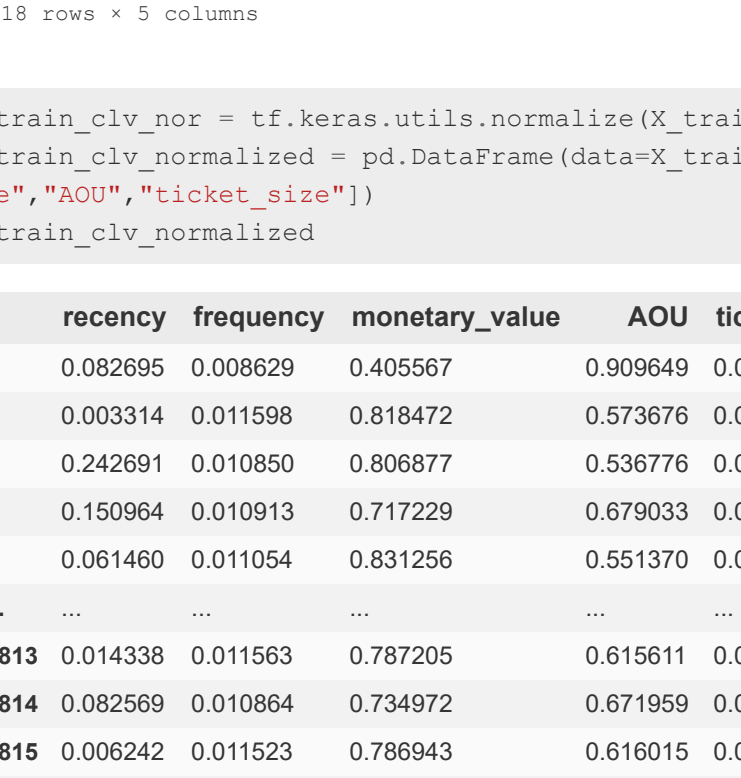
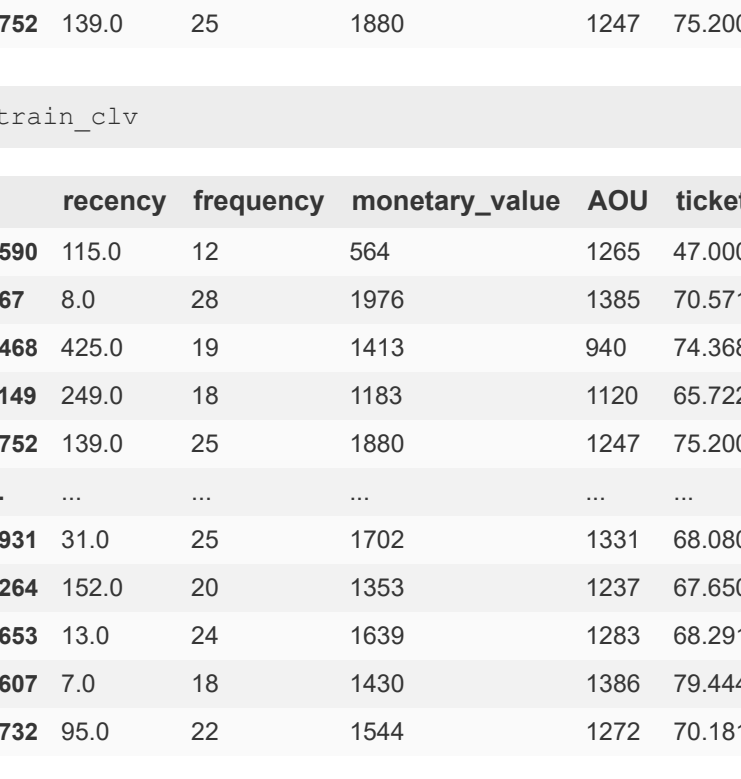
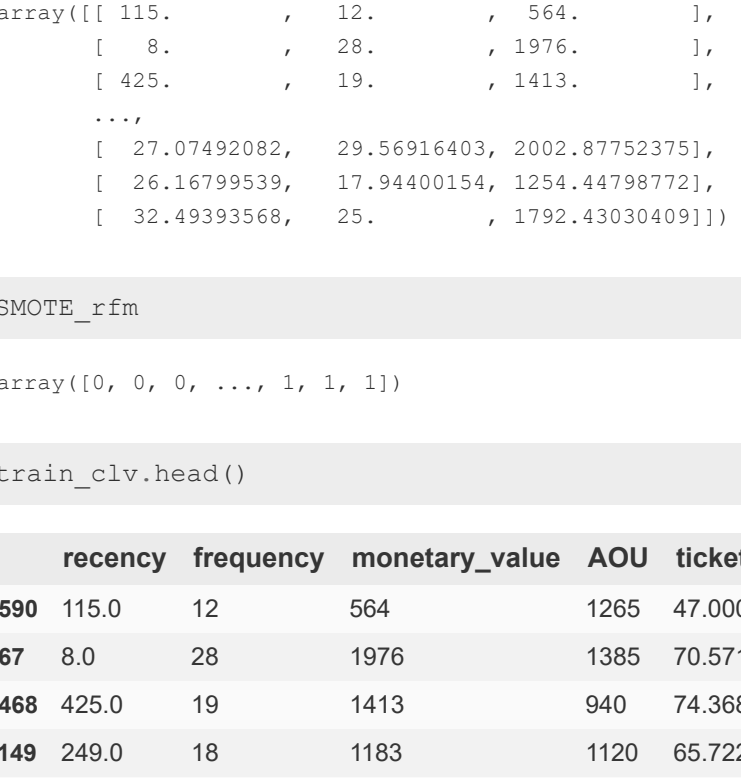
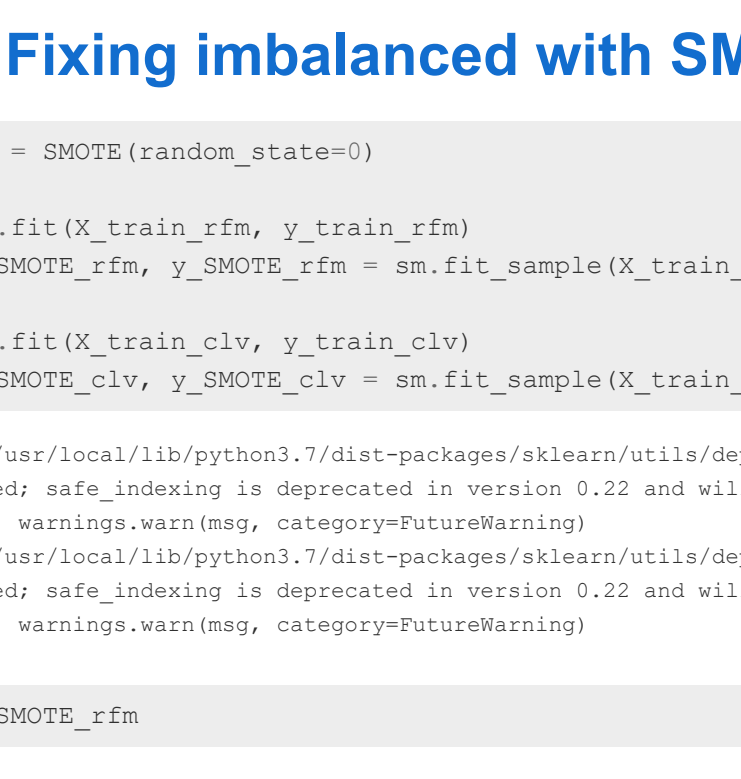
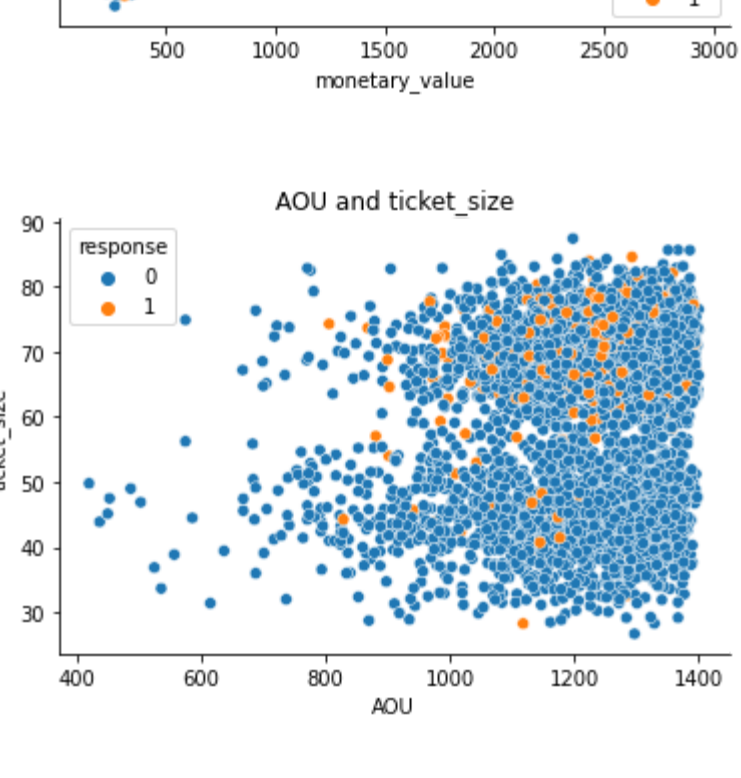
In [17]: # creating train and test dataset

X_train_rfm, X_test_rfm, y_train_rfm, y_test_rfm = train_test_split(X_rfm, y_rfm, test_size=0.3, random_state=0)
X_train_cvv, X_test_cvv, y_train_cvv, y_test_cvv = train_test_split(X_cvv, y_cvv, test_size=0.3, random_state=0)

In [18]: for i, col_i in enumerate(df_modeling_rfm[['recency', 'frequency', 'monetary_value']].columns):
for j, col_j in enumerate(df_modeling_cvv[['recency', 'frequency', 'monetary_value']].columns):
if i < j:
plt.title(col_i + ' and ' + col_j)
sns.scatterplot(data=df_modeling_rfm, x=col_i, y=col_j, hue='response')
sns.despine()
plt.show()
```



```
In [19]: for i, col_i in enumerate(df_modeling_cvv[['recency', 'frequency', 'monetary_value', 'AOU', 'ticket_size']].columns):
for j, col_j in enumerate(df_modeling_cvv[['recency', 'frequency', 'monetary_value', 'AOU', 'ticket_size']].columns):
if i < j:
plt.title(col_i + ' and ' + col_j)
sns.scatterplot(data=df_modeling_cvv, x=col_i, y=col_j, hue='response')
sns.despine()
plt.show()
```



Fixing imbalanced with SMOTE

```
In [138]: sm = SMOTE(random_state=0)

sm.fit(X_train_rfm, y_train_rfm)
X_SMOTE_rfm, y_SMOTE_rfm = sm.fit_sample(X_train_rfm, y_train_rfm)

sm.fit(X_train_cvv, y_train_cvv)
X_SMOTE_cvv, y_SMOTE_cvv = sm.fit_sample(X_train_cvv, y_train_cvv)

/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function safe_indexing is deprecated as of 0.22.0. Use indexing with integer arrays instead.
warnings.warn(msg, category=FutureWarning)

/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function safe_indexing is deprecated as of 0.22.0. Use indexing with integer arrays instead.
warnings.warn(msg, category=FutureWarning)
```

```
In [139]: X_SMOTE_rfm

array([[ 115.,    12.,    564.,    ],
       [  8.,    28.,    1976.,    ],
       [ 425.,    19.,    1413.,    ],
       ...,
       [ 270740082., 293691485., 2002.87753379],
       [ 2610799358., 1394400154., 1234.44798772],
       [ 32.49393568., 25.,    1792.43030409]])
```

```
In [140]: y_SMOTE_rfm

array([0, 0, ..., 1, 1, 1])
```

```
In [141]: X_train_cvv.head()
```

recency	frequency	monetary_value	AOU	ticket_size
6590	115.0	12	564	1265
567	8.0	28	1976	1385
2468	425.0	19	1413	940
1149	249.0	18	1183	1120
2752	139.0	25	1880	1247

```
In [142]: X_train_cvv

recency frequency monetary_value AOU ticket_size
0 0.082695 0.008629 0.405567 0.909649 0.037971
1 0.003314 0.011598 0.818472 0.573876 0.028113
2 0.060520 0.012738 0.862518 0.500910 0.027823
3 0.000528 0.008742 0.406679 0.881580 0.035988
4 0.010592 0.010034 0.644409 0.763703 0.034801
...
2061 0.101866 0.010502 0.786573 0.603845 0.078657
2062 0.018660 0.011070 0.809043 0.586456 0.032362
2063 0.038258 0.010340 0.744989 0.664855 0.037249
2064 0.115401 0.011480 0.760682 0.637426 0.040059
2065 0.008540 0.009780 0.656943 0.752709 0.041638
2066 rows x 5 columns
```

```
In [143]: X_test_cvv

recency frequency monetary_value AOU ticket_size
5546 28.0 18 788 1342 43.777778
3027 158.0 31 2099 1219 67.709677
6234 90.0 13 694 1311 53.84615
4308 19.0 18 1156 1370 64.222222
3035 25.0 25 1739 1362 69.560000
...
2266 97.0 10 749 575 74.900000
1303 42.0 25 1821 1320 72.840000
474 74.0 20 1441 1286 72.050000
1789 191.0 19 1259 1055 66.263158
3617 14.0 16 1077 1234 67.312500
2066 rows x 5 columns
```

```
In [144]: y_train_cvv

6590 0
567 0
2468 0
1149 0
2752 0
...
4931 0
2266 0
1303 0
474 0
1789 0
3617 0
Name: response, Length: 4818, dtype: int64
```

```
In [145]: y_test_cvv

5546 0
3027 1
6234 0
4308 0
3035 0
...
2266 0
1303 0
474 0
1789 0
3617 0
Name: response, Length: 2066, dtype: int64
```

TF

```
In [150]: inputs = tf.keras.Input(shape=(X_train_cvv.normalized.shape[1],))
x = tf.keras.layers.Dense(128, activation='relu')(inputs)
y = tf.keras.layers.Dense(128, activation='relu')(x)
outputs = tf.keras.layers.Dense(1, activation='sigmoid')(x)

In [151]: model = tf.keras.Model(inputs=inputs, outputs=outputs)
model.compile(
optimizer='adam',
loss='binary_crossentropy',
metrics=[
'accuracy',
tf.keras.metrics.AUC(name='auc')
])

history = model.fit(
X_train_cvv.normalized,
y_train_cvv,
validation_split=0.2,
batch_size=32,
epochs=100,
callbacks=[
tf.keras.callbacks.EarlyStopping(
monitor='val_loss',
patience=3,
restore_best_weights=True
)
])
```

```
Epoch 1/100
121/121 [=====] - 1s 5ms/step - loss: 0.2768 - accuracy: 0.9102 - auc: 0.7113 - val_loss: 0.2714
- val_accuracy: 0.9139 - val_auc: 0.7227
Epoch 2/100
121/121 [=====] - 0s 3ms/step - loss: 0.2769 - accuracy: 0.9102 - auc: 0.7106 - val_loss: 0.2707
- val_accuracy: 0.9139 - val_auc: 0.7246
Epoch 3/100
121/121 [=====] - 0s 2ms/step - loss: 0.2760 - accuracy: 0.9102 - auc: 0.7124 - val_loss: 0.2704
- val_accuracy: 0.9139 - val_auc: 0.7273
Epoch 4/100
121/121 [=====] - 0s 3ms/step - loss: 0.2760 - accuracy: 0.9102 - auc: 0.7129 - val_loss: 0.2699
- val_accuracy: 0.9139 - val_auc: 0.7279
Epoch 5/100
121/121 [=====] - 0s 3ms/step - loss: 0.2763 - accuracy: 0.9102 - auc: 0.7106 - val_loss: 0.2699
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 6/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 7/100
121/121 [=====] - 0s 2ms/step - loss: 0.2772 - accuracy: 0.9102 - auc: 0.7063 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7289
Epoch 8/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 9/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 10/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 11/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 12/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 13/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 14/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 15/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 16/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 17/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 18/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 19/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 20/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 21/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 22/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 23/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 24/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 25/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 26/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 27/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 28/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 29/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 30/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 31/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 32/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 33/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 34/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 35/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 36/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 37/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 38/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 39/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 40/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 41/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 42/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 43/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 44/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 45/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 46/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 47/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 48/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 49/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 50/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 51/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 52/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 53/100
121/121 [=====] - 0s 3ms/step - loss: 0.2776 - accuracy: 0.9102 - auc: 0.7040 - val_loss: 0.2702
- val_accuracy: 0.9139 - val_auc: 0.7292
Epoch 54/100

```



```
results = model.evaluate(X_train_cliv, y_train_cliv, verbose=0)
print("    Train Loss: {:.5f}".format(results[0]))
print("Train Accuracy: {:.2f}%".format(results[1] * 100))
print("    Train AUC: {:.5f}".format(results[2]))
```

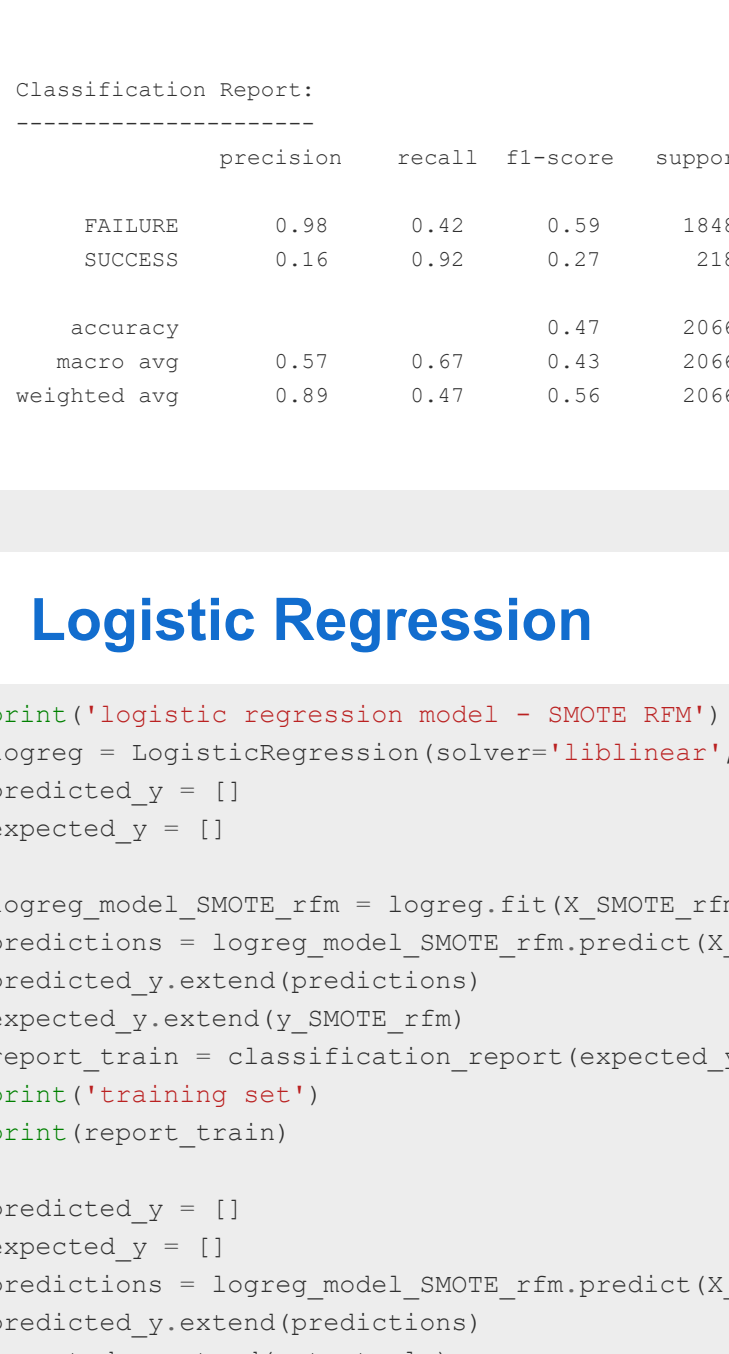
Train Loss: 0.28952
Train Accuracy: 86.454
Train AUC: 0.71426

```
In [142]: y_pred = np.array(model.predict(X_test_cliv) >= 0.05, dtype=np.int)

cm = confusion_matrix(y_test_cliv, y_pred)
clr = classification_report(y_test_cliv, y_pred, target_names=["FAILURE", "SUCCESS"])

plt.figure(figsize=(6, 6))
sns.heatmap(cm, annot=True, fmt='g', vmin=0, cmap='Blues', cbar=False)
plt.xticks(ticks=np.arange(2) + 0.5, labels=["FAILURE", "SUCCESS"])
plt.yticks(ticks=np.arange(2) + 0.5, labels=["FAILURE", "SUCCESS"])
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.title("Confusion Matrix")
plt.show()
```

print("Classification Report:\n-----\n", clr)



Logistic Regression

```
In [221]: print('logistic regression model - SMOTE RFM')
logreg = LogisticRegression(solver='liblinear', class_weight='balanced')
predicted_y = []
expected_y = []

logreg_model_SMOTE_rfm = logreg.fit(X_SMOTE_rfm, y_SMOTE_rfm)
predictions = logreg_model_SMOTE_rfm.predict(X_SMOTE_rfm)
predicted_y.extend(predictions)
expected_y.extend(y_test_rfm)
report_train = classification_report(expected_y, predicted_y)
print('training set')
print(report_train)

predicted_y = []
expected_y = []
predictions = logreg_model_SMOTE_rfm.predict(X_test_rfm)
predicted_y.extend(predictions)
expected_y.extend(y_test_cliv)
report_test = classification_report(expected_y, predicted_y)
print('test set')
print(report_test)
```

logistic regression model - SMOTE RFM

training set	precision	recall	f1-score	support
0	0.68	0.62	0.65	4389
1	0.65	0.71	0.68	4389
accuracy	0.67	0.67	0.67	8778
macro avg	0.67	0.67	0.67	8778
weighted avg	0.67	0.67	0.67	8778

test set	precision	recall	f1-score	support
0	0.95	0.63	0.75	1848
1	0.18	0.71	0.29	218
accuracy	0.57	0.67	0.64	2066
macro avg	0.57	0.67	0.52	2066
weighted avg	0.87	0.64	0.71	2066

```
In [222]: print('logistic regression model - SMOTE CLV')
logreg = LogisticRegression(solver='liblinear', class_weight='balanced')
predicted_y = []
expected_y = []

logreg_model_SMOTE_clv = logreg.fit(X_SMOTE_clv, y_SMOTE_clv)
predictions = logreg_model_SMOTE_clv.predict(X_SMOTE_clv)
predicted_y.extend(predictions)
expected_y.extend(y_test_cliv)
report_train = classification_report(expected_y, predicted_y)
print('training set')
print(report_train)

predicted_y = []
expected_y = []
predictions = logreg_model_SMOTE_clv.predict(X_test_cliv)
predicted_y.extend(predictions)
expected_y.extend(y_test_cliv)
report_test = classification_report(expected_y, predicted_y)
print('test set')
print(report_test)
```

logistic regression model - SMOTE CLV

training set	precision	recall	f1-score	support
0	0.68	0.62	0.65	4389
1	0.65	0.71	0.68	4389
accuracy	0.67	0.67	0.67	8778
macro avg	0.67	0.67	0.67	8778
weighted avg	0.67	0.67	0.67	8778

test set	precision	recall	f1-score	support
0	0.95	0.63	0.75	1848
1	0.18	0.72	0.29	218
accuracy	0.57	0.67	0.63	2066
macro avg	0.57	0.67	0.52	2066
weighted avg	0.87	0.64	0.71	2066

XGBoost

```
In [231]: print('XGBoost model - SMOTE RFM')

xgb_model = xgb.XGBClassifier(objective='binary:logistic', eval_metric='auc',
                              learning_rate=0.01,
                              n_estimators=100,
                              max_depth=2,
                              gamma=0.0,
                              colsample_bytree=0.6)

predicted_y = []
expected_y = []

xgb_model_SMOTE_rfm = xgb_model.fit(X_SMOTE_rfm, y_SMOTE_rfm, early_stopping_rounds=5, eval_set=[(X_test_rfm.to_numpy(), y_test_rfm)])
predictions = xgb_model_SMOTE_rfm.predict(X_SMOTE_rfm)
predicted_y.extend(predictions)
expected_y.extend(y_test_rfm)
report_train = classification_report(expected_y, predicted_y)
print('training set')
print(report_train)

predicted_y = []
expected_y = []
predictions = xgb_model_SMOTE_rfm.predict(X_test_rfm.to_numpy())
predicted_y.extend(predictions)
expected_y.extend(y_test_rfm)
report_test = classification_report(expected_y, predicted_y)
print('test set')
print(report_test)
```

XGBoost model - SMOTE RFM

[0] validation_0-auc:0.567615
Will train until validation_0-auc hasn't improved in 5 rounds.

[1]	validation_0-auc:0.713189
[2]	validation_0-auc:0.724408
[3]	validation_0-auc:0.69828
[4]	validation_0-auc:0.711449
[5]	validation_0-auc:0.713189
[6]	validation_0-auc:0.70826

Stopping. Best iteration:
[1] validation_0-auc:0.713189

training set	precision	recall	f1-score	support
0	0.75	0.57	0.65	4389
1	0.65	0.81	0.72	4389
accuracy	0.70	0.69	0.69	8778
macro avg	0.70	0.69	0.69	8778
weighted avg	0.70	0.69	0.69	8778

test set	precision	recall	f1-score	support
0	0.96	0.57	0.71	1848
1	0.18	0.79	0.29	218
accuracy	0.57	0.68	0.59	2066
macro avg	0.57	0.68	0.50	2066
weighted avg	0.88	0.59	0.67	2066

```
In [241]: print('XGBoost model - SMOTE CLV')

xgb_model = xgb.XGBClassifier(objective='binary:logistic', eval_metric='auc',
                              learning_rate=0.01,
                              n_estimators=100,
                              max_depth=2,
                              gamma=0.0,
                              colsample_bytree=0.6)

predicted_y = []
expected_y = []

xgb_model_SMOTE_clv = xgb_model.fit(X_SMOTE_clv, y_SMOTE_clv, early_stopping_rounds=5, eval_set=[(X_test_clv.to_numpy(), y_test_cliv)])
predictions = xgb_model_SMOTE_clv.predict(X_SMOTE_clv)
predicted_y.extend(predictions)
expected_y.extend(y_test_cliv)
report_train = classification_report(expected_y, predicted_y)
print('training set')
print(report_train)

predicted_y = []
expected_y = []
predictions = xgb_model_SMOTE_clv.predict(X_test_clv.to_numpy())
predicted_y.extend(predictions)
expected_y.extend(y_test_cliv)
report_test = classification_report(expected_y, predicted_y)
print('test set')
print(report_test)
```

XGBoost model - SMOTE CLV

[0] validation_0-auc:0.68482
Will train until validation_0-auc hasn't improved in 5 rounds.

[1]	validation_0-auc:0.717808
[2]	validation_0-auc:0.724408
[3]	validation_0-auc:0.725975
[4]	validation_0-auc:0.723803
[5]	validation_0-auc:0.726727
[6]	validation_0-auc:0.72539
[7]	validation_0-auc:0.724527
[8]	validation_0-auc:0.723888
[9]	validation_0-auc:0.724839
[10]	validation_0-auc:0.724921

Stopping. Best iteration:
[5] validation_0-auc:0.726727

training set	precision	recall	f1-score	support
0	0.73	0.55	0.64	4389
1	0.65	0.84	0.73	4389
accuracy	0.71	0.70	0.70	8778
macro avg	0.71	0.70	0.69	8778
weighted avg	0.71	0.70	0.69	8778

test set	precision	recall	f1-score	support
0	0.96	0.55	0.70	1848
1	0.17	0.80	0.28	218
accuracy	0.57	0.67	0.58	2066
macro avg	0.57	0.67	0.49	2066
weighted avg	0.88	0.58	0.66	2066

```
In [251]: ## building pipeline for hyperparameter tuning

from sklearn.pipeline import Pipeline
from sklearn.feature_selection import SelectKBest, chi2

# Create a pipeline
pipe = Pipeline([
    ('fs', SelectKBest()),
    ('clf', xgb.XGBClassifier(objective='binary:logistic', scale_pos_weight=9))
])

In [261]: ## hyper parameter tuning - grid search

from sklearn.model_selection import KFold, GridSearchCV
from sklearn.metrics import accuracy_score, make_scorer
# Define our search space for grid search
search_space = [
    (
        'clf_n_estimators': [100, 300],
        'clf_learning_rate': [0.01, 0.1],
        'clf_max_depth': range(2, 5),
        'clf_colsample_bytree': [1/10.0 for i in range(4, 7)],
        'fs_score_func': [chi2],
        'fs_k': [2],
    )
]
# Define cross validation
kfold = KFold(n_splits=5, random_state=42)
# AUC and F1 as score
scoring = ('AUC': 'roc_auc', 'F1 score': 'f1_micro')
# Define grid search
grid = GridSearchCV(
    pipe,
    param_grid=search_space,
    cv=kfold,
    scoring=scoring,
    refit='AUC',
    verbose=1,
    n_jobs=-1
)
# Fit grid search
xgb_model_clv_gs = grid.fit(X_train_cliv, y_train_cliv)
```

Fitting 5 folds for each of 108 candidates, totalling 540 fits

/usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_split.py:236: FutureWarning: Setting a random_state has no effect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None).
r.set_shuffle=True.
FutureWarning
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 2 concurrent workers.
[Parallel(n_jobs=-1)]: Done 46 tasks | elapsed: 9.8s
[Parallel(n_jobs=-1)]: Done 196 tasks | elapsed: 37.9s
[Parallel(n_jobs=-1)]: Done 446 tasks | elapsed: 1.4min
[Parallel(n_jobs=-1)]: Done 540 out of 540 | elapsed: 1.7min finished

```
In [271]: predicted_y = []
expected_y = []
predictions = xgb_model_clv_gs.predict(X_test_cliv)
print('Best AUC Score: {}'.format(xgb_model_clv_gs.best_score_))
print(confusion_matrix(y_test_cliv, predictions))

predicted_y.extend(predictions)
expected_y.extend(y_test_cliv)
report_test = classification_report(expected_y, predicted_y)
print('test set')
print(report_test)
```

Best AUC Score: 0.7089054673735633

[1089 759]
[51 167]

test set	precision	recall	f1-score	support
0	0.96	0.59	0.73	1848
1	0.18	0.77	0.29	218
accuracy	0.57	0.68	0.61	2066
macro avg	0.57	0.68	0.51	2066
weighted avg	0.87	0.61	0.68	2066

```
In [281]: print(xgb_model_clv_gs.best_params_)

{'clf_colsample_bytree': 0.4, 'clf_gamma': 0.0, 'clf_learning_rate': 0.01, 'clf_max_depth': 2, 'clf_n_estimators': 100, 'fs_k': 2, 'fs_score_func': <function chi2 at 0x7fddc96a70>}
```

```
In [281]:
```