

Project Title: Home Credit Default Risk (HCDR)



Deepak Kasi
Nathan
(dekasi@iu.edu)



Sai Sumanth Muvva
(saimuvva@iu.edu)



Viswa Suhaas Penugonda
(vpenugon@iu.edu)



Teja Naidu Chintha
(tnchinth@iu.edu)

The course project is based on the [Home Credit Default Risk \(HCDR\) Kaggle Competition](#). The goal of this project is to predict whether or not a client will repay a loan. In order to make sure that people who struggle to get loans due to insufficient or non-existent credit histories have a positive loan experience, Home Credit makes use of a variety of alternative data--including telco and transactional information--to predict their clients' repayment abilities.

Some of the challenges

1. Dataset size
 - (688 meg compressed) with millions of rows of data
 - 2.71 Gig of data uncompressed
- Dealing with missing data
- Imbalanced datasets
- Summarizing transaction data

Dataset and how to download

Back ground Home Credit Group

Many people struggle to get loans due to insufficient or non-existent credit histories. And, unfortunately, this population is often taken advantage of by untrustworthy lenders.

Home Credit Group

Home Credit strives to broaden financial inclusion for the unbanked population by providing a positive and safe borrowing experience. In order to make sure this underserved population has a positive loan experience, Home Credit makes use of a variety of alternative data--including telco and transactional information--to predict their clients' repayment abilities.

While Home Credit is currently using various statistical and machine learning methods to make these predictions, they're challenging Kagglers to help them unlock the full potential of their data. Doing so will ensure that clients capable of repayment are not rejected and that loans are given with a principal, maturity, and repayment calendar that will empower their clients to be successful.

Background on the dataset

Home Credit is a non-banking financial institution, founded in 1997 in the Czech Republic.

The company operates in 14 countries (including United States, Russia, Kazakhstan, Belarus, China, India) and focuses on lending primarily to people with little or no credit history which will either not obtain loans or became victims of untrustworthy lenders.

Home Credit group has over 29 million customers, total assets of 21 billions Euro, over 160 millions loans, with the majority in Asia and almost half of them in China (as of 19-05-2018).

While Home Credit is currently using various statistical and machine learning methods to make these predictions, they're challenging Kagglers to help them unlock the full potential of their data. Doing so will ensure that clients capable of repayment are not rejected and that loans are given with a principal, maturity, and repayment calendar that will empower their clients to be successful.

Data files overview

The `HomeCredit_columns_description.csv` acts as a data dictionary.

There are 7 different sources of data:

- **application_train/application_test (307k rows, and 48k rows):** the main training and testing data with information about each loan application at Home Credit. Every loan has its own row and is identified by the feature `SK_ID_CURR`. The training application data comes with the TARGET indicating **0: the loan was repaid or 1: the**

loan was not repaid. The target variable defines if the client had payment difficulties meaning he/she had late payment more than X days on at least one of the first Y installments of the loan. Such case is marked as 1 while other all other cases as 0.

- **bureau (1.7 Million rows):** data concerning client's previous credits from other financial institutions. Each previous credit has its own row in bureau, but one loan in the application data can have multiple previous credits.
- **bureau_balance (27 Million rows):** monthly data about the previous credits in bureau. Each row is one month of a previous credit, and a single previous credit can have multiple rows, one for each month of the credit length.
- **previous_application (1.6 Million rows):** previous applications for loans at Home Credit of clients who have loans in the application data. Each current loan in the application data can have multiple previous loans. Each previous application has one row and is identified by the feature SK_ID_PREV.
- **POS_CASH_BALANCE (10 Million rows):** monthly data about previous point of sale or cash loans clients have had with Home Credit. Each row is one month of a previous point of sale or cash loan, and a single previous loan can have many rows.
- **credit_card_balance:** monthly data about previous credit cards clients have had with Home Credit. Each row is one month of a credit card balance, and a single credit card can have many rows.
- **installments_payment (13.6 Million rows):** payment history for previous loans at Home Credit. There is one row for every made payment and one row for every missed payment.

Table sizes

name	[rows cols]	MegaBytes
application_train	: [307,511, 122]:	158MB
application_test	: [48,744, 121]:	25MB
bureau	: [1,716,428, 17]	162MB
bureau_balance	: [27,299,925, 3]:	358MB
credit_card_balance	: [3,840,312, 23]	405MB
installments_payments	: [13,605,401, 8]	690MB
previous_application	: [1,670,214, 37]	386MB
POS_CASH_balance	: [10,001,358, 8]	375MB

Imports

In [1]:

```
from scipy import stats
# import latexify
import time
import numpy as np
import pandas as pd
from sklearn.preprocessing import LabelEncoder
import os
import zipfile
import pickle
import json
from sklearn.base import BaseEstimator, TransformerMixin
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import GridSearchCV
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import MinMaxScaler
from sklearn.pipeline import Pipeline, FeatureUnion
from pandas.plotting import scatter_matrix
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import OneHotEncoder
from sklearn.pipeline import make_pipeline
from sklearn.model_selection import ShuffleSplit
from sklearn.metrics import roc_auc_score
from sklearn.metrics import accuracy_score, confusion_matrix, f1_score
from sklearn.metrics import log_loss, classification_report, roc_auc_score,
from scipy import stats
from sklearn.svm import SVC
from xgboost import XGBClassifier
import warnings
warnings.filterwarnings('ignore')
```

Data files overview

Data Dictionary

As part of the data download comes a Data Dictionary. It named
`HomeCredit_columns_description.csv`

Application train

In [2]:

```

def load_data(in_path, name):
    df = pd.read_csv(in_path)
    print(f'{name}: shape is {df.shape}')
    print(df.info())
    display(df.head(5))
    return df
datasets={}
ds_name = 'application_train'
DATA_DIR=f"/Users/deepak/Desktop/AML/home-credit-default-risk/"
datasets[ds_name] = load_data(os.path.join(DATA_DIR, f'{ds_name}.csv'), ds_name)
datasets['application_train'].shape

```

application_train: shape is (307511, 122)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 307511 entries, 0 to 307510
Columns: 122 entries, SK_ID_CURR to AMT_REQ_CREDIT_BUREAU_YEAR
dtypes: float64(65), int64(41), object(16)
memory usage: 286.2+ MB
None

	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_
0	100002	1	Cash loans	M	N	
1	100003	0	Cash loans	F	N	
2	100004	0	Revolving loans	M	Y	
3	100006	0	Cash loans	F	N	
4	100007	0	Cash loans	M	N	

5 rows × 122 columns

Out[2]: (307511, 122)

Application test

- **application_train/application_test:** the main training and testing data with information about each loan application at Home Credit. Every loan has its own row and is identified by the feature SK_ID_CURR. The training application data comes with the TARGET indicating **0: the loan was repaid** or **1: the loan was not repaid**. The target variable defines if the client had payment difficulties meaning he/she had late payment more than X days on at least one of the first Y installments of the loan. Such case is marked as 1 while other all other cases as 0.

In [3]:

```
ds_name = 'application_test'  
datasets[ds_name] = load_data(os.path.join(DATA_DIR, f'{ds_name}.csv'), ds_
```

```
application_test: shape is (48744, 121)  
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 48744 entries, 0 to 48743  
Columns: 121 entries, SK_ID_CURR to AMT_REQ_CREDIT_BUREAU_YEAR  
dtypes: float64(65), int64(40), object(16)  
memory usage: 45.0+ MB  
None
```

	SK_ID_CURR	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REA
0	100001	Cash loans	F	N	
1	100005	Cash loans	M	N	
2	100013	Cash loans	M	Y	
3	100028	Cash loans	F	N	
4	100038	Cash loans	M	Y	

5 rows × 121 columns

The application dataset has the most information about the client: Gender, income, family status, education ...

The Other datasets

- **bureau:** data concerning client's previous credits from other financial institutions. Each previous credit has its own row in bureau, but one loan in the application data can have multiple previous credits.
- **bureau_balance:** monthly data about the previous credits in bureau. Each row is one month of a previous credit, and a single previous credit can have multiple rows, one for each month of the credit length.
- **previous_application:** previous applications for loans at Home Credit of clients who have loans in the application data. Each current loan in the application data can have multiple previous loans. Each previous application has one row and is identified by the feature SK_ID_PREV.
- **POS_CASH_BALANCE:** monthly data about previous point of sale or cash loans clients have had with Home Credit. Each row is one month of a previous point of sale or cash loan, and a single previous loan can have many rows.
- **credit_card_balance:** monthly data about previous credit cards clients have had with Home Credit. Each row is one month of a credit card balance, and a single credit card can have many rows.
- **installments_payment:** payment history for previous loans at Home Credit. There is one row for every made payment and one row for every missed payment.

In [4]:

```
%%time
ds_names = ("application_train", "application_test", "bureau", "bureau_balance",
             "previous_application", "POS_CASH_balance")

for ds_name in ds_names:
    datasets[ds_name] = load_data(os.path.join(DATA_DIR, f'{ds_name}.csv'),
```



```
application_train: shape is (307511, 122)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 307511 entries, 0 to 307510
Columns: 122 entries, SK_ID_CURR to AMT_REQ_CREDIT_BUREAU_YEAR
dtypes: float64(65), int64(41), object(16)
memory usage: 286.2+ MB
None
```

SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_!
0	100002	1	Cash loans	M	N
1	100003	0	Cash loans	F	N
2	100004	0	Revolving loans	M	Y
3	100006	0	Cash loans	F	N
4	100007	0	Cash loans	M	N

5 rows × 122 columns

```
application_test: shape is (48744, 121)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 48744 entries, 0 to 48743
Columns: 121 entries, SK_ID_CURR to AMT_REQ_CREDIT_BUREAU_YEAR
dtypes: float64(65), int64(40), object(16)
memory usage: 45.0+ MB
None
```

SK_ID_CURR	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REA
0	100001	Cash loans	F	N
1	100005	Cash loans	M	N
2	100013	Cash loans	M	Y
3	100028	Cash loans	F	N
4	100038	Cash loans	M	Y

5 rows × 121 columns

```
bureau: shape is (1716428, 17)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1716428 entries, 0 to 1716427
Data columns (total 17 columns):
 #   Column           Dtype  
 --- 
 0   SK_ID_CURR       int64  
 1   SK_ID_BUREAU     int64  
 2   CREDIT_ACTIVE     object  
 3   CREDIT_CURRENCY   object  
 4   DAYS_CREDIT       int64  
 5   CREDIT_DAY_OVERDUE int64  
 6   DAYS_CREDIT_ENDDATE float64 
 7   DAYS_ENDDATE_FACT float64 
 8   AMT_CREDIT_MAX_OVERDUE float64 
 9   CNT_CREDIT_PROLONG int64  
 10  AMT_CREDIT_SUM     float64 
 11  AMT_CREDIT_SUM_DEBT float64 
 12  AMT_CREDIT_SUM_LIMIT float64 
 13  AMT_CREDIT_SUM_OVERDUE float64 
 14  CREDIT_TYPE       object  
 15  DAYS_CREDIT_UPDATE int64  
 16  AMT_ANNUITY       float64 
dtypes: float64(8), int64(6), object(3)
memory usage: 222.6+ MB
None
```

	SK_ID_CURR	SK_ID_BUREAU	CREDIT_ACTIVE	CREDIT_CURRENCY	DAYS_CREDIT	CREDI
0	215354	5714462	Closed	currency 1	-497	
1	215354	5714463	Active	currency 1	-208	
2	215354	5714464	Active	currency 1	-203	
3	215354	5714465	Active	currency 1	-203	
4	215354	5714466	Active	currency 1	-629	

```
bureau_balance: shape is (27299925, 3)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 27299925 entries, 0 to 27299924
Data columns (total 3 columns):
 #   Column           Dtype  
 --- 
 0   SK_ID_BUREAU     int64  
 1   MONTHS_BALANCE   int64  
 2   STATUS            object  
dtypes: int64(2), object(1)
memory usage: 624.8+ MB
None
```

	SK_ID_BUREAU	MONTHS_BALANCE	STATUS
0	5715448	0	C
1	5715448	-1	C
2	5715448	-2	C
3	5715448	-3	C
4	5715448	-4	C

```
credit_card_balance: shape is (3840312, 23)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3840312 entries, 0 to 3840311
Data columns (total 23 columns):
 #   Column           Dtype  
 --- 
 0   SK_ID_PREV        int64  
 1   SK_ID_CURR        int64  
 2   MONTHS_BALANCE    int64  
 3   AMT_BALANCE       float64 
 4   AMT_CREDIT_LIMIT_ACTUAL int64  
 5   AMT_DRAWINGS_ATM_CURRENT float64 
 6   AMT_DRAWINGS_CURRENT  float64 
 7   AMT_DRAWINGS_OTHER_CURRENT float64 
 8   AMT_DRAWINGS_POS_CURRENT float64 
 9   AMT_INST_MIN_REGULARITY float64 
 10  AMT_PAYMENT_CURRENT  float64 
 11  AMT_PAYMENT_TOTAL_CURRENT float64 
 12  AMT_RECEIVABLE_PRINCIPAL float64 
 13  AMT_RECVABLE        float64 
 14  AMT_TOTAL_RECEIVABLE float64 
 15  CNT_DRAWINGS_ATM_CURRENT float64 
 16  CNT_DRAWINGS_CURRENT  int64  
 17  CNT_DRAWINGS_OTHER_CURRENT float64 
 18  CNT_DRAWINGS_POS_CURRENT float64 
 19  CNT_INSTALMENT_MATURE_CUM float64 
 20  NAME_CONTRACT_STATUS  object  
 21  SK_DPD             int64  
 22  SK_DPD_DEF         int64  
dtypes: float64(15), int64(7), object(1)
memory usage: 673.9+ MB
None
```

	SK_ID_PREV	SK_ID_CURR	MONTHS_BALANCE	AMT_BALANCE	AMT_CREDIT_LIMIT_ACTL
0	2562384	378907	-6	56.970	1351
1	2582071	363914	-1	63975.555	451
2	1740877	371185	-7	31815.225	4501
3	1389973	337855	-4	236572.110	2251
4	1891521	126868	-1	453919.455	4501

5 rows × 23 columns

```
installments_payments: shape is (13605401, 8)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 13605401 entries, 0 to 13605400
Data columns (total 8 columns):
 #   Column           Dtype  
--- 
 0   SK_ID_PREV      int64  
 1   SK_ID_CURR      int64  
 2   NUM_INSTALMENT_VERSION  float64
 3   NUM_INSTALMENT_NUMBER   int64  
 4   DAYS_INSTALMENT    float64
 5   DAYS_ENTRY_PAYMENT float64
 6   AMT_INSTALMENT     float64
 7   AMT_PAYMENT        float64
dtypes: float64(5), int64(3)
memory usage: 830.4 MB
None
```

	SK_ID_PREV	SK_ID_CURR	NUM_INSTALMENT_VERSION	NUM_INSTALMENT_NUMBER	DA
0	1054186	161674		1.0	6
1	1330831	151639		0.0	34
2	2085231	193053		2.0	1
3	2452527	199697		1.0	3
4	2714724	167756		1.0	2

```

previous_application: shape is (1670214, 37)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1670214 entries, 0 to 1670213
Data columns (total 37 columns):
 #   Column           Non-Null Count   Dtype  
--- 
 0   SK_ID_PREV       1670214 non-null   int64  
 1   SK_ID_CURR       1670214 non-null   int64  
 2   NAME_CONTRACT_TYPE 1670214 non-null   object  
 3   AMT_ANNUITY      1297979 non-null   float64 
 4   AMT_APPLICATION  1670214 non-null   float64 
 5   AMT_CREDIT        1670213 non-null   float64 
 6   AMT_DOWN_PAYMENT  774370 non-null   float64 
 7   AMT_GOODS_PRICE   1284699 non-null   float64 
 8   WEEKDAY_APPR_PROCESS_START 1670214 non-null   object  
 9   HOUR_APPR_PROCESS_START 1670214 non-null   int64  
 10  FLAG_LAST_APPL_PER_CONTRACT 1670214 non-null   object  
 11  NFLAG_LAST_APPL_IN_DAY    1670214 non-null   int64  
 12  RATE_DOWN_PAYMENT     774370 non-null   float64 
 13  RATE_INTEREST_PRIMARY 5951 non-null   float64 
 14  RATE_INTEREST_PRIVILEGED 5951 non-null   float64 
 15  NAME_CASH_LOAN_PURPOSE 1670214 non-null   object  
 16  NAME_CONTRACT_STATUS  1670214 non-null   object  
 17  DAYS_DECISION       1670214 non-null   int64  
 18  NAME_PAYMENT_TYPE   1670214 non-null   object  
 19  CODE_REJECT_REASON  1670214 non-null   object  
 20  NAME_TYPE_SUITE     849809 non-null   object  
 21  NAME_CLIENT_TYPE   1670214 non-null   object  
 22  NAME_GOODS_CATEGORY 1670214 non-null   object  
 23  NAME_PORTFOLIO      1670214 non-null   object  
 24  NAME_PRODUCT_TYPE   1670214 non-null   object  
 25  CHANNEL_TYPE        1670214 non-null   object  
 26  SELLERPLACE_AREA    1670214 non-null   int64  
 27  NAME_SELLER_INDUSTRY 1670214 non-null   object  
 28  CNT_PAYMENT         1297984 non-null   float64 
 29  NAME_YIELD_GROUP   1670214 non-null   object  
 30  PRODUCT_COMBINATION 1669868 non-null   object  
 31  DAYS_FIRST_DRAWING 997149 non-null   float64 
 32  DAYS_FIRST_DUE     997149 non-null   float64 
 33  DAYS_LAST_DUE_1ST_VERSION 997149 non-null   float64 
 34  DAYS_LAST_DUE      997149 non-null   float64 
 35  DAYS_TERMINATION   997149 non-null   float64 
 36  NFLAG_INSURED_ON_APPROVAL 997149 non-null   float64 

dtypes: float64(15), int64(6), object(16)
memory usage: 471.5+ MB
None

```

	SK_ID_PREV	SK_ID_CURR	NAME_CONTRACT_TYPE	AMT_ANNUITY	AMT_APPLICATION
0	2030495	271877	Consumer loans	1730.430	17145.0
1	2802425	108129	Cash loans	25188.615	607500.0
2	2523466	122040	Cash loans	15060.735	112500.0
3	2819243	176158	Cash loans	47041.335	450000.0
4	1784265	202054	Cash loans	31924.395	337500.0

5 rows × 37 columns

```
POS_CASH_balance: shape is (10001358, 8)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10001358 entries, 0 to 10001357
Data columns (total 8 columns):
 #   Column           Dtype  
--- 
 0   SK_ID_PREV      int64  
 1   SK_ID_CURR      int64  
 2   MONTHS_BALANCE int64  
 3   CNT_INSTALMENT float64 
 4   CNT_INSTALMENT_FUTURE float64
 5   NAME_CONTRACT_STATUS object 
 6   SK_DPD          int64  
 7   SK_DPD_DEF      int64  
dtypes: float64(2), int64(5), object(1)
memory usage: 610.4+ MB
None
```

	SK_ID_PREV	SK_ID_CURR	MONTHS_BALANCE	CNT_INSTALMENT	CNT_INSTALMENT_FU
0	1803195	182943		-31	48.0
1	1715348	367990		-33	36.0
2	1784872	397406		-32	12.0
3	1903291	269225		-35	48.0
4	2341044	334279		-35	36.0

```
CPU times: user 15.3 s, sys: 2.26 s, total: 17.6 s
Wall time: 17.8 s
```

In [5]:

```
for ds_name in datasets.keys():
    print(f'dataset {ds_name}: {24}: [{datasets[ds_name].shape[0]:10}, {data
```

```
dataset application_train      : [    307,511, 122]
dataset application_test       : [     48,744, 121]
dataset bureau                  : [ 1,716,428, 17]
dataset bureau_balance          : [ 27,299,925, 3]
dataset credit_card_balance     : [ 3,840,312, 23]
dataset installments_payments   : [ 13,605,401, 8]
dataset previous_application    : [ 1,670,214, 37]
dataset POS_CASH_balance         : [ 10,001,358, 8]
```

Undersampling

Undersampling is performed when the data is highly biased. Here we can see that in the target variable, the defaulters are very less as compared to other people who have successfully repaid the loan. Hence we perform undersampling and take random entries from the good population and keep all the entries from the defaulters.

Undersampling data randomly

```
In [6]: # application_train = datasets['application_train']
# # und_application_train = application_train[application_train['TARGET']==1]
# und_application_train = und_application_train.append(application_train[ap

In [7]: # Access the 'application_train' dataset from the 'datasets' container
application_train = datasets['application_train']

# Select the minority class instances (TARGET = 1) from the training database
minority_application_train = application_train[application_train['TARGET']==1]

# Append a randomly sampled subset of majority class instances (TARGET = 0)
undersampled_application_train = minority_application_train.append(
    application_train[application_train['TARGET']==0].reset_index(drop=True)
)

In [8]: # Assign the undersampled training dataset to a new key in the 'datasets' dictionary
datasets["undersampled_application_train"] = undersampled_application_train

# Count the number of instances in each class
class_distribution = undersampled_application_train['TARGET'].value_counts()

# Print the class distribution
print("Class distribution in the undersampled training dataset:")
print(class_distribution)

Class distribution in the undersampled training dataset:
0    75000
1    24825
Name: TARGET, dtype: int64
```

In [9]:

```
# datasets["und_application_train"] = und_application_train
# und_application_train['TARGET'].value_counts()
```

Undersampling by keeping similar ratio of non-defaulters to defaulters and also maintaining loan types of non-defaulters uniformly

In [10]:

```
# Assuming this is a dictionary where you store your datasets

# Filtering rows with TARGET == 1 and creating a new DataFrame
datasets["undersampled_application_train_2"] = datasets["application_train"]
datasets["undersampled_application_train_2"]["weight"] = 1

# Undersampling Cash loans
num_default_cashloans = len(datasets["undersampled_application_train_2"][(datasets["application_train"]['TARGET'] == 1) & (df_sample_cash['loan_type'] == 'cash')]) * 1
df_sample_cash = datasets["application_train"][(datasets["application_train"]['TARGET'] == 1) & (df_sample_cash['loan_type'] == 'cash')]
df_sample_cash['weight'] = 1

# Undersampling Revolving loans
num_default_revolvingloans = len(datasets["undersampled_application_train_2"][(datasets["application_train"]['TARGET'] == 1) & (df_sample_revolving['loan_type'] == 'revolving')]) * 1
df_sample_revolving = datasets["application_train"][(datasets["application_train"]['TARGET'] == 1) & (df_sample_revolving['loan_type'] == 'revolving')]
df_sample_revolving['weight'] = 1

# Combining undersampled cash loans and revolving loans with the initial Dataframe
datasets["undersampled_application_train_2"] = pd.concat([datasets["undersampled_application_train_2"], df_sample_cash, df_sample_revolving])

# Check the distribution of the TARGET variable
print(datasets["undersampled_application_train_2"].TARGET.value_counts())
```

```
1    24825
0    24825
Name: TARGET, dtype: int64
```

In [11]:

```
# Assuming this is a dictionary where you store your datasets

# Filtering rows with TARGET == 1 and creating a new DataFrame
undersampled_application_train_2 = datasets["application_train"][datasets["TARGET"] == 1]
undersampled_application_train_2['weight'] = 1

# Undersampling Cash loans
num_default_cashloans = len(undersampled_application_train_2[(undersampled_application_train_2['TARGET'] == 1) & (datasets["application_train"]["NAME_CONTRACT_STATUS"] == "Cash")])
df_sample_cash = datasets["application_train"][(datasets["application_train"]["NAME_CONTRACT_STATUS"] == "Cash")]
df_sample_cash['weight'] = 1

# Undersampling Revolving loans
num_default_revolvingloans = len(undersampled_application_train_2[(undersampled_application_train_2['TARGET'] == 1) & (datasets["application_train"]["NAME_CONTRACT_STATUS"] == "Revolving")])
df_sample_revolving = datasets["application_train"][(datasets["application_train"]["NAME_CONTRACT_STATUS"] == "Revolving")]
df_sample_revolving['weight'] = 1

# Combining undersampled cash loans and revolving loans with the initial Dataframe
undersampled_application_train_2 = pd.concat([undersampled_application_train_2, df_sample_cash, df_sample_revolving])

# Check the distribution of the TARGET variable
print(undersampled_application_train_2.TARGET.value_counts())
```

```
1    24825
0    24825
Name: TARGET, dtype: int64
```

Correlation Analysis

The correlation coefficient is not the best way to identify the relevancs between the features. But it gives an idea about possible relationship between the data.

Correlation with the target column

In [12]:

```
correlations = datasets["application_train"].corr()['TARGET'].sort_values()
print('Most Positive Correlations:\n', correlations.tail(10))
print('\nMost Negative Correlations:\n', correlations.head(10))
```

Most Positive Correlations:

FLAG_DOCUMENT_3	0.044346
REG_CITY_NOT_LIVE_CITY	0.044395
FLAG_EMP_PHONE	0.045982
REG_CITY_NOT_WORK_CITY	0.050994
DAYS_ID_PUBLISH	0.051457
DAYS_LAST_PHONE_CHANGE	0.055218
REGION_RATING_CLIENT	0.058899
REGION_RATING_CLIENT_W_CITY	0.060893
DAYS_BIRTH	0.078239
TARGET	1.000000

Name: TARGET, dtype: float64

Most Negative Correlations:

EXT_SOURCE_3	-0.178919
EXT_SOURCE_2	-0.160472
EXT_SOURCE_1	-0.155317
DAYS_EMPLOYED	-0.044932
FLOORSMAX_AVG	-0.044003
FLOORSMAX_MEDI	-0.043768
FLOORSMAX_MODE	-0.043226
AMT_GOODS_PRICE	-0.039645
REGION_POPULATION_RELATIVE	-0.037227
ELEVATORS_AVG	-0.034199

Name: TARGET, dtype: float64

In [13]:

```
corr_application_train = application_train.corr()['TARGET'].sort_values()
corr_application_train = corr_application_train.reset_index().rename(columns={0: 'correlation'})
corr_application_train
```

Out[13]:

	Attributes	Correlation
0	EXT_SOURCE_3	-0.178919
1	EXT_SOURCE_2	-0.160472
2	EXT_SOURCE_1	-0.155317
3	DAYS_EMPLOYED	-0.044932
4	FLOORSMAX_AVG	-0.044003
...
101	DAYSLAST_PHONE_CHANGE	0.055218
102	REGION_RATING_CLIENT	0.058899
103	REGION_RATING_CLIENT_W_CITY	0.060893
104	DAYSBIRTH	0.078239
105	TARGET	1.000000

106 rows × 2 columns

In [14]:

```
corr = datasets["undersampled_application_train"].corr()['TARGET']
corr=corr.sort_values(ascending=False)
```

In [15]:

```
print('NEGATIVE CORRELATIONS:\n', corr.tail(10))
```

NEGATIVE CORRELATIONS:

AMT_GOODS_PRICE	-0.063789
FLOORSMAX_MODE	-0.072400
FLOORSMAX_MEDI	-0.073826
DAYS_EMPLOYED	-0.074151
FLOORSMAX_AVG	-0.074277
EXT_SOURCE_1	-0.243774
EXT_SOURCE_2	-0.244929
EXT_SOURCE_3	-0.273157
FLAG_MOBIL	NaN
FLAG_DOCUMENT_12	NaN

Name: TARGET, dtype: float64

In [16]:

```
\n\nPOSITIVE CORRELATIONS\n', corr.head(10))
```

POSITIVE CORRELATIONS

	1.000000
TARGET	1.000000
DAYS_BIRTH	0.124241
REGION_RATING_CLIENT_W_CITY	0.097582
REGION_RATING_CLIENT	0.094539
DAYS_LAST_PHONE_CHANGE	0.087981
DAYS_ID_PUBLISH	0.080646
REG_CITY_NOT_WORK_CITY	0.079599
FLAG_EMP_PHONE	0.075783
FLAG_DOCUMENT_3	0.070787
REG_CITY_NOT_LIVE_CITY	0.068640

Name: TARGET, dtype: float64

In [17]:

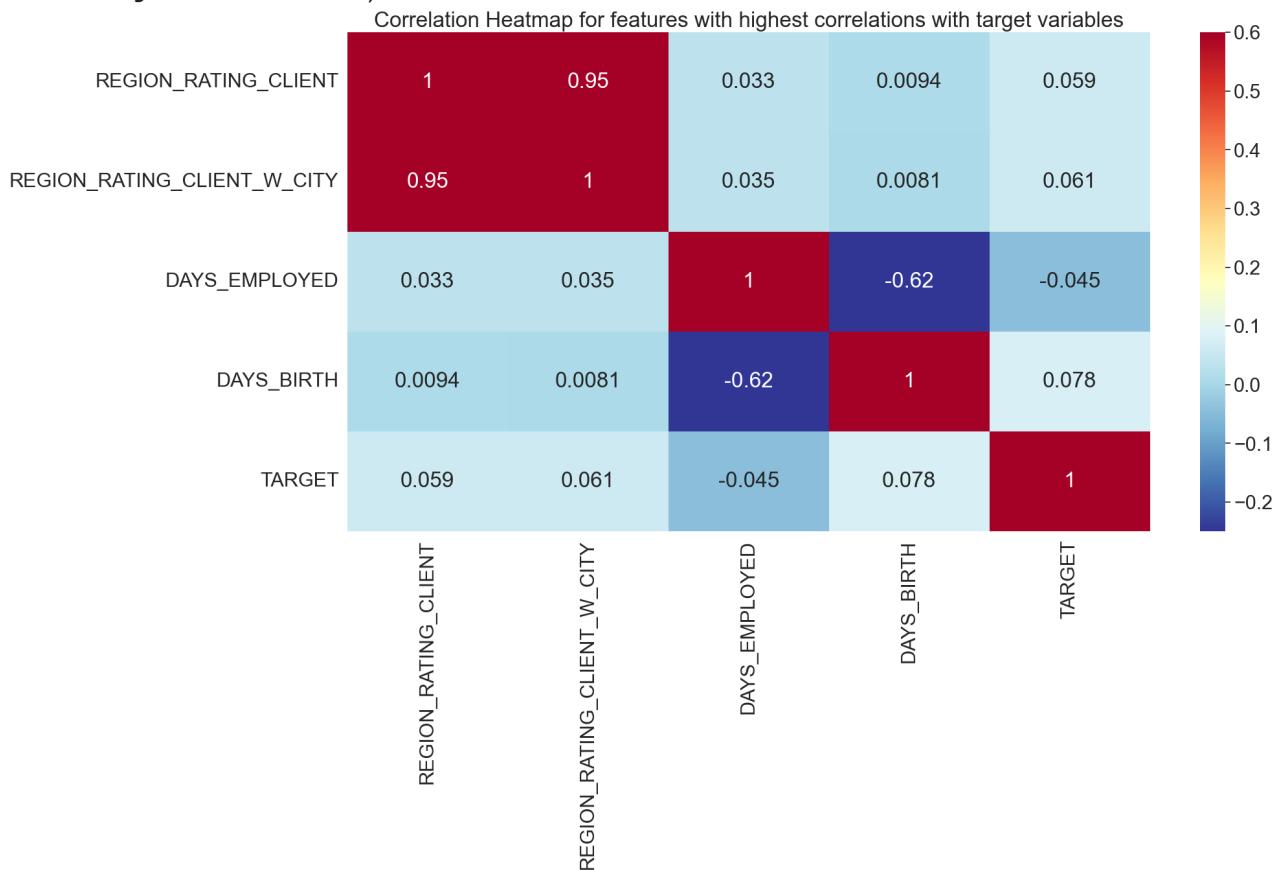
```
most_corr=datasets[ "application_train" ][['REGION_RATING_CLIENT',
                                             'REGION_RATING_CLIENT_W_CITY', 'DAYS_EMPLOYED', 'DAYS_B
most_corr_corr = most_corr.corr()

sns.set_style("dark")
sns.set_context("notebook", font_scale=2.0, rc={"lines.linewidth": 1.0})
fig, axes = plt.subplots(figsize = (20,10),sharey=True)
sns.heatmap(most_corr_corr,cmap=plt.cm.RdYlBu_r,vmin=-0.25,vmax=0.6,annot=T
plt.title('Correlation Heatmap for features with highest correlations with t

```

Out[17]:

Text(0.5, 1.0, 'Correlation Heatmap for features with highest correlations with target variables')



Feature Engineering:

In the process of feature engineering, we have utilized three tables namely, "Previous Applications", "Installment Payments", and "Credit Card Balance", from the secondary tables.

To identify the best customers of an organization, RFM features are employed. These features are based on three metrics: Recency, Frequency, and Monetary Value.

- **Recency** measures how recent the customer made a purchase
- **Frequency** determines how often they make purchases
- **Monetary Value** denotes how much money they spend on each purchase.

The frequency and monetary value metrics are indicative of the customer's engagement and their lifetime value, while recency is an indicator of retention and engagement.

Since we are analyzing the spending patterns of customers in this project, we use the RFM method to create features.

These features are generated by applying various functions such as **min, max, mean, sum, and count** to the relevant columns of the tables, thus producing new features that are significant for analysis.

Joining secondary tables with the primary table

In the case of the HCDR competition (and many other machine learning problems that involve multiple tables in 3NF or not) we need to join these datasets (denormalize) when using a machine learning pipeline. Joining the secondary tables with the primary table will lead to lots of new features about each loan application; these features will tend to be aggregate type features or meta data about the loan or its application. How can we do this when using Machine Learning Pipelines?

Joining previous_application with application_x

We refer to the `application_train` data (and also `application_test` data also)

as the **primary table** and the other files as the **secondary tables** (e.g., `previous_application` dataset). All tables can be joined using the primary key `SK_ID_PREV`.

Let's assume we wish to generate a feature based on previous application attempts. In this case, possible features here could be:

- A simple feature could be the number of previous applications.
- Other summary features of original features such as `AMT_APPLICATION`, `AMT_CREDIT` could be based on average, min, max, median, etc.

To build such features, we need to join the `application_train` data (and also `application_test` data also) with the 'previous_application' dataset (and the other available datasets).

When joining this data in the context of pipelines, different strategies come to mind with various tradeoffs:

1. Preprocess each of the non-application data sets, thereby generating many new (derived) features, and then joining (aka merge) the results with the `application_train` data (the labeled dataset) and with the `application_test` data (the unlabeled submission dataset) prior to processing the data (in a train, valid, test partition) via your machine learning pipeline. [This approach is recommended for this HCDR competition. WHY?]
- Do the joins as part of the transformation steps. [Not recommended here. WHY?]. How can this be done? Will it work?
 - This would be necessary if we had dataset wide features such as IDF (inverse document frequency) which depend on the entire subset of data as opposed to a single loan application (e.g., a feature about the relative amount applied for such as the percentile of the loan amount being applied for).

I want you to think about this section and build on this.

Roadmap for secondary table processing

1. Transform all the secondary tables to features that can be joined into the main table the application table (labeled and unlabeled)
 - 'bureau', 'bureau_balance', 'credit_card_balance', 'installments_payments',
 - 'previous_application', 'POS_CASH_balance'
- Merge the transformed secondary tables with the primary tables (i.e., the `application_train` data (the labeled dataset) and with the `application_test` data (the unlabeled submission dataset)), thereby leading to `X_train`, `y_train`, `X_valid`, etc.
- Proceed with the learning pipeline using `X_train`, `y_train`, `X_valid`, etc.
- Generate a submission file using the learnt model

Feature transformer

In [18]:

```
# create aggregate features (via pipeline)
class FeaturesAggregater(BaseEstimator, TransformerMixin):

    def __init__(self, features=None, agg_needed=["mean"]): # no *args or *
        self.agg_needed = agg_needed
        self.agg_op_features = {}
        for f in features:
            self.agg_op_features[f] = self.agg_needed[:]
    def fit(self, X, y=None):
        return self
    def transform(self, X, y=None):
        result = X.groupby(["SK_ID_CURR"]).agg(self.agg_op_features)
        df_result = pd.DataFrame()
        for x1, x2 in result.columns:
            new_col = x1 + " " + x2
            df_result[new_col] = result[x1][x2]
        df_result = df_result.reset_index(level=["SK_ID_CURR"])
        return df_result
```

Fixing Column names after Pandas agg() function to summarize grouped data

Since we have both the variable name and the operation performed in two rows in the Multi-Index dataframe, we can use that and name our new columns correctly.

For more details unstacking groupby results and examples please see [here](#)

For more details and examples please see [here](#)

Missing values in previous application table

In [19]:

```
# appsDF = datasets["previous_application"]
# appsDF.isna().sum()
```

In [20]:

```
# Access the 'previous_application' dataset from the 'datasets' container
previous_application_data = datasets["previous_application"]

# Apply the 'isna()' method on the 'previous_application_data' DataFrame to
# and then apply the 'sum()' method to count the number of missing values in
missing_values_count_per_column = previous_application_data.isna().sum()
missing_values_count_per_column
```

```
Out[20]: SK_ID_PREV          0  
          SK_ID_CURR         0  
          NAME_CONTRACT_TYPE 0  
          AMT_ANNUITY        372235  
          AMT_APPLICATION     0  
          AMT_CREDIT          1  
          AMT_DOWN_PAYMENT    895844  
          AMT_GOODS_PRICE      385515  
          WEEKDAY_APPR_PROCESS_START 0  
          HOUR_APPR_PROCESS_START 0  
          FLAG_LAST_APPL_PER_CONTRACT 0  
          NFLAG_LAST_APPL_IN_DAY 0  
          RATE_DOWN_PAYMENT   895844  
          RATE_INTEREST_PRIMARY 1664263  
          RATE_INTEREST_PRIVILEGED 1664263  
          NAME_CASH_LOAN_PURPOSE 0  
          NAME_CONTRACT_STATUS 0  
          DAYS_DECISION       0  
          NAME_PAYMENT_TYPE   0  
          CODE_REJECT_REASON 0  
          NAME_TYPE_SUITE     820405  
          NAME_CLIENT_TYPE   0  
          NAME_GOODS_CATEGORY 0  
          NAME_PORTFOLIO      0  
          NAME_PRODUCT_TYPE   0  
          CHANNEL_TYPE        0  
          SELLERPLACE_AREA    0  
          NAME_SELLER_INDUSTRY 0  
          CNT_PAYMENT         372230  
          NAME_YIELD_GROUP   0  
          PRODUCT_COMBINATION 346  
          DAYS_FIRST_DRAWING 673065  
          DAYS_FIRST_DUE      673065  
          DAYS_LAST_DUE_1ST_VERSION 673065  
          DAYS_LAST_DUE       673065  
          DAYS_TERMINATION    673065  
          NFLAG_INSURED_ON_APPROVAL 673065  
          dtype: int64
```

Feature engineering for prevApp table

The groupby output will have an index or multi-index on rows corresponding to your chosen grouping variables. To avoid setting this index, pass "as_index=False" to the groupby operation.

```
import pandas as pd
import dateutil

# Load data from csv file
data = pd.DataFrame.from_csv('phone_data.csv')
# Convert date from string to date times
data['date'] = data['date'].apply(dateutil.parser.parse,
dayfirst=True)

data.groupby('month', as_index=False).agg({"duration": "sum"})
```

Pandas `reset_index()` to convert Multi-Index to Columns We can simplify the multi-index dataframe using `reset_index()` function in Pandas. By default, Pandas `reset_index()` converts the indices to columns.

Columns on which the feature engineering is performed include "AMT_APPLICATION", "AMT_CREDIT", "AMT_ANNUITY", "approved_credit_ratio", "AMT_ANNUITY_credit_ratio", "Interest_ratio", "LTV_ratio", "SK_ID_PREV", "approved".

We have derived five new features from the previously mentioned features.

- The first feature is the **approved_credit_ratio**, which is the ratio of the credit amount the client requested on their previous application to the final credit amount approved on that application.
- The second feature is the **AMT_ANNUITY_credit_ratio**, which is the ratio of the annuity amount of the previous application to the final credit amount approved on that application.
- The third feature is the **Interest_ratio**, which is the ratio of the annuity amount of the previous application to the final credit amount approved on that application.
- The fourth feature is the **LTV_ratio**, which is the ratio of the final credit amount approved on the previous application to the goods price of the product the client applied for (if applicable) on the previous application.
- The fifth and final feature is the **approved**, which takes a value of 1 if the credit amount approved on the previous application is greater than 0, indicating that the application was approved.

In [21]:

```

previous_feature = ["AMT_APPLICATION", "AMT_CREDIT", "AMT_ANNUITY", "approved"]
agg_needed = ["min", "max", "mean", "count", "sum"]

agg_needed = ["min", "max", "mean", "count", "sum"]

def previous_feature_aggregation(df, feature, agg_needed):
    df['approved_credit_ratio'] = (df['AMT_APPLICATION']/df['AMT_CREDIT']).r
    # installment over credit approved ratio
    df['AMT_ANNUITY_credit_ratio'] = (df['AMT_ANNUITY']/df['AMT_CREDIT']).r
    # total interest payment over credit ratio
    df['Interest_ratio'] = (df['AMT_ANNUITY']/df['AMT_CREDIT']).replace(np.
    # loan cover ratio
    df['LTV_ratio'] = (df['AMT_CREDIT']/df['AMT_GOODS_PRICE']).replace(np.i
    df['approved'] = np.where(df.AMT_CREDIT >0 ,1, 0)

    test_pipeline = make_pipeline(FeaturesAggregater(feature, agg_needed))
    return(test_pipeline.fit_transform(df))

datasets['previous_application_agg'] = previous_feature_aggregation(dataset)

```

Missing value after the feature engineering

In [22]:

```
datasets["previous_application_agg"].isna().sum()
```

Out[22]:

SK_ID_CURR	0
AMT_APPLICATION_min	0
	dtype: int64

Missing values in Installment payments

In [23]:

```
datasets["installments_payments"].isna().sum()
```

Out[23]:

SK_ID_PREV	0
SK_ID_CURR	0
NUM_INSTALMENT_VERSION	0
NUM_INSTALMENT_NUMBER	0
DAYS_INSTALMENT	0
DAYS_ENTRY_PAYMENT	2905
AMT_INSTALMENT	0
AMT_PAYMENT	2905
	dtype: int64

Feature Engineering for Installment payments

Columns on which the feature engineering is performed include "**DAYSTINSTALMENTDIFF**", "**AMTPATMENTPCT**".

From the previous features, we have generated two additional features.

- The first feature is called **DAYSTINSTALMENTDIFF**, which is the difference between the date when the installment of the previous credit was due and the actual date when it was paid.
- The second feature is the **AMTPATMENTPCT**, which represents the percentage of the prescribed installment amount of the previous credit that the client actually paid on a particular installment, for every entry in the dataset.

Feature Engineering for Installment payments

In [24]:

```
payments_features = ["DAYSTINSTALMENTDIFF", "AMTPATMENTPCT"]

agg_needed = ["mean"]

def payments_feature_aggregation(df, feature, agg_needed):
    df['DAYSTINSTALMENTDIFF'] = df['DAYSTINSTALMENT'] - df['DAYSENTRYPAY']
    df['AMTPATMENTPCT'] = [x/y if (y != 0) & pd.notnull(y) else np.nan for x, y in zip(df['DAYSTINSTALMENT'], df['AMTPATMENT'])]
    test_pipeline = make_pipeline(FeaturesAggregator(feature, agg_needed))
    return(test_pipeline.fit_transform(df))

datasets['installments_payments_agg'] = payments_feature_aggregation(datasets, payments_features, agg_needed)
```

Missing value after the feature engineering

In [25]:

```
datasets['installments_payments_agg'].isna().sum()
```

Out [25]:

SK_ID_CURR	0
DAYSTINSTALMENTDIFF_mean	9
dtype:	int64

Missing value in Credit card balance

In [26]:

```
datasets['credit_card_balance'].isna().sum()
```

```
Out[26]: SK_ID_PREV          0
           SK_ID_CURR          0
           MONTHS_BALANCE       0
           AMT_BALANCE          0
           AMT_CREDIT_LIMIT_ACTUAL 0
           AMT_DRAWINGS_ATM_CURRENT 749816
           AMT_DRAWINGS_CURRENT    0
           AMT_DRAWINGS_OTHER_CURRENT 749816
           AMT_DRAWINGS_POS_CURRENT 749816
           AMT_INST_MIN_REGULARITY 305236
           AMT_PAYMENT_CURRENT     767988
           AMT_PAYMENT_TOTAL_CURRENT 0
           AMT_RECEIVABLE_PRINCIPAL 0
           AMT_RECVABLE            0
           AMT_TOTAL_RECEIVABLE    0
           CNT_DRAWINGS_ATM_CURRENT 749816
           CNT_DRAWINGS_CURRENT    0
           CNT_DRAWINGS_OTHER_CURRENT 749816
           CNT_DRAWINGS_POS_CURRENT 749816
           CNT_INSTALMENT_MATURE_CUM 305236
           NAME_CONTRACT_STATUS     0
           SK_DPD                  0
           SK_DPD_DEF               0
dtype: int64
```

Feature Engineering for Credit card balance

Columns on which the feature engineering is performed include "**AMT_BALANCE**", "**AMT_DRAWINGS_PCT**", "**AMT_DRAWINGS_ATM_PCT**", "**AMT_DRAWINGS_OTHER_PCT**", "**AMT_DRAWINGS_POS_PCT**", "**AMT_PRINCIPAL_RECEIVABLE_PCT**", "**CNT_DRAWINGS_ATM_CURRENT**", "**CNT_DRAWINGS_CURRENT**", "**CNT_DRAWINGS_OTHER_CURRENT**", "**CNT_DRAWINGS_POS_CURRENT**", "**SK_DPD**", "**SK_DPD_DEF**".

We have generated five new features using the previous features mentioned.

- The first feature is called **AMT_DRAWINGS_PCT**, which represents the ratio of the amount drawn during the previous credit month to the credit card limit during that month.
- The second feature is **AMT_DRAWINGS_ATM_PCT**, which is the ratio of the amount drawn at an ATM during the previous credit month to the credit card limit during that month.
- The third feature is **AMT_DRAWINGS_OTHER_PCT**, which is the ratio of the amount drawn for other purposes during the previous credit month to the credit card limit during that month.
- The fourth feature is **AMT_DRAWINGS_POS_PCT**, which is the ratio of the amount drawn or spent on goods during the previous credit month to the credit card limit during that month.
- Finally, the fifth feature is **MT_PRINCIPAL_RECEIVABLE_PCT**, which represents the ratio of the amount receivable for principal on the previous credit to the total amount receivable on that credit.

Feature Engineering for Credit card balance

In [27]:

```

credit_features = [
    "AMT_BALANCE",
    "AMT_DRAWINGS_PCT",
    "AMT_DRAWINGS_ATM_PCT",
    "AMT_DRAWINGS_OTHER_PCT",
    "AMT_DRAWINGS_POS_PCT",
    "AMT_PRINCIPAL_RECEIVABLE_PCT",
    "CNT_DRAWINGS_ATM_CURRENT",
    "CNT_DRAWINGS_CURRENT",
    "CNT_DRAWINGS_OTHER_CURRENT",
    "CNT_DRAWINGS_POS_CURRENT",
    "SK_DPD",
    "SK_DPD_DEF",
]

agg_needed = ["mean"]

def calculate_pct(x, y):
    return x / y if (y != 0) & pd.notnull(y) else np.nan
#def pct(x, y):
#    return x / y if (y != 0) & pd.notnull(y) else np.nan

def credit_feature_aggregation(df, feature, agg_needed):
    pct_columns = [
        ("AMT_DRAWINGS_CURRENT", "AMT_DRAWINGS_PCT"),
        ("AMT_DRAWINGS_ATM_CURRENT", "AMT_DRAWINGS_ATM_PCT"),
        ("AMT_DRAWINGS_OTHER_CURRENT", "AMT_DRAWINGS_OTHER_PCT"),
        ("AMT_DRAWINGS_POS_CURRENT", "AMT_DRAWINGS_POS_PCT"),
        ("AMT_RECEIVABLE_PRINCIPAL", "AMT_PRINCIPAL_RECEIVABLE_PCT"),
    ]

    for col_x, col_pct in pct_columns:
        df[col_pct] = [calculate_pct(x, y) for x, y in zip(df[col_x], df["A"])]
    pipeline = make_pipeline(FeaturesAggregater(feature, agg_needed))
    return pipeline.fit_transform(df)

datasets["credit_card_balance_agg"] = credit_feature_aggregation(
    datasets["credit_card_balance"], credit_features, agg_needed
)

```

Missing values after feature engineering

In [28]:

```
datasets["credit_card_balance_agg"].isna().sum()
```

```
Out[28]: SK_ID_CURR      0
          AMT_BALANCE_mean    0
          dtype: int64
```

Join the feature engineered datasets with under-sampled datasets

```
In [29]: datasets.keys()
```

```
Out[29]: dict_keys(['application_train', 'application_test', 'bureau', 'bureau_balance',
       'credit_card_balance', 'installments_payments', 'previous_application',
       'POS_CASH_balance', 'undersampled_application_train', 'undersampled_application_train_2',
       'previous_application_agg', 'installments_payments_agg', 'credit_card_balance_agg'])
```

```
In [30]: # Load the train dataset
train_data = datasets["application_train"]

# Compute the distribution of the target variable
target_counts = train_data['TARGET'].value_counts()

# Display the target distribution
print("Target variable distribution:\n")
print(target_counts)
print("\n")

# Compute the percentage of positive and negative examples in the dataset
positive_count = target_counts[1]
negative_count = target_counts[0]
total_count = positive_count + negative_count
positive_percentage = (positive_count / total_count) * 100
negative_percentage = (negative_count / total_count) * 100

# Display the percentages of positive and negative examples
print(f"Percentage of positive examples: {positive_percentage:.2f}%")
print(f"Percentage of negative examples: {negative_percentage:.2f}%")
```

Target variable distribution:

0	282686
1	24825
Name: TARGET, dtype: int64	

Percentage of positive examples: 8.07%
 Percentage of negative examples: 91.93%

```
In [31]: train_dataset= datasets["undersampled_application_train"] #primary dataset  
merge_all_data = True  
  
# merge primary table and secondary tables using features based on meta data  
if merge_all_data:  
    # 1. Join/Merge in prevApps Data  
    train_dataset = train_dataset.merge(datasets["previous_application_agg"]  
  
    # 2. Join/Merge in Installments Payments Data  
    train_dataset = train_dataset.merge(datasets["installments_payments_agg"]  
  
    # 3. Join/Merge in Credit Card Balance Data  
    train_dataset = train_dataset.merge(datasets["credit_card_balance_agg"])
```

```
In [32]: datasets["undersampled_application_train_4"] = train_dataset
```

```
In [33]: train_dataset.shape
```

```
Out[33]: (99825, 125)
```

```
In [34]: train_dataset = datasets["undersampled_application_train_2"]  
train_dataset = train_dataset.merge(datasets["previous_application_agg"], h  
train_dataset = train_dataset.merge(datasets["installments_payments_agg"],  
train_dataset = train_dataset.merge(datasets["credit_card_balance_agg"], ho  
train_dataset = train_dataset.drop(columns = 'weight')  
datasets["undersampled_application_train_4_2"] = train_dataset
```

```
In [35]: train_dataset.shape
```

```
Out[35]: (49650, 125)
```

```
In [36]: train_dataset.to_csv('train_dataset.csv', index=False)
```

Join the unlabeled dataset (i.e., the submission file)

In [37]:

```
x_kaggle_test= datasets["application_test"]

# merge primary table and secondary tables using features based on meta data
if merge_all_data:
    # 1. Join/Merge in prevApps Data
    X_kaggle_test = X_kaggle_test.merge(datasets["previous_application_agg"])

    # 2. Join/Merge in Installments Payments Data
    X_kaggle_test = X_kaggle_test.merge(datasets["installments_payments_agg"])

    # 3. Join/Merge in Credit Card Balance Data
    X_kaggle_test = X_kaggle_test.merge(datasets["credit_card_balance_agg"])
```

In [38]:

```
X_kaggle_test.to_csv('X_kaggle_test.csv', index=False)
```

Loss Functions

LOG LOSS

It is a metric that reflects the degree of proximity between the predicted probability and the true value (which is typically 0 or 1 in the context of binary classification). The log-loss value increases as the predicted probability increasingly deviates from the true value.

$$\text{LogLoss} = -\frac{1}{m} \sum_{i=1}^m (y_i \cdot \log(p_i) + (1 - y_i) \cdot \log(1 - p_i)) + \lambda \sum_{j=1}^n |\theta_j|$$

Squared Hinge Loss

Hinge loss is primarily employed in support vector machine algorithms. The quintessential boundary in any classification challenge endeavors to minimize misclassifications. To embody this concept algorithmically, hinge loss imposes penalties on each erroneous classification. The squared hinge loss is essentially the hinge loss squared. This variant is employed when it is desirable to accentuate the penalization of substantial errors.

$$\text{SquareHingeLoss} = \frac{1}{m} \sum_{i=1}^m \left[\max(0, 1 - y_i (\theta^T \cdot x_i + b))^2 \right] + \lambda \sum_{j=1}^n \theta_j^2$$

Entropy:

Entropy is a metric for the degree of disorderliness within a node. In the context of decision trees, disorderliness refers to a scenario in which all classes are uniformly represented in the data.

$$H(p) = - \sum_{i=1}^c p_i \log_2(p_i)$$

The symbol 'pi' in the above context pertains to the ratio of instances belonging to a specific class 'i', while 'c' in the formula denotes the count of unique values that can be attributed to the response variable.

Gini Impurity:

Gini Impurity is a quantification of the dispersion among distinct categories.

$$S = \sum_{i=1}^c p_i(1 - p_i)$$

The 'pi' above refers to the fraction of the observations that belong to a particular class 'i' and the 'c' given in the formula is the number of different possible values of the response variable.

Evaluation metrics

The evaluation of submissions is conducted through the calculation of the area under the ROC curve, which measures the relationship between the predicted probability and the observed target. The SkLearn `roc_auc_score` function is utilized to compute the AUC or AUROC, effectively summarizing the information contained in the ROC curve into a single numerical value.

Submissions are evaluated on [area under the ROC curve](#) between the predicted probability and the observed target.

The SkLearn `roc_auc_score` function computes the area under the receiver operating characteristic (ROC) curve, which is also denoted by AUC or AUROC. By computing the area under the roc curve, the curve information is summarized in one number.

```
from sklearn.metrics import roc_auc_score
>>> y_true = np.array([0, 0, 1, 1])
>>> y_scores = np.array([0.1, 0.4, 0.35, 0.8])
>>> roc_auc_score(y_true, y_scores)
0.75
```

ACCURACY

It refers to the proportion of accurately classified data instances in relation to the overall number of data instances.

$$\text{Accuracy} = \frac{TN + TP}{TN + FP + TP + FN}$$

PRECISION:

Precision refers to the ratio of true positives to the sum of true positives and false positives.

$$\text{Precision} = \frac{TP}{TP + FP}$$

RECALL

It denotes the fraction of positive instances that are correctly identified as positive by the model. This metric is equivalent to the TPR (True Positive Rate).

$$\text{Recall} = \frac{TP}{TP + FN}$$

F1 SCORE

It is the harmonic mean of accuracy and recall, taking into account both false positives and false negatives. It is a useful metric for evaluating models on imbalanced datasets.

$$\text{F1Score} = \frac{Precision * Recall}{Precision + Recall}$$

A CONFUSION MATRIX

It is a tabular representation consisting of two axes - one representing the actual values and the other representing the predicted values. The matrix is of size 2x2 and is commonly used in classification tasks to assess the performance of a model.

Confusion Matrix

Create confusion matrix for baseline model

In [39]:

```
class_labels = ["No Default", "Default"]
import numpy as np
from sklearn.metrics import confusion_matrix

def confusion_matrix_normalized(model, X_train, y_train, X_test, y_test):
    y_train_pred = model.predict(X_train)
    y_test_pred = model.predict(X_test)

    cm_train = confusion_matrix(y_train, y_train_pred, normalize='true').astype('float')
    cm_test = confusion_matrix(y_test, y_test_pred, normalize='true').astype('float')

    return cm_train, cm_test
```

Processing pipeline

In [40]:

```
# Create a class to select numerical or categorical columns
from sklearn.base import BaseEstimator, TransformerMixin

# Create a transformer to select numerical or categorical columns
class ColumnSelector(BaseEstimator, TransformerMixin):
    def __init__(self, columns):
        self.columns = columns

    def fit(self, X, y=None):
        return self

    def transform(self, X):
        return X[self.columns].values
```

In [41]:

```
def pct(x):
    return round(100*x,3)
```

Empirical Findings

Every trial in this undertaking shall be identified by Pipeline configuration

I intend to formulate five distinct experiments, predicated on the ensuing criteria:

all normal features

undersampled features

baseline machine learning models.

Prior to embarking on model selection through the process of model comparison and model design, it is imperative to initialize a logbook and establish a precise loss metric framework. This shall allow for the systematic tracking and monitoring of model performance across various iterations, and facilitate the comprehensive evaluation of model efficacy in accordance with pre-defined objectives. Such a structured approach is essential to ensure robust and accurate analyses, and to enable the identification of optimal model configurations that deliver maximal predictive power and generalizability.

Tracking results in dataframe

In [42]:

```
try:  
    expLog  
except NameError:  
    expLog = pd.DataFrame(columns=[ "exp_name",  
                                    "description",  
                                    "Train Time (sec)",  
                                    "Test Time (sec)",  
                                    "Train Acc",  
                                    "Valid Acc",  
                                    "Test Acc",  
                                    "Train AUC",  
                                    "Valid AUC",  
                                    "Test AUC",  
                                    "Train F1 Score",  
                                    "Valid F1 Score",  
                                    "Test F1 Score"  
                                ])
```

In [43]:

```
def get_results(expLog, exp_name, description, model, train_time, test_time  
expLog.loc[len(expLog)] = [f"{exp_name}", description] + list(np.round(  
    [train_time, test_time,  
     accuracy_score(y_train, model.predict(X_train)),  
     accuracy_score(y_valid, model.predict(X_valid)),  
     accuracy_score(y_test, model.predict(X_test)),  
     roc_auc_score(y_train, model.predict_proba(X_train)[:, 1]),  
     roc_auc_score(y_valid, model.predict_proba(X_valid)[:, 1]),  
     roc_auc_score(y_test, model.predict_proba(X_test)[:, 1]),  
     f1_score(y_train, model.predict(X_train)),  
     f1_score(y_valid, model.predict(X_valid)),  
     f1_score(y_test, model.predict(X_test))],  
    4))  
return expLog
```

OHE when previously unseen unique values in the test/validation set

Train, validation and Test sets (and the leakage problem we have mentioned previously):

Let's look at a small usecase to tell us how to deal with this:

- The OneHotEncoder is fitted to the training set, which means that for each unique value present in the training set, for each feature, a new column is created. Let's say we have 39 columns after the encoding up from 30 (before preprocessing).
- The output is a numpy array (when the option sparse=False is used), which has the disadvantage of losing all the information about the original column names and values.
- When we try to transform the test set, after having fitted the encoder to the training set, we obtain a `ValueError`. This is because there are new, previously unseen unique values in the test set and the encoder doesn't know how to handle these values. In order to use both the transformed training and test sets in machine learning algorithms, we need them to have the same number of columns.

This last problem can be solved by using the option `handle_unknown='ignore'` of the OneHotEncoder, which, as the name suggests, will ignore previously unseen values when transforming the test set.

Here is a example that in action:

```
# Identify the categorical features we wish to consider.
cat_attribs = ['CODE_GENDER',
'FLAG_OWN_REALTY', 'FLAG_OWN_CAR', 'NAME_CONTRACT_TYPE',
'NAME_EDUCATION_TYPE', 'OCCUPATION_TYPE', 'NAME_INCOME_TYPE']

# Notice handle_unknown="ignore" in OHE which ignore values from
# the validation/test that
# do NOT occur in the training set
cat_pipeline = Pipeline([
    ('selector', DataFrameSelector(cat_attribs)),
    ('imputer', SimpleImputer(strategy='most_frequent')),
    ('ohe', OneHotEncoder(sparse=False,
handle_unknown="ignore"))])
```

MACHINE LEARNING PIPELINES

HCDR preprocessing with all columns

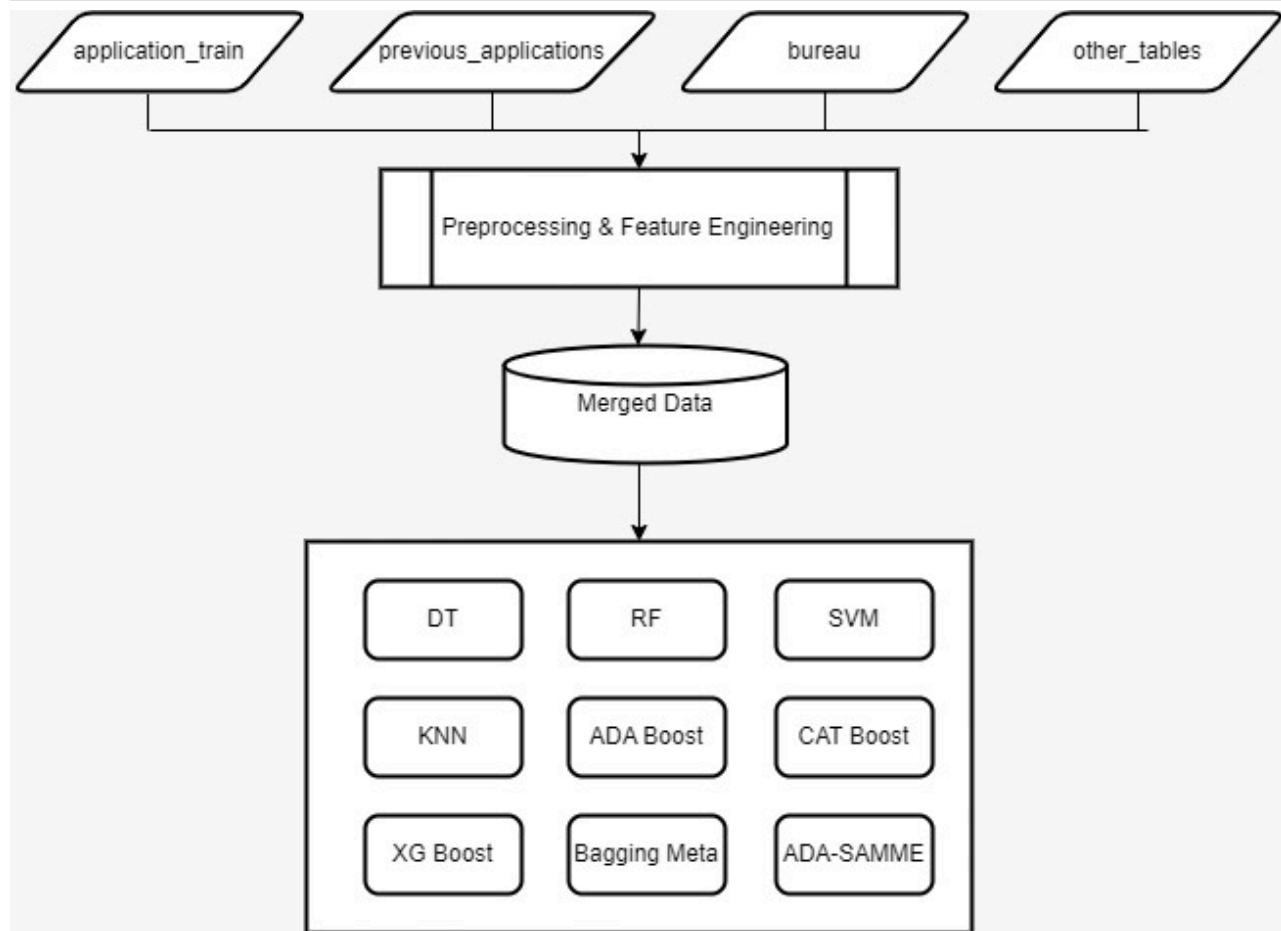
By opting to incorporate a variety of data sources, including the Previous Application, Installment Payments, and Credit Card Applications tables, our analysis benefits from the diversification of our dataset, circumventing the limitations posed by solely utilizing the application_train data. Given the imbalanced nature of the data, which prompted the implementation of undersampling techniques to offset the surplus of non-defaulters (those classified with a target variable of 0), leveraging additional tables with relevant features is imperative for establishing a robust predictive model. Furthermore, our correlation analysis has confirmed the salience of features contained within these other tables, further substantiating our decision to utilize this multi-faceted approach.

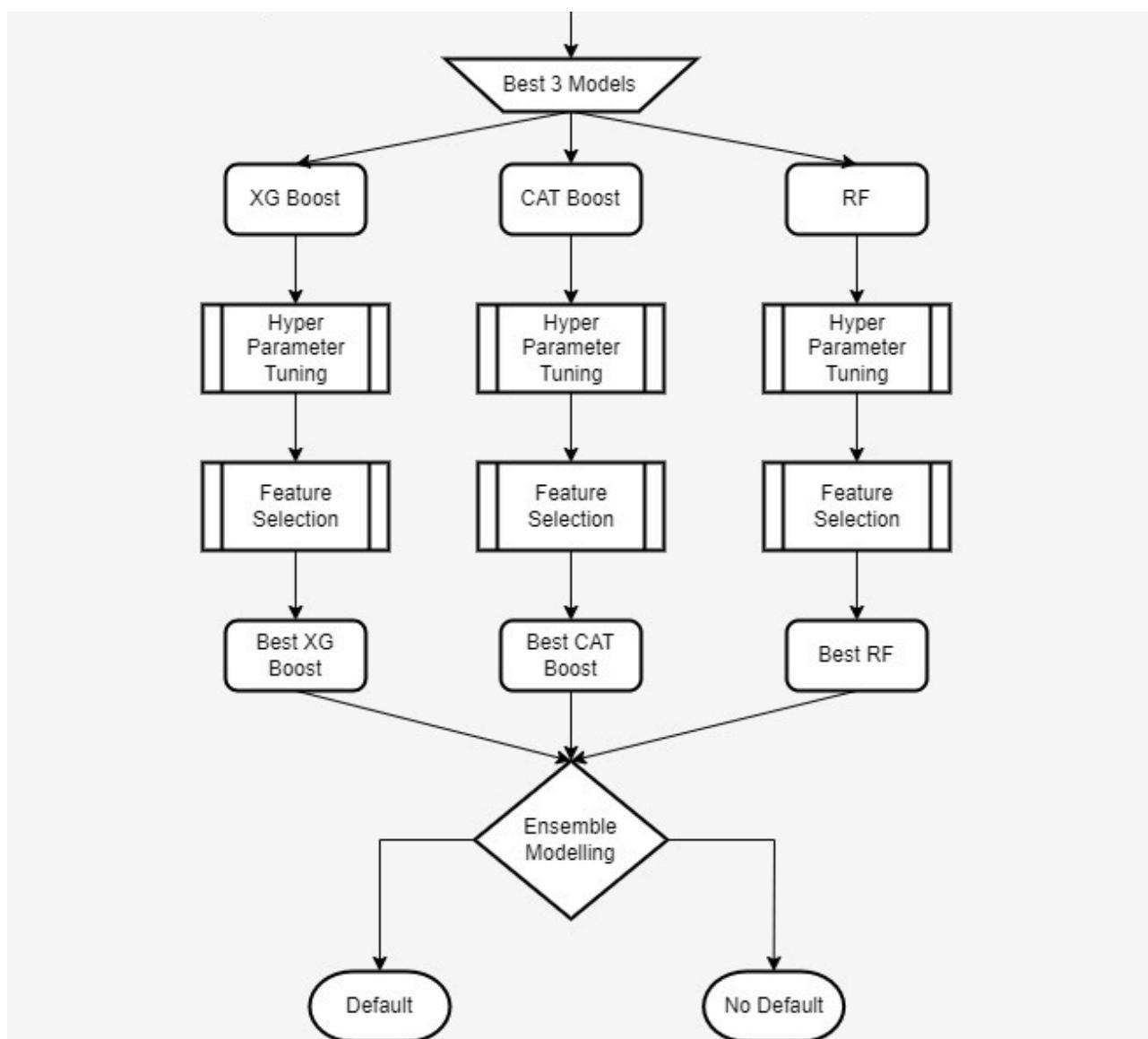
Block Diagram

In [71]:

```
from IPython.display import Image
Image(filename='Block_dee_pipeline.jpeg')
```

Out [71]:





Selecting the numerical and categorical features

In [44]:

```

# Load the undersampled training dataset
train_dataset = datasets["undersampled_application_train_4"]

# Separate numerical and categorical features
numerical_features = []
categorical_features = []

for feature_name in train_dataset:
    # Check if feature is numerical or categorical
    if train_dataset[feature_name].dtype in [np.float64, np.int64]:
        numerical_features.append(feature_name)
    else:
        categorical_features.append(feature_name)

# Remove target and ID columns from numerical features
numerical_features.remove('TARGET')
numerical_features.remove('SK_ID_CURR')

# Define pipelines for categorical and numerical features
categorical_pipeline = Pipeline([
    ('selector', ColumnSelector(categorical_features)), # Select categorical
    ('imputer', SimpleImputer(strategy='most_frequent')), # Impute missing
    ('one_hot_encoder', OneHotEncoder(sparse=False, handle_unknown="ignore"))
])

numerical_pipeline = Pipeline([
    ('selector', ColumnSelector(numerical_features)), # Select numerical features
    ('imputer', SimpleImputer(strategy='mean')), # Impute missing values with mean
    ('standard_scaler', StandardScaler()), # Standardize numerical features
])

# Combine pipelines for numerical and categorical features
data_prep_pipeline = FeatureUnion(transformer_list=[
    ("numerical_pipeline", numerical_pipeline),
    ("categorical_pipeline", categorical_pipeline),
])

# Compute the total number of features, as well as the number of numerical
selected_features = numerical_features + categorical_features + ["SK_ID_CURR"]
total_features = f"Total Features: {len(selected_features)} - Numerical: {len(numerical_features)} - Categorical: {len(categorical_features)}"

print(total_features) # Print the total number of features and their breakdown

```

Total Features: 124 - Numerical: 107, Categorical: 16

Create Train and Test Datasets

In [45]:

```
y_train = train_dataset['TARGET']
X_train = train_dataset[selected_features]
X_train, X_valid, y_train, y_valid = train_test_split(X_train, y_train, test_size=0.2, random_state=42)
X_train, X_test, y_train, y_test = train_test_split(X_train, y_train, test_size=0.2, random_state=42)
X_kaggle_test = X_kaggle_test[selected_features]

print(f"X train shape: {X_train.shape}")
print(f"X validation shape: {X_valid.shape}")
print(f"X test shape: {X_test.shape}")
print(f"X kaggle test shape: {X_kaggle_test.shape}")
```

```
X train shape: (72123, 124)
X validation shape: (14974, 124)
X test shape: (12728, 124)
X kaggle test shape: (48744, 124)
```

LEARNER MODELS

In order to establish a benchmark for comparison, we shall employ the utilization of a logistic regression model, which will make use of certain preprocessed features, processed by our established pipeline.

For the HCDR project, a total of 10 machine learning models were constructed to accurately classify the credit defaluters . Among these models, the three best-performing ones were selected based on their superior performance metrics. These three models were then subject to hyperparameter tuning and feature selection to optimize their performance. The hyperparameter tuning involved fine-tuning the various parameters of the models to identify the optimal set of parameters that produce the highest accuracy, precision, and recall scores. Meanwhile, feature selection aimed to identify the most relevant features that are highly correlated with the target variable to eliminate irrelevant variables that may affect the accuracy of the models.

Model-1 baseline using training columns for the application with LogisticRegression()

In [46]:

```
from sklearn.metrics import roc_curve
# Logistic Regression model with under-sampled data
np.random.seed(42)

# Define a pipeline that includes data preparation and logistic regression
full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline), # Data preparation pipeline
    ("linear", LogisticRegression()) # Logistic Regression model
])

# Train the model and measure the training time
start_time = time.time()
model = full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start_time, 4)

# Evaluate the model on the test set and measure the test time
start_time = time.time()
test_score = full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start_time, 4)

# Define an experiment name based on the number of selected features
experiment_name = f"Model-1 Baseline LR"
experiment_description = f"Logistic regression with undersampled data {len(s) - 1} selected features"

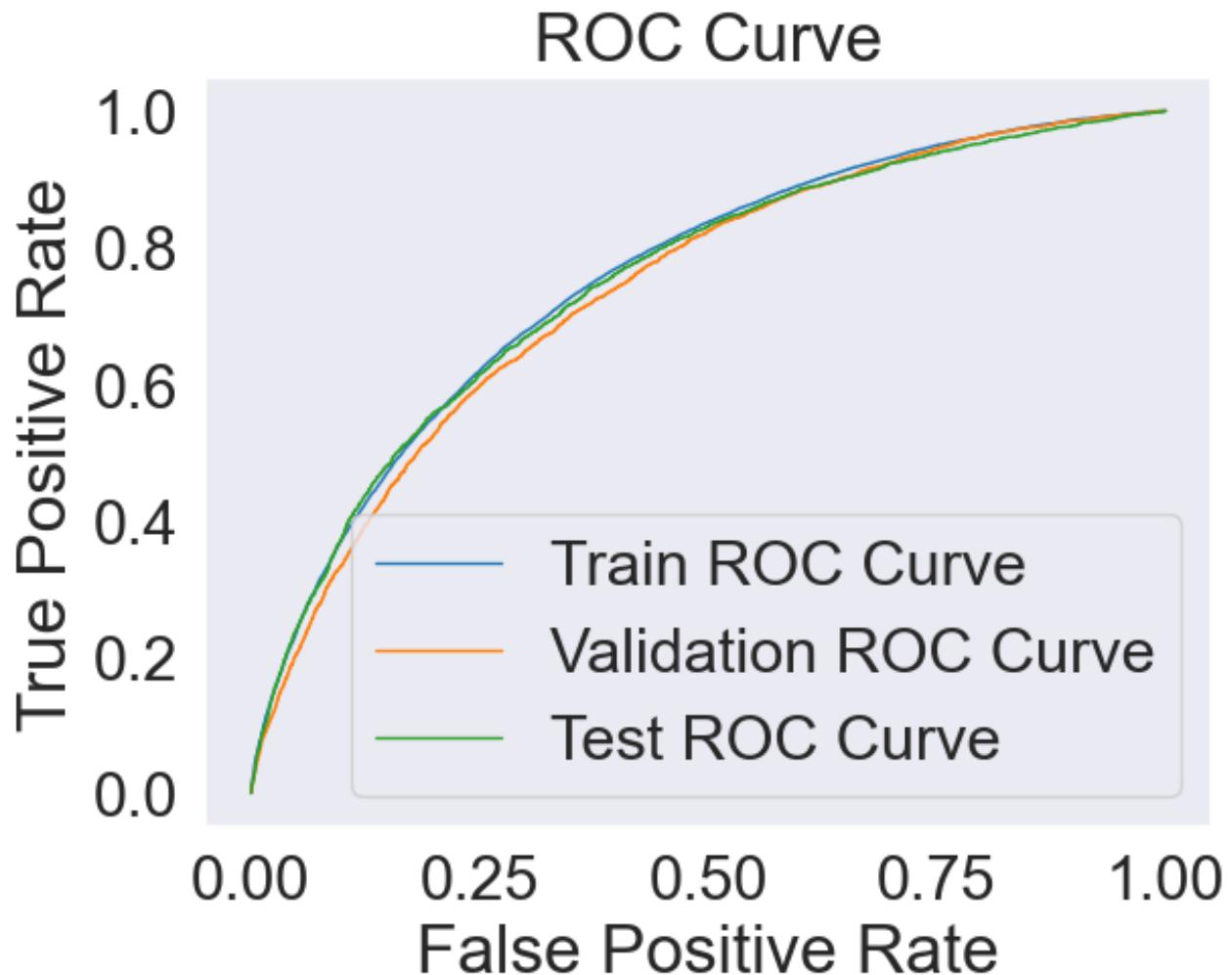
# Log the results of the experiment
expLog = get_results(expLog, experiment_name, experiment_description, model, train_time, test_time)

# Model Training and Validation
model.fit(X_train, y_train)

# Model Predictions
train_preds = model.predict_proba(X_train)[:, 1]
valid_preds = model.predict_proba(X_valid)[:, 1]
test_preds = model.predict_proba(X_test)[:, 1]

# Compute Metrics
train_fpr, train_tpr, _ = roc_curve(y_train, train_preds)
valid_fpr, valid_tpr, _ = roc_curve(y_valid, valid_preds)
test_fpr, test_tpr, _ = roc_curve(y_test, test_preds)
# Plot ROC Curve
plt.plot(train_fpr, train_tpr, label="Train ROC Curve")
plt.plot(valid_fpr, valid_tpr, label="Validation ROC Curve")
plt.plot(test_fpr, test_tpr, label="Test ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()

expLog
```



Out [46]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC	Test AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379	0.7481

After performing correlation analysis, specific columns were selected from the "application train" dataset

In [47]:

```
## Model-2 baseline using training columns for the application with Logisti
```

Creating data preparation pipeline from given columns

In [48]:

```

def get_pipeline(num_cols = None):
    # Load the undersampled training dataset and join with additional features
    train_dataset = datasets["undersampled_application_train_2"]
    train_dataset = train_dataset.merge(datasets["previous_application_agg"])
    train_dataset = train_dataset.merge(datasets["installments_payments_agg"])
    train_dataset = train_dataset.merge(datasets["credit_card_balance_agg"])

    # Separate numerical and categorical features
    numerical_features = []
    categorical_features = []
    for feature_name in train_dataset:
        if train_dataset[feature_name].dtype in [np.float64, np.int64]:
            numerical_features.append(feature_name)
        else:
            categorical_features.append(feature_name)

    # Remove unnecessary features
    numerical_features.remove('TARGET')
    numerical_features.remove('weight')
    numerical_features.remove('SK_ID_CURR')

    # Define pipelines for categorical and numerical features
    categorical_pipeline = Pipeline([
        ('selector', ColumnSelector(categorical_features)), # Select categorical features
        ('imputer', SimpleImputer(strategy='most_frequent')), # Impute missing values
        ('one_hot_encoder', OneHotEncoder(sparse=False, handle_unknown="ignore"))
    ])

    # If columns are provided, use only those columns for numerical pipeline
    if num_cols == None:
        final_numerical_features = numerical_features
    else:
        final_numerical_features = num_cols

    numerical_pipeline = Pipeline([
        ('selector', ColumnSelector(final_numerical_features)), # Select numerical features
        ('imputer', SimpleImputer(strategy='mean')), # Impute missing values
        ('standard_scaler', StandardScaler()), # Standardize numerical features
    ])

    # Combine pipelines for numerical and categorical features
    data_prep_pipeline = FeatureUnion(transformer_list=[
        ("numerical_pipeline", numerical_pipeline),
        ("categorical_pipeline", categorical_pipeline),
    ])

    # Compute the total number of features, as well as the number of numerical features
    selected_features = final_numerical_features + categorical_features + [len(categorical_pipeline)]
    total_features = f"Total Features: {len(selected_features)} - Numerical Features: {len(numerical_pipeline)}"

    # Print the total number of features and their breakdown
    print(total_features)

```

```
    return data_prep_pipeline, selected_features
```

Creating Train, Test, and Validation datasets

In [49]:

```
# Load the undersampled training dataset and join with additional feature d
train_dataset = datasets["undersampled_application_train_2"]
train_dataset = train_dataset.merge(datasets["previous_application_agg"], how="left")
train_dataset = train_dataset.merge(datasets["installments_payments_agg"], how="left")
train_dataset = train_dataset.merge(datasets["credit_card_balance_agg"], how="left")

# Select the target variable
y_train = train_dataset['TARGET']

# Select the features for the training set
X_train = train_dataset[selected_features]

# Split the data into training and validation sets
# The training set will be used to train the model, and the validation set will be used to validate the model
X_train, X_valid, y_train, y_valid = train_test_split(X_train, y_train, test_size=0.2, random_state=42)

# Split the training set into training and test sets
# The training set will be used to train the model, and the test set will be used to test the model
X_train, X_test, y_train, y_test = train_test_split(X_train, y_train, test_size=0.3, random_state=42)

# Print the shapes of the training, validation, and test sets
print(f"Training set shape: {X_train.shape}")
print(f"Validation set shape: {X_valid.shape}")
print(f"Test set shape: {X_test.shape}")
```

```
Training set shape: (35871, 124)
Validation set shape: (7448, 124)
Test set shape: (6331, 124)
```

Model-2 baseline using LogisticRegression()

For establishing a baseline, we shall employ certain processed features stemming from the pipeline. The logistic regression model shall serve as the rudimentary benchmark model 2 . We will use two undersampled data here .

In [50]:

```
import matplotlib.pyplot as plt
import matplotlib.patheffects as path_effects

data = [len(numerical_features), len(categorical_features)]
labels = ['Numerical Features ', 'Categorical Features']

fig, ax = plt.subplots()
bars = ax.bar(labels, data, color=['#0072B2', '#E69F00'], edgecolor='black')

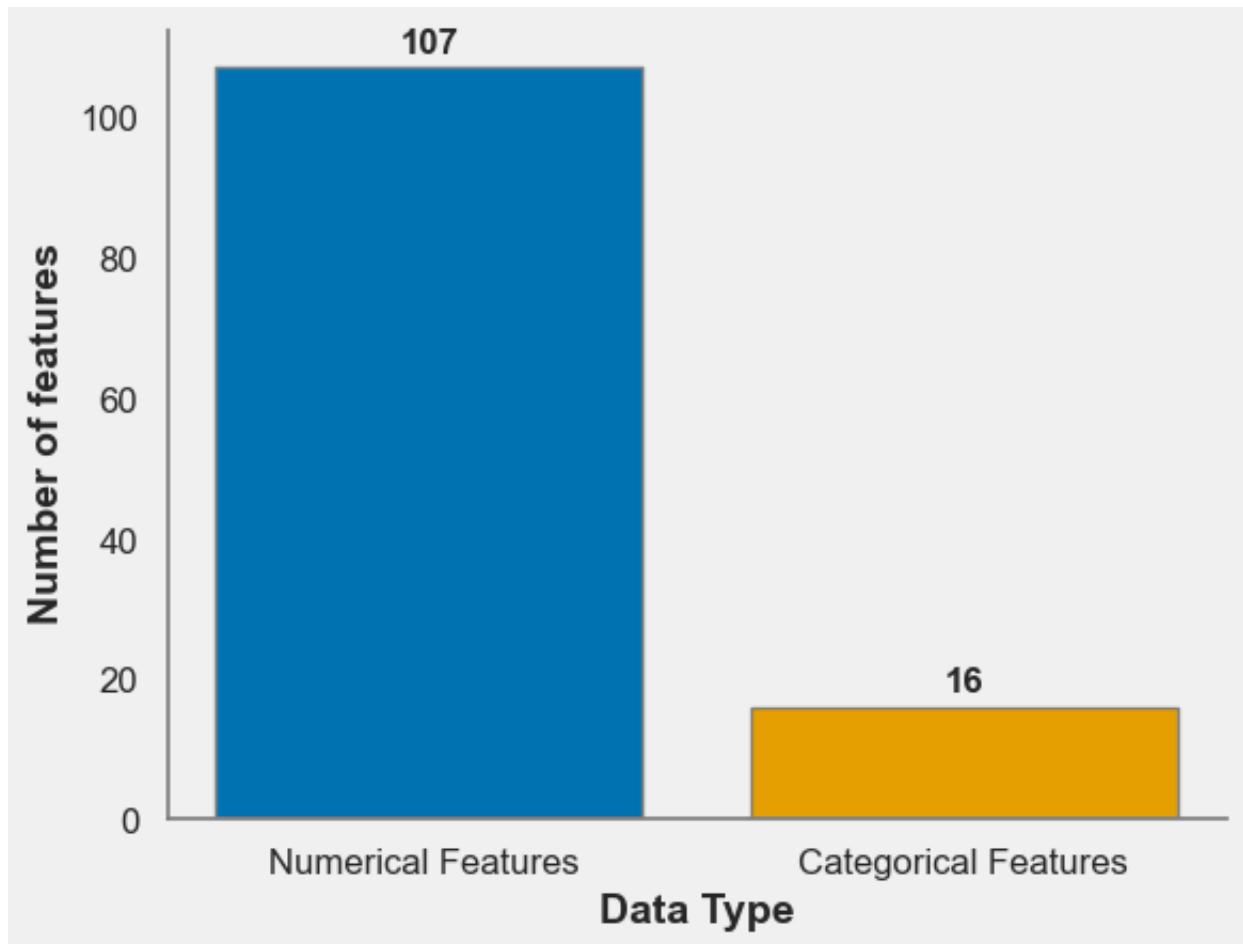
# Add shadows to the bars
for bar in bars:
    bar.set_edgecolor('gray')
    bar.set linewidth(1)
    bar.set zorder(0)

# Add labels to the bars
height = bar.get_height()
ax.annotate(f'{height:.0f}', xy=(bar.get_x() + bar.get_width() / 2, hei
    xytext=(0, 3), textcoords='offset points', ha='center', va=
    fontsize=12, fontweight='bold')

# Customize the axis labels and ticks
ax.set_xlabel('Data Type', fontsize=14, fontweight='bold')
ax.set_ylabel('Number of features ', fontsize=14, fontweight='bold')
ax.tick_params(axis='both', labelsize=12)

# Customize the plot background
ax.set_facecolor('#F0F0F0')
fig.set_facecolor('#F0F0F0')
ax.spines['bottom'].set_color('gray')
ax.spines['left'].set_color('gray')
ax.spines['top'].set_visible(False)
ax.spines['right'].set_visible(False)

plt.show()
```



In [51]:

```
# Set the random seed for reproducibility
np.random.seed(42)

# Create pipeline for preparing the data and select features
data_prep_pipeline, selected_features = get_pipeline()

# Join the preparation pipeline with logistic regression model
full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("logistic_regression", LogisticRegression())
])

# Train the model on the training set
start = time.time()
model = full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

# Compute cross-validation scores
cv_splits = ShuffleSplit(n_splits=3, test_size=0.7, random_state=42)
logit_scores = cross_val_score(full_pipeline_with_predictor, X_train, y_train)

print("Cross-validation scores:", logit_scores)

# Compute the test score
```

```
start = time.time()
test_score = full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

print("Test score:", test_score)

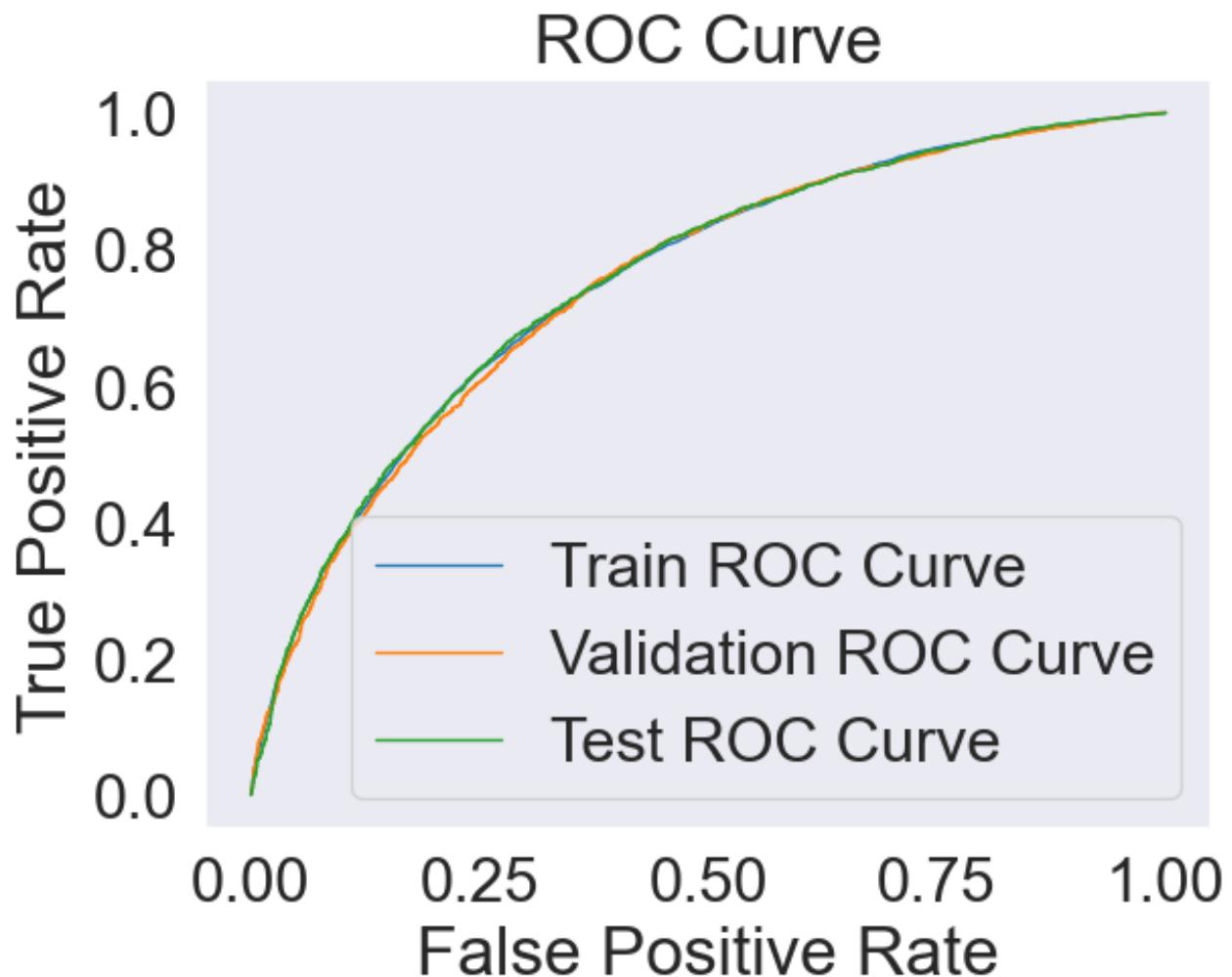
# Save the experiment results
exp_name = f"Model-2 Baseline LR"
experiment_description = f"Logistic regression with undersampled data-2 {len(expLog)} rows"
get_results(expLog, exp_name, experiment_description, model, train_time, test_time)

# Model Training and Validation
model.fit(X_train, y_train)

# Model Predictions
train_preds = model.predict_proba(X_train)[:, 1]
valid_preds = model.predict_proba(X_valid)[:, 1]
test_preds = model.predict_proba(X_test)[:, 1]

# Compute Metrics
train_fpr, train_tpr, _ = roc_curve(y_train, train_preds)
valid_fpr, valid_tpr, _ = roc_curve(y_valid, valid_preds)
test_fpr, test_tpr, _ = roc_curve(y_test, test_preds)
# Plot ROC Curve
plt.plot(train_fpr, train_tpr, label="Train ROC Curve")
plt.plot(valid_fpr, valid_tpr, label="Validation ROC Curve")
plt.plot(test_fpr, test_tpr, label="Test ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()
expLog
```

Total Features: 124 - Numerical: 107, Categorical: 16
Cross-validation scores: [0.67953007 0.67992832 0.67905217]
Test score: 0.6904122571473701



Out [51]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC	Test AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379	0.7481
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4864	0.0465	0.6876	0.6843	0.6904	0.7525	0.7489	0.7535

Model-3 KNN model

In [52]:

```
from sklearn.neighbors import KNeighborsClassifier
```

```
# Set the random seed for reproducibility
np.random.seed(42)

# Create pipeline for preparing the data and select features
data_prep_pipeline, selected_features = get_pipeline()

# Join the preparation pipeline with KNN model
knn_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("knn", KNeighborsClassifier(n_neighbors=11, p=2))
])

# Train the model on the training set
start = time.time()
model = knn_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

# Compute the test score
start = time.time()
test_score = knn_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

print("Test score:", test_score)

# Results
# Save the experiment results
exp_name = f"Model-3 KNN"
experiment_description = f"KNN with undersampled data-2 {len(selected_features)} features"
expLog = get_results(expLog, exp_name, experiment_description, model, train_time, test_time)

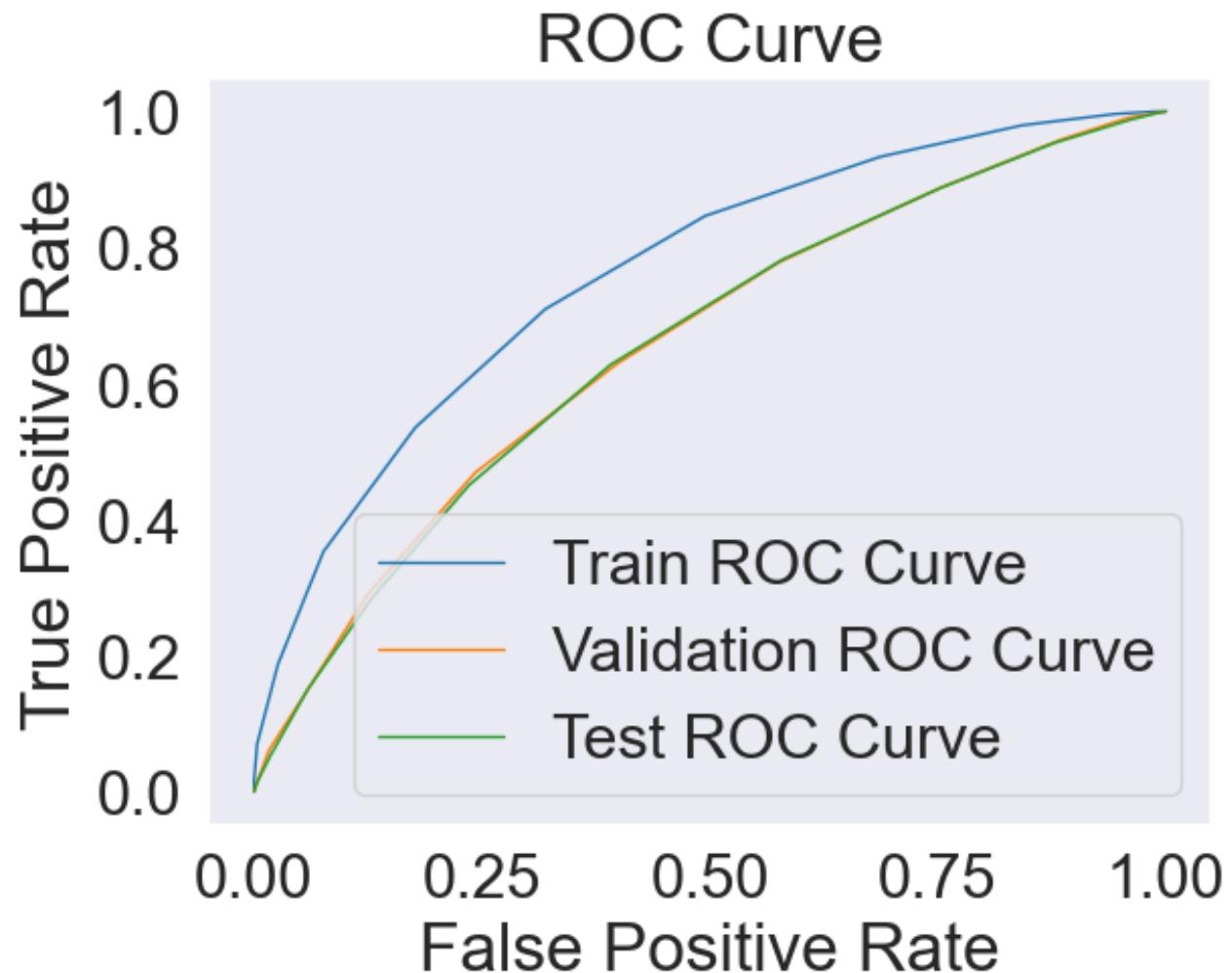
# Model Training and Validation
model.fit(X_train, y_train)

# Model Predictions
train_preds = model.predict_proba(X_train)[:, 1]
valid_preds = model.predict_proba(X_valid)[:, 1]
test_preds = model.predict_proba(X_test)[:, 1]

# Compute Metrics
train_fpr, train_tpr, _ = roc_curve(y_train, train_preds)
valid_fpr, valid_tpr, _ = roc_curve(y_valid, valid_preds)
test_fpr, test_tpr, _ = roc_curve(y_test, test_preds)
# Plot ROC Curve
plt.plot(train_fpr, train_tpr, label="Train ROC Curve")
plt.plot(valid_fpr, valid_tpr, label="Validation ROC Curve")
plt.plot(test_fpr, test_tpr, label="Test ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()
```

expLog

Total Features: 124 – Numerical: 107, Categorical: 16
Test score: 0.6183857210551256



Out[52]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC	Test AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379	0.7481
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4864	0.0465	0.6876	0.6843	0.6904	0.7525	0.7489	0.7535
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1799	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571	0.6550

Model 4 - SVM model

In [56]:

```
from sklearn.svm import SVC
from sklearn.pipeline import Pipeline
import time
import numpy as np

from sklearn.svm import SVC
from sklearn.pipeline import Pipeline
from sklearn.model_selection import GridSearchCV
import time
import numpy as np

# Create pipeline for preparing the data and select features
data_prep_pipeline, selected_features = get_pipeline()

# Join the preparation pipeline with SVM model
svm_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("svm", SVC(random_state=42, C=0.1, degree=1, kernel="poly", probability=True))
])

# Define a simple parameter grid
param_grid = {
    'svm_C': [0.1],
    'svm_degree': [1],
    'svm_kernel': ['poly'],
}

# Create a GridSearchCV object with n_jobs=-1 to use all available CPU core
grid_search = GridSearchCV(svm_full_pipeline_with_predictor, param_grid, n_jobs=-1)

# Train the model on the training set
start = time.time()
grid_search.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

# Compute the test score
start = time.time()
test_score = grid_search.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

print("Test score:", test_score)

# Results
exp_name = f"Model-4 SVM"
experiment_description = f"SVM with undersampled data-2 {len(selected_features)} features"
expLog = get_results(expLog, exp_name, experiment_description, model, train_time, test_time)
```

Total Features: 124 - Numerical: 107, Categorical: 16

```

/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
Test score: 0.6834623282261886

```

Out[56]:

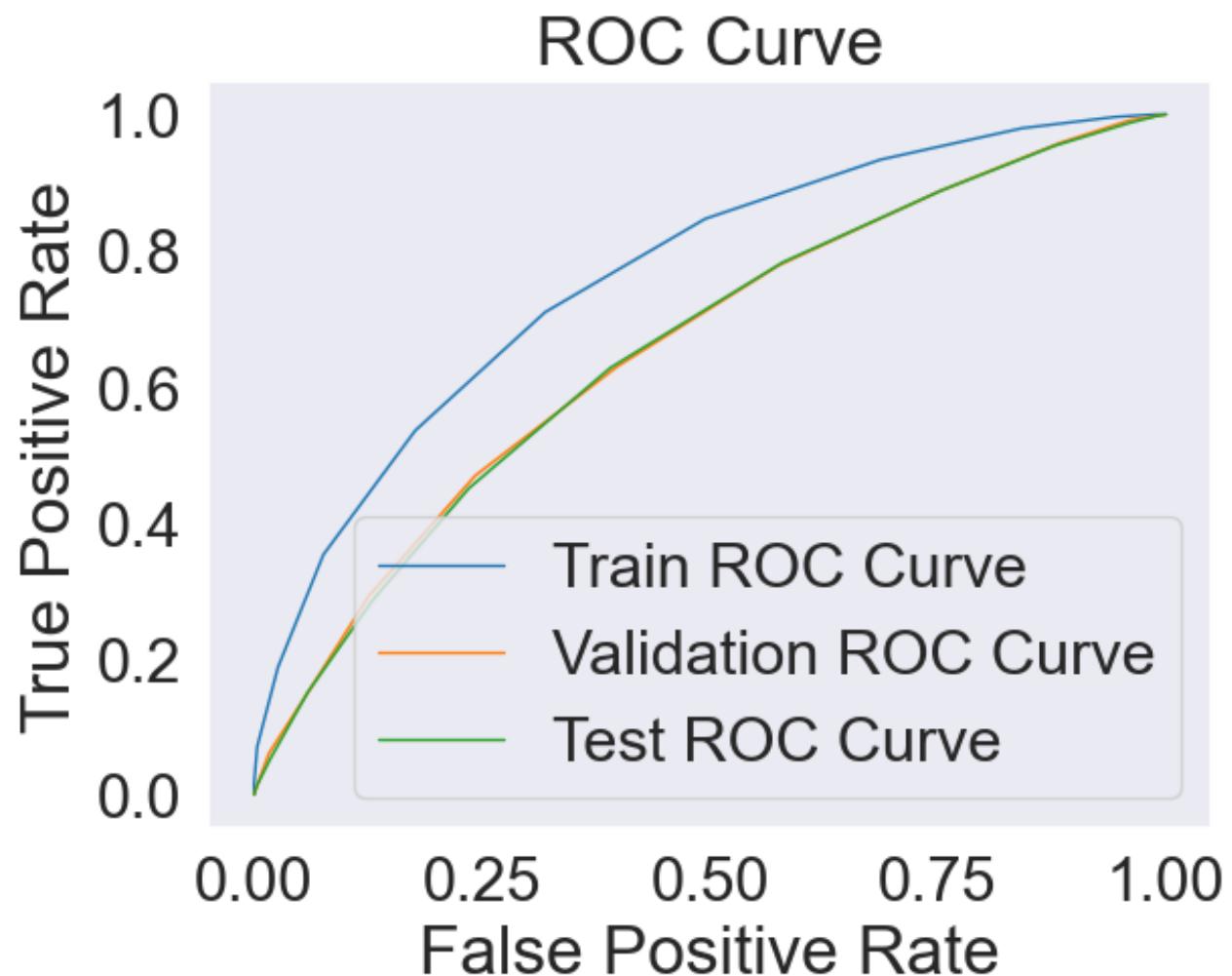
	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline	Logistic regression with LR undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379 0.7489
1	Model-2 Baseline	Logistic regression with LR undersampled data-2 1...	1.4864	0.0465	0.6876	0.6843	0.6904	0.7525	0.7489 0.7489
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1799	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571 0.6571
3	Model-4 SVM	SVM with undersampled data-2 124 features	1993.0196	15.9503	0.6950	0.6155	0.6184	0.7625	0.6571 0.6571

In [58]:

```
# Model Training and Validation
model.fit(X_train, y_train)

# Model Predictions
train_preds = model.predict_proba(X_train)[:, 1]
valid_preds = model.predict_proba(X_valid)[:, 1]
test_preds = model.predict_proba(X_test)[:, 1]

# Compute Metrics
train_fpr, train_tpr, _ = roc_curve(y_train, train_preds)
valid_fpr, valid_tpr, _ = roc_curve(y_valid, valid_preds)
test_fpr, test_tpr, _ = roc_curve(y_test, test_preds)
# Plot ROC Curve
plt.plot(train_fpr, train_tpr, label="Train ROC Curve")
plt.plot(valid_fpr, valid_tpr, label="Validation ROC Curve")
plt.plot(test_fpr, test_tpr, label="Test ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()
expLog
```



Out [58]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379 0.7
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4864	0.0465	0.6876	0.6843	0.6904	0.7525	0.7489 0.7
2	Model-3 KNN KNN	KNN with undersampled data-2 124 features	0.1799	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571 0.6
3	Model-4 SVM SVM	SVM with undersampled data-2 124 features	1993.0196	15.9503	0.6950	0.6155	0.6184	0.7625	0.6571 0.6

In [57]:

```
print("Done")
```

Done

Model 5-Decision Tree

In [59]:

```
from sklearn.tree import DecisionTreeClassifier

np.random.seed(42)

# Create pipeline for preparing the data and select features
data_prep_pipeline, selected_features = get_pipeline()

# Join the preparation pipeline with Decision Tree model
decision_tree_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("decision tree", DecisionTreeClassifier(random_state=42, criterion='en
])]

# Train the model on the training set
start = time.time()
model = decision_tree_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

# Compute the test score
```

```
start = time.time()
test_score = decision_tree_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Print the test score and other details
print(f"Test score: {test_score:.4f}")
print(f"Training time: {train_time} sec")
print(f"Test time: {test_time} sec")
#print(f"Selected features: {selected_features}")
print(f"Number of features: {len(selected_features)}")

# Results
# Save the experiment results
exp_name = f"Model-5 Decision Tree"
experiment_description = f"Decision tree with undersampled data-2 {len(selected_features)} features"
expLog = get_results(expLog, exp_name, experiment_description, model, train_time, test_time, test_accuracy)

# Model Training and Validation
model.fit(X_train, y_train)

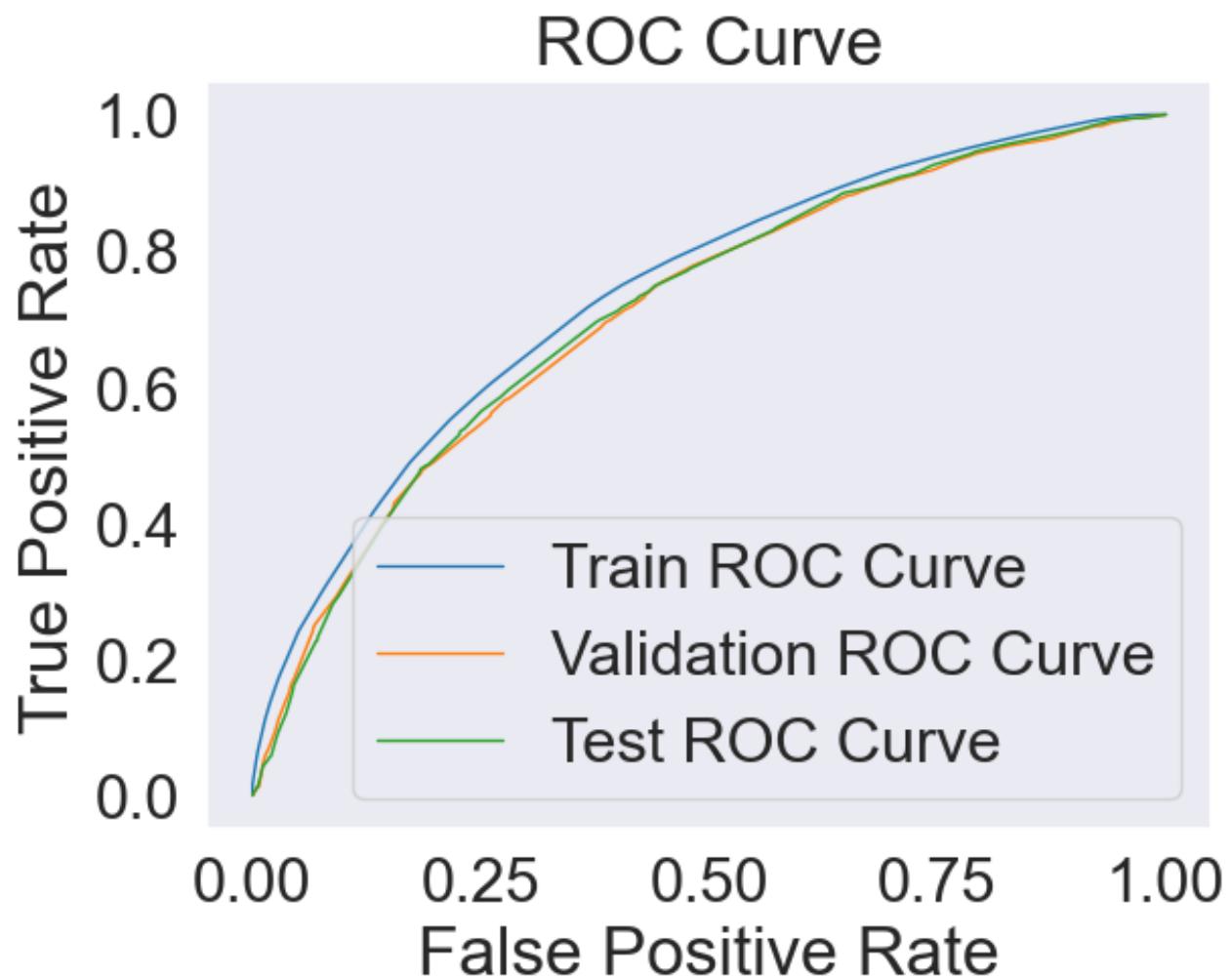
# Model Predictions
train_preds = model.predict_proba(X_train)[:, 1]
valid_preds = model.predict_proba(X_valid)[:, 1]
test_preds = model.predict_proba(X_test)[:, 1]

# Compute Metrics
train_fpr, train_tpr, _ = roc_curve(y_train, train_preds)
valid_fpr, valid_tpr, _ = roc_curve(y_valid, valid_preds)
test_fpr, test_tpr, _ = roc_curve(y_test, test_preds)

# Plot ROC Curve
plt.plot(train_fpr, train_tpr, label="Train ROC Curve")
plt.plot(valid_fpr, valid_tpr, label="Validation ROC Curve")
plt.plot(test_fpr, test_tpr, label="Test ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()

expLog
```

Total Features: 124 - Numerical: 107, Categorical: 16
Test score: 0.6591
Training time: 0.8042 sec
Test time: 0.0229 sec
Number of features: 124



Out[59]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379 0.7
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 124...	1.4864	0.0465	0.6876	0.6843	0.6904	0.7525	0.7489 0.7
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1799	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571 0.6
3	Model-4 SVM	SVM with undersampled data-2 124 features	1993.0196	15.9503	0.6950	0.6155	0.6184	0.7625	0.6571 0.6
4	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.8042	0.0229	0.6749	0.6535	0.6591	0.7380	0.7105 0.7

Model 6- Random Forest

In [60]:

```

from sklearn.ensemble import RandomForestClassifier

np.random.seed(42)

# Creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline()

# Attaching random forest model to the above pipeline
random_forest_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("random forest", RandomForestClassifier(random_state=42, bootstrap=True,
                                             max_features=5, min_samples_leaf=10, min_samples_split=2))
])

# Training the model
print("Training the model...")
start_time = time.time()
model = random_forest_full_pipeline_with_predictor.fit(X_train, y_train)

```

```
train_time = np.round(time.time() - start_time, 4)
print("Training time: ", train_time)

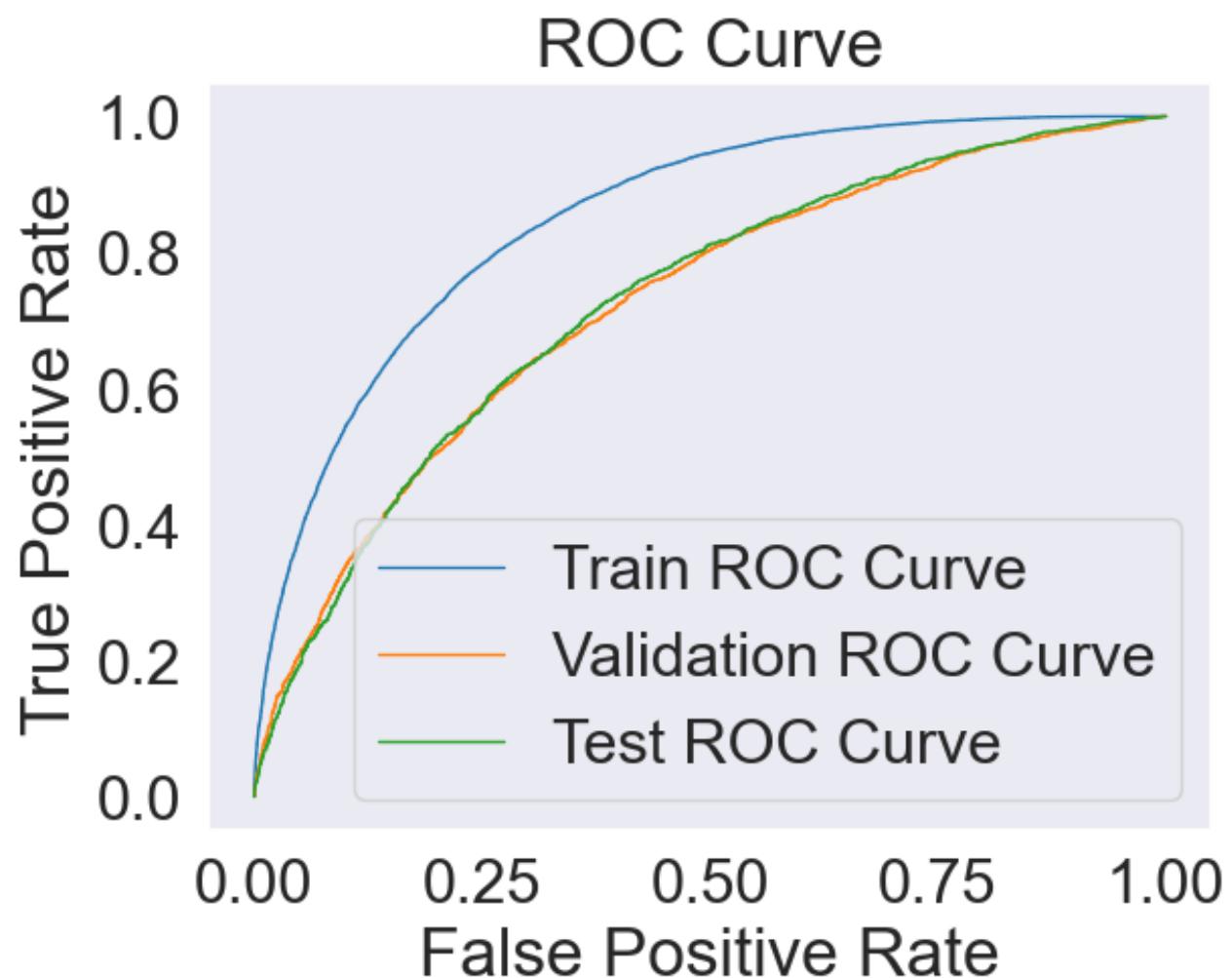
# Evaluate the model on the test set and measure the test time
print("Evaluating the model on the test set...")
start_time = time.time()
score_test = random_forest_full_pipeline_with_predictor.score(X_test, y_te
test_time = np.round(time.time() - start_time, 4)
print("Test score: ", score_test)
print("Test time: ", test_time)

# Results
exp_name = f"Model-6 Random Forest"
experiment_description = f"Random Forest with undersampled data-2 {len(selec
expLog = get_results(expLog, exp_name, experiment_description, model, train_t
expLog
# Model Training and Validation
model.fit(X_train, y_train)

# Model Predictions
train_preds = model.predict_proba(X_train)[:, 1]
valid_preds = model.predict_proba(X_valid)[:, 1]
test_preds = model.predict_proba(X_test)[:, 1]

# Compute Metrics
train_fpr, train_tpr, _ = roc_curve(y_train, train_preds)
valid_fpr, valid_tpr, _ = roc_curve(y_valid, valid_preds)
test_fpr, test_tpr, _ = roc_curve(y_test, test_preds)
# Plot ROC Curve
plt.plot(train_fpr, train_tpr, label="Train ROC Curve")
plt.plot(valid_fpr, valid_tpr, label="Validation ROC Curve")
plt.plot(test_fpr, test_tpr, label="Test ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()
expLog
```

Total Features: 124 - Numerical: 107, Categorical: 16
Training the model...
Training time: 10.1732
Evaluating the model on the test set...
Test score: 0.6660875059232348
Test time: 0.306



Out[60]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379 0.
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4864	0.0465	0.6876	0.6843	0.6904	0.7525	0.7489 0.
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1799	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571 0.
3	Model-4 SVM	SVM with undersampled data-2 124 features	1993.0196	15.9503	0.6950	0.6155	0.6184	0.7625	0.6571 0.
4	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.8042	0.0229	0.6749	0.6535	0.6591	0.7380	0.7105 0.
5	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.1732	0.3060	0.7664	0.6654	0.6661	0.8504	0.7243 0.

Model 7- Extra Trees

In [61]:

```

from sklearn.ensemble import ExtraTreesClassifier

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline()

# Attaching Extra-Trees model to the above pipeline
extra_trees_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("extra trees", ExtraTreesClassifier(random_state=42, bootstrap=True, m
                                         max_features=5, min_samples_leaf=10, min_samples_spli

```

```
])

# Training the model
start = time.time()
model = extra_trees_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = extra_trees_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

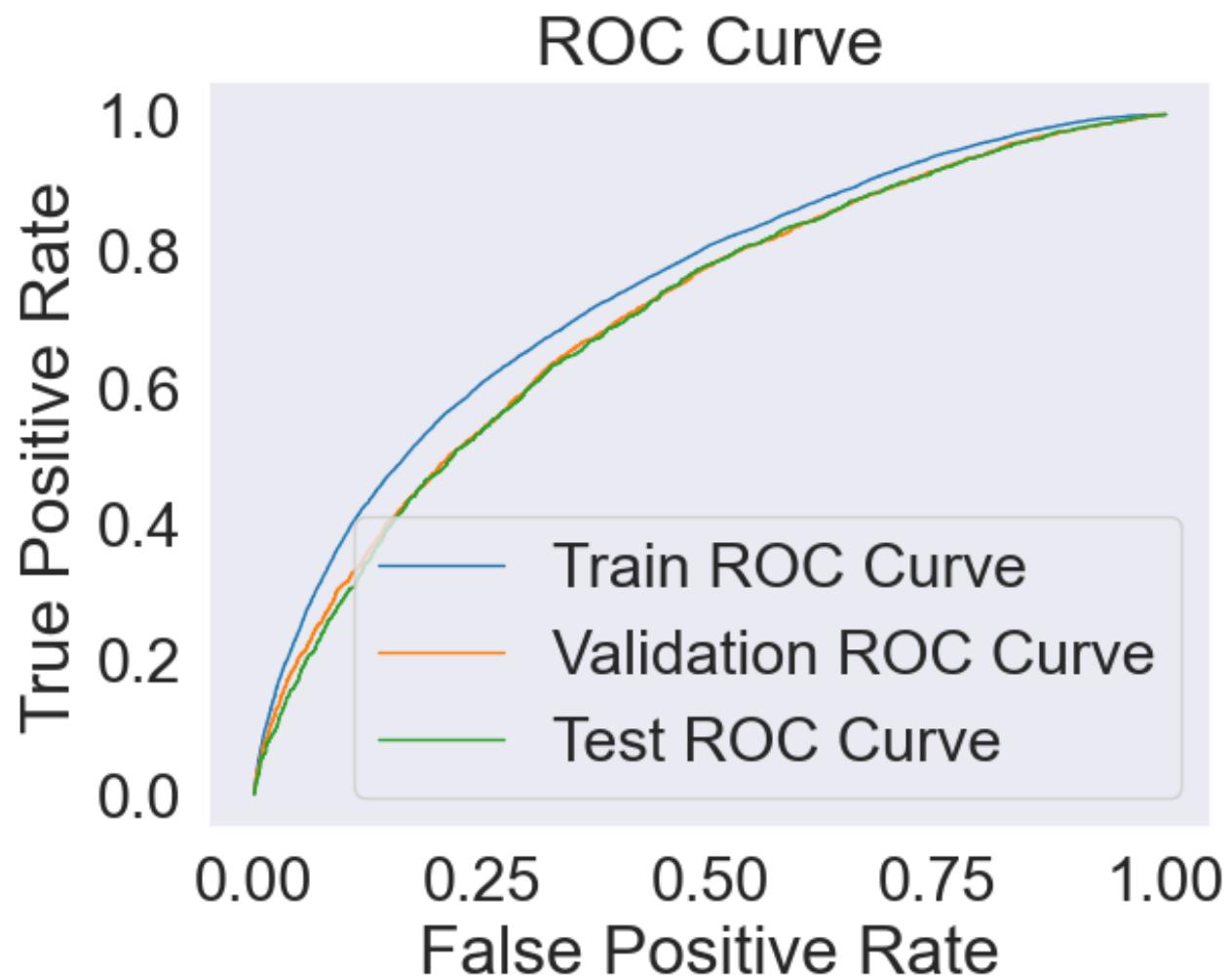
# Results
exp_name = f"Model-7 Extra Trees "
experiment_description = f"Extra Trees with undersampled data-2 {len(selecte
expLog = get_results(expLog, exp_name, experiment_description, model, train_t
expLog

# Model Training and Validation
model.fit(X_train, y_train)

# Model Predictions
train_preds = model.predict_proba(X_train)[:, 1]
valid_preds = model.predict_proba(X_valid)[:, 1]
test_preds = model.predict_proba(X_test)[:, 1]

# Compute Metrics
train_fpr, train_tpr, _ = roc_curve(y_train, train_preds)
valid_fpr, valid_tpr, _ = roc_curve(y_valid, valid_preds)
test_fpr, test_tpr, _ = roc_curve(y_test, test_preds)
# Plot ROC Curve
plt.plot(train_fpr, train_tpr, label="Train ROC Curve")
plt.plot(valid_fpr, valid_tpr, label="Validation ROC Curve")
plt.plot(test_fpr, test_tpr, label="Test ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()
expLog
```

Total Features: 124 - Numerical: 107, Categorical: 16



Out[61]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379 0.
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4864	0.0465	0.6876	0.6843	0.6904	0.7525	0.7489 0.
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1799	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571 0.
3	Model-4 SVM	SVM with undersampled data-2 124 features	1993.0196	15.9503	0.6950	0.6155	0.6184	0.7625	0.6571 0.
4	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.8042	0.0229	0.6749	0.6535	0.6591	0.7380	0.7105 0.
5	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.1732	0.3060	0.7664	0.6654	0.6661	0.8504	0.7243 0.
6	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features	1.3084	0.0857	0.6738	0.6502	0.6482	0.7405	0.7059 0.

Model 8-Bagging Meta Estimator

In [62]:

```
from sklearn.ensemble import BaggingClassifier

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline()
```

```

# Attaching bagging meta-estimator to the above pipeline
bagging_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("bagging", BaggingClassifier(base_estimator=None, n_estimators=10, max_
        bootstrap=True, bootstrap_features=False
        verbose=0, warm_start=False))
])

# Training the model
start = time.time()
model = bagging_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = bagging_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
exp_name = f"Model-8 Bagging Meta Estimator"
experiment_description =f"Bagging Meta Estimator with undersampled data-2 {expLog = get_results(expLog, exp_name, experiment_description, model, train_t
expLog

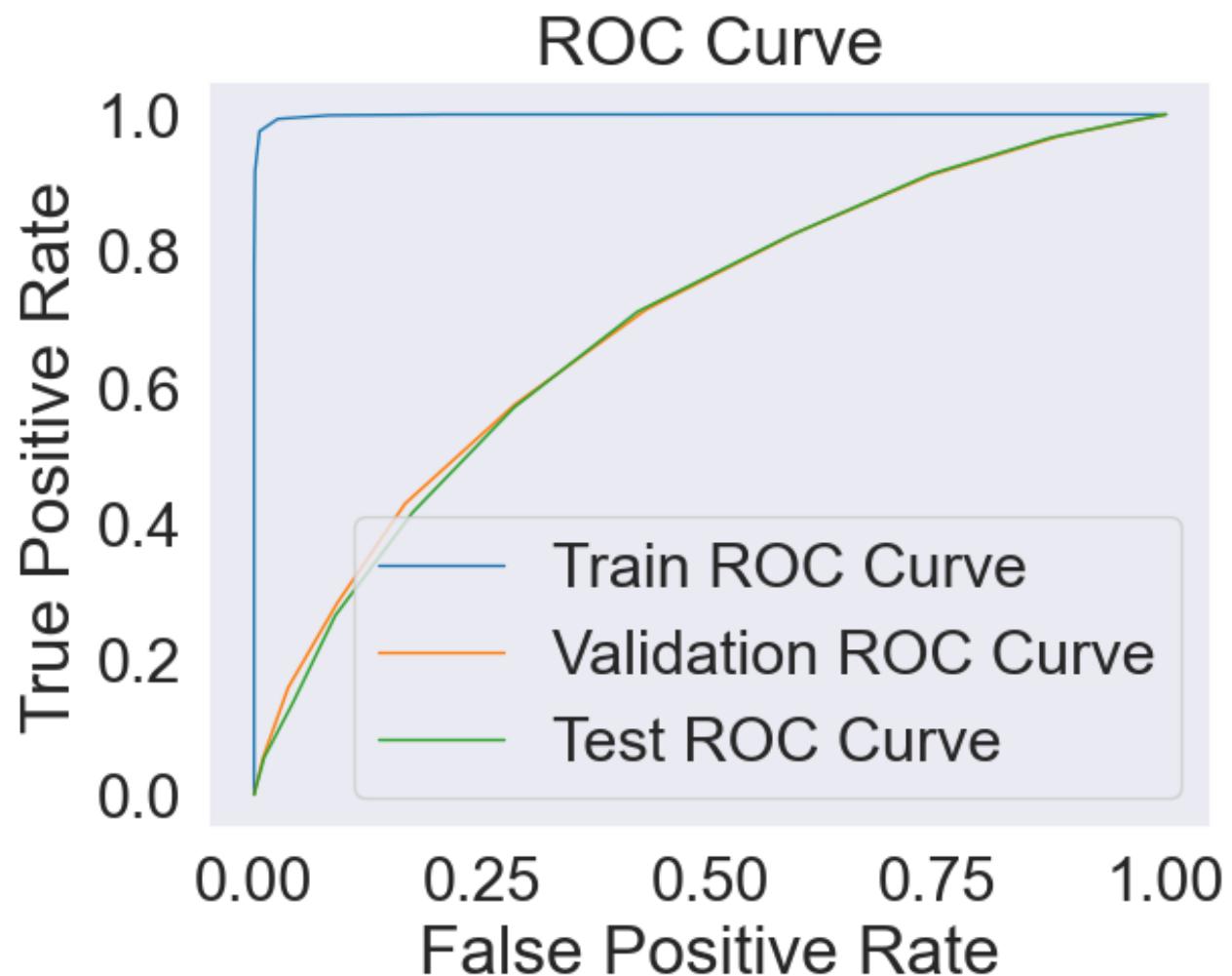
# Model Training and Validation
model.fit(X_train, y_train)

# Model Predictions
train_preds = model.predict_proba(X_train)[:, 1]
valid_preds = model.predict_proba(X_valid)[:, 1]
test_preds = model.predict_proba(X_test)[:, 1]

# Compute Metrics
train_fpr, train_tpr, _ = roc_curve(y_train, train_preds)
valid_fpr, valid_tpr, _ = roc_curve(y_valid, valid_preds)
test_fpr, test_tpr, _ = roc_curve(y_test, test_preds)
# Plot ROC Curve
plt.plot(train_fpr, train_tpr, label="Train ROC Curve")
plt.plot(valid_fpr, valid_tpr, label="Validation ROC Curve")
plt.plot(test_fpr, test_tpr, label="Test ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()
expLog

```

Total Features: 124 - Numerical: 107, Categorical: 16



Out[62]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379 0.
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4864	0.0465	0.6876	0.6843	0.6904	0.7525	0.7489 0.
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1799	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571 0.
3	Model-4 SVM	SVM with undersampled data-2 124 features	1993.0196	15.9503	0.6950	0.6155	0.6184	0.7625	0.6571 0.
4	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.8042	0.0229	0.6749	0.6535	0.6591	0.7380	0.7105 0.
5	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.1732	0.3060	0.7664	0.6654	0.6661	0.8504	0.7243 0.
6	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features	1.3084	0.0857	0.6738	0.6502	0.6482	0.7405	0.7059 0.
7	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...	2.5579	0.1725	0.9844	0.6477	0.6446	0.9990	0.6978 0.

Model 9 - ADABOOST SAMME

In [63]:

```
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import BaggingClassifier
from sklearn.ensemble import AdaBoostClassifier

# Define data preparation pipeline
data_prep_pipeline, selected_features = get_pipeline()

# Define base estimator
base_estimator = DecisionTreeClassifier(max_depth=5)

# Define AdaBoost classifier with SAMME using Bagging and enable paralleliz
adaboost_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("Bagging", BaggingClassifier(base_estimator=AdaBoostClassifier(base_es
                                random_state=None), n_jobs=-1))
])

# Train AdaBoost classifier and log results
start = time.time()
model = adaboost_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = adaboost_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
exp_name = f"Model-9 ADABoost SAMME "
experiment_description = f"ADABoost SAMME with undersampled data-2 {len(sele
expLog = get_results(expLog, exp_name, experiment_description, model, train_t
#expLog
```

```
Total Features: 124 - Numerical: 107, Categorical: 16
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio
n 1.2 and will be removed in 1.4.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio
n 1.2 and will be removed in 1.4.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio
n 1.2 and will be removed in 1.4.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio
n 1.2 and will be removed in 1.4.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio
n 1.2 and will be removed in 1.4.
```

```
n 1.2 and will be removed in 1.4.  
    warnings.warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.  
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio  
n 1.2 and will be removed in 1.4.  
    warnings.warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.  
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio  
n 1.2 and will be removed in 1.4.  
    warnings.warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.  
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio  
n 1.2 and will be removed in 1.4.  
    warnings.warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.  
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio  
n 1.2 and will be removed in 1.4.  
    warnings.warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.  
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio  
n 1.2 and will be removed in 1.4.  
    warnings.warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh  
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh  
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh  
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh  
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh  
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh  
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh  
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh  
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh  
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh  
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh  
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```

In [64]:

```
# Model Training and Validation
model.fit(X_train, y_train)

# Model Predictions
train_preds = model.predict_proba(X_train)[:, 1]
valid_preds = model.predict_proba(X_valid)[:, 1]
test_preds = model.predict_proba(X_test)[:, 1]

# Compute Metrics
train_fpr, train_tpr, _ = roc_curve(y_train, train_preds)
valid_fpr, valid_tpr, _ = roc_curve(y_valid, valid_preds)
test_fpr, test_tpr, _ = roc_curve(y_test, test_preds)
# Plot ROC Curve
plt.plot(train_fpr, train_tpr, label="Train ROC Curve")
plt.plot(valid_fpr, valid_tpr, label="Validation ROC Curve")
plt.plot(test_fpr, test_tpr, label="Test ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()
expLog
```

```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio
n 1.2 and will be removed in 1.4.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio
n 1.2 and will be removed in 1.4.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio
n 1.2 and will be removed in 1.4.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio
n 1.2 and will be removed in 1.4.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.
py:166: FutureWarning: `base_estimator` was renamed to `estimator` in versio
n 1.2 and will be removed in 1.4.
    warnings.warn(
```

```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.py:166: FutureWarning: `base_estimator` was renamed to `estimator` in version 1.2 and will be removed in 1.4.  
    warnings.warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.py:166: FutureWarning: `base_estimator` was renamed to `estimator` in version 1.2 and will be removed in 1.4.  
    warnings.warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.py:166: FutureWarning: `base_estimator` was renamed to `estimator` in version 1.2 and will be removed in 1.4.  
    warnings.warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.py:166: FutureWarning: `base_estimator` was renamed to `estimator` in version 1.2 and will be removed in 1.4.  
    warnings.warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_base.py:166: FutureWarning: `base_estimator` was renamed to `estimator` in version 1.2 and will be removed in 1.4.  
    warnings.warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weight_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weight_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weight_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weight_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weight_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weight_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weight_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weight_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weight_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weight_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)
```

```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



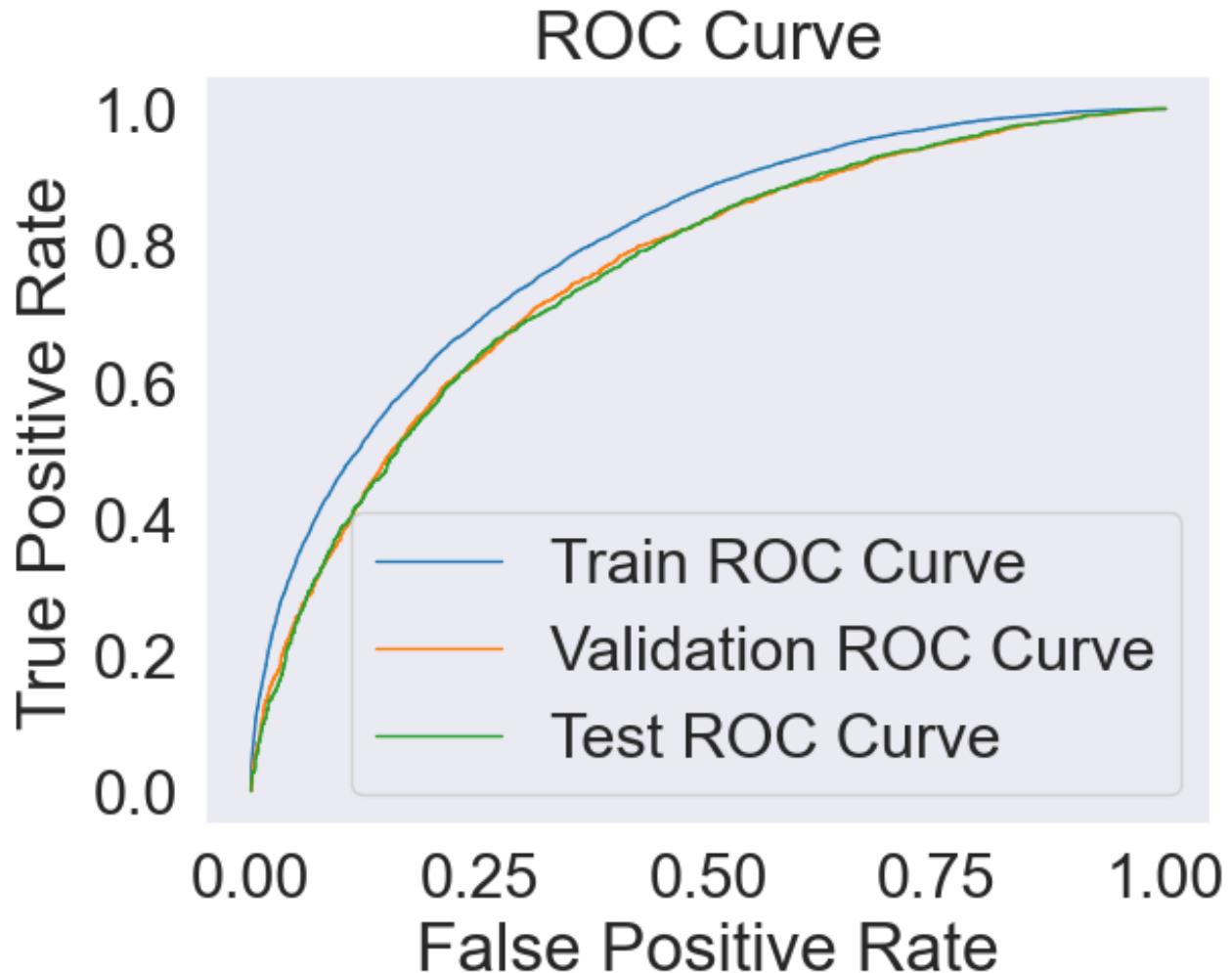
```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log
    np.log(sample_weight)
```



```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh  
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_weigh  
t_boosting.py:677: RuntimeWarning: divide by zero encountered in log  
    np.log(sample_weight)
```



Out[64]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379 0.
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4864	0.0465	0.6876	0.6843	0.6904	0.7525	0.7489 0.
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1799	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571 0.
3	Model-4 SVM	SVM with undersampled data-2 124 features	1993.0196	15.9503	0.6950	0.6155	0.6184	0.7625	0.6571 0.
4	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.8042	0.0229	0.6749	0.6535	0.6591	0.7380	0.7105 0.
5	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.1732	0.3060	0.7664	0.6654	0.6661	0.8504	0.7243 0.
6	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features	1.3084	0.0857	0.6738	0.6502	0.6482	0.7405	0.7059 0.
7	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...	2.5579	0.1725	0.9844	0.6477	0.6446	0.9990	0.6978 0.
8	Model-9 ADABoost SAMME	ADABoost SAMME with undersampled data-2 124 fe...	21.9628	0.3113	0.7185	0.6966	0.6950	0.7989	0.7580 0.

In [65]:

expLog

Out[65]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379 0.
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4864	0.0465	0.6876	0.6843	0.6904	0.7525	0.7489 0.
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1799	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571 0.
3	Model-4 SVM	SVM with undersampled data-2 124 features	1993.0196	15.9503	0.6950	0.6155	0.6184	0.7625	0.6571 0.
4	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.8042	0.0229	0.6749	0.6535	0.6591	0.7380	0.7105 0.
5	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.1732	0.3060	0.7664	0.6654	0.6661	0.8504	0.7243 0.
6	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features	1.3084	0.0857	0.6738	0.6502	0.6482	0.7405	0.7059 0.
7	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...	2.5579	0.1725	0.9844	0.6477	0.6446	0.9990	0.6978 0.
8	Model-9 ADABoost SAMME	ADABoost SAMME with undersampled data-2 124 fe...	21.9628	0.3113	0.7185	0.6966	0.6950	0.7989	0.7580 0.

Model 10 - XGBoost

In [66]:

```
np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline()

# Attaching XGBoost model to the above pipeline
xgboost_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("xgboost", XGBClassifier(random_state=42,
                               objective='binary:logistic', max_depth=5, eta=0.001,
                               learning_rate=0.01, colsample_bytree=0.7, n_estimators=1000
    ))
]

# Training the model
start = time.time()
model = xgboost_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = xgboost_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

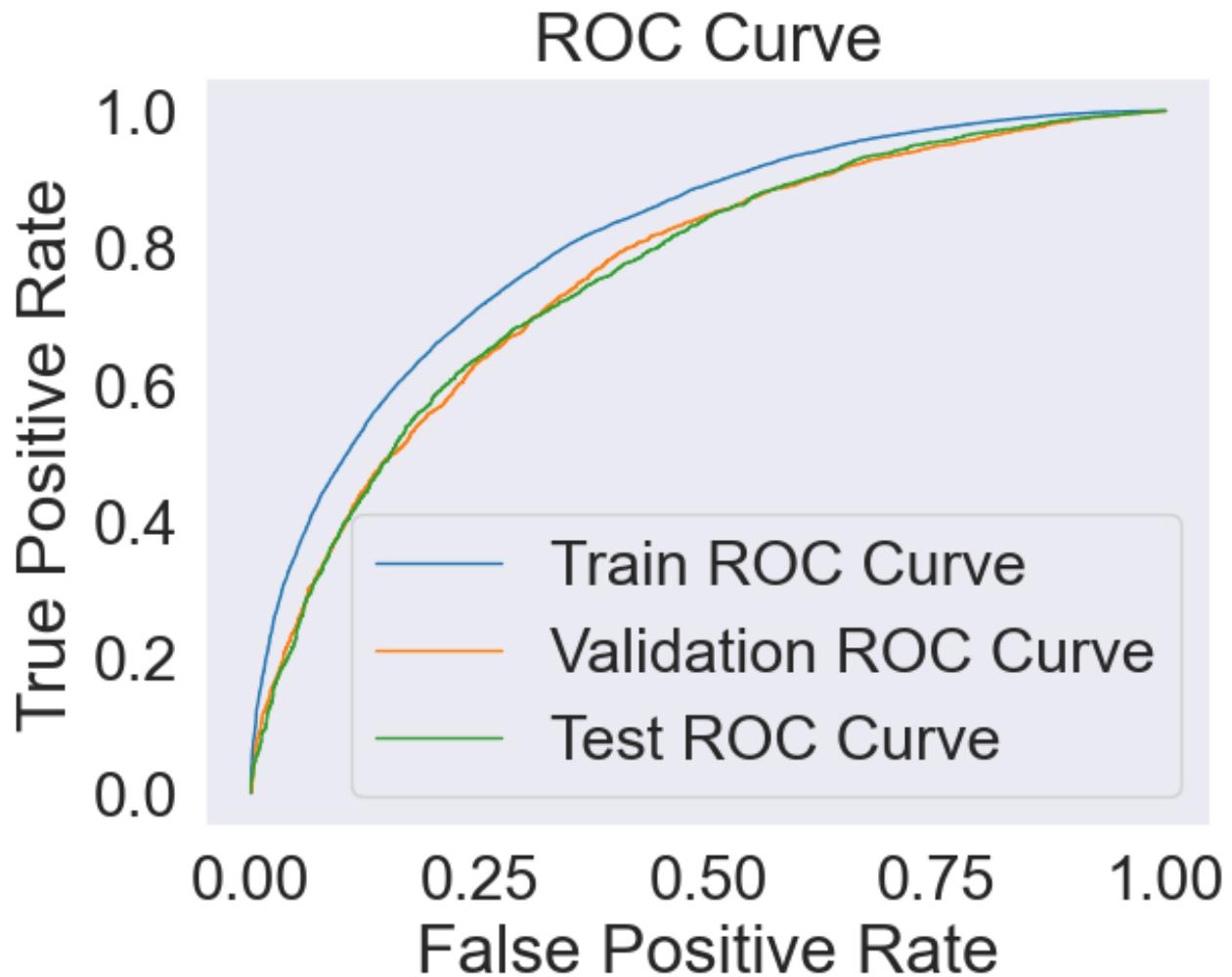
# Results
exp_name = f"Model-10 XGBoost"
experiment_description = f"XGBoost SAMME with undersampled data-2 {len(selected_features)} features"
expLog = get_results(expLog, exp_name, experiment_description, model, train_time, test_time)
expLog

# Model Training and Validation
model.fit(X_train, y_train)

# Model Predictions
train_preds = model.predict_proba(X_train)[:, 1]
valid_preds = model.predict_proba(X_valid)[:, 1]
test_preds = model.predict_proba(X_test)[:, 1]

# Compute Metrics
train_fpr, train_tpr, _ = roc_curve(y_train, train_preds)
valid_fpr, valid_tpr, _ = roc_curve(y_valid, valid_preds)
test_fpr, test_tpr, _ = roc_curve(y_test, test_preds)
# Plot ROC Curve
plt.plot(train_fpr, train_tpr, label="Train ROC Curve")
plt.plot(valid_fpr, valid_tpr, label="Validation ROC Curve")
plt.plot(test_fpr, test_tpr, label="Test ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()
expLog
```

Total Features: 124 – Numerical: 107, Categorical: 16



Out[66]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4864	0.0465	0.6876	0.6843	0.6904	0.7525	0.7489
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1799	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571
3	Model-4	SVM with undersampled	1993.0196	15.9503	0.6950	0.6155	0.6184	0.7625	0.6571

	SVM	data-2 124 features									
4	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.8042	0.0229	0.6749	0.6535	0.6591	0.7380	0.7105	0.	
5	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.1732	0.3060	0.7664	0.6654	0.6661	0.8504	0.7243	0.	
6	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features	1.3084	0.0857	0.6738	0.6502	0.6482	0.7405	0.7059	0.	
7	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...	2.5579	0.1725	0.9844	0.6477	0.6446	0.9990	0.6978	0.	
8	Model-9 ADABOOST SAMMME	ADABOOST SAMMME with undersampled data-2 124 fe...	21.9628	0.3113	0.7185	0.6966	0.6950	0.7989	0.7580	0.	
9	Model-10 XGBoost	XGBoost SAMMME with undersampled data-2 124 fea...	24.6443	0.0366	0.7312	0.6937	0.6940	0.8116	0.7612	0	

Model 11-ADABOOST

In [67]:

```

from sklearn.ensemble import AdaBoostClassifier, BaggingClassifier
from sklearn.utils import parallel_backend
from joblib import parallel_backend

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline()

# Define base estimator
base_estimator = AdaBoostClassifier(random_state=42, n_estimators=1000, lea

```

```

# Enable parallelization with Bagging
with parallel_backend('multiprocessing', n_jobs=10):
    adaboost_full_pipeline_with_predictor = Pipeline([
        ("preparation", data_prep_pipeline),
        ("Bagging", BaggingClassifier(base_estimator=base_estimator, n_jobs=1))
    ])

# Training the model
start = time.time()
model = adaboost_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = adaboost_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
exp_name = f"Model-11 ADABoost"
experiment_description = f"ADABoost with undersampled data-2 {len(selected_f_exLog = get_results(expLog, exp_name, experiment_description, model, train_t_exLog

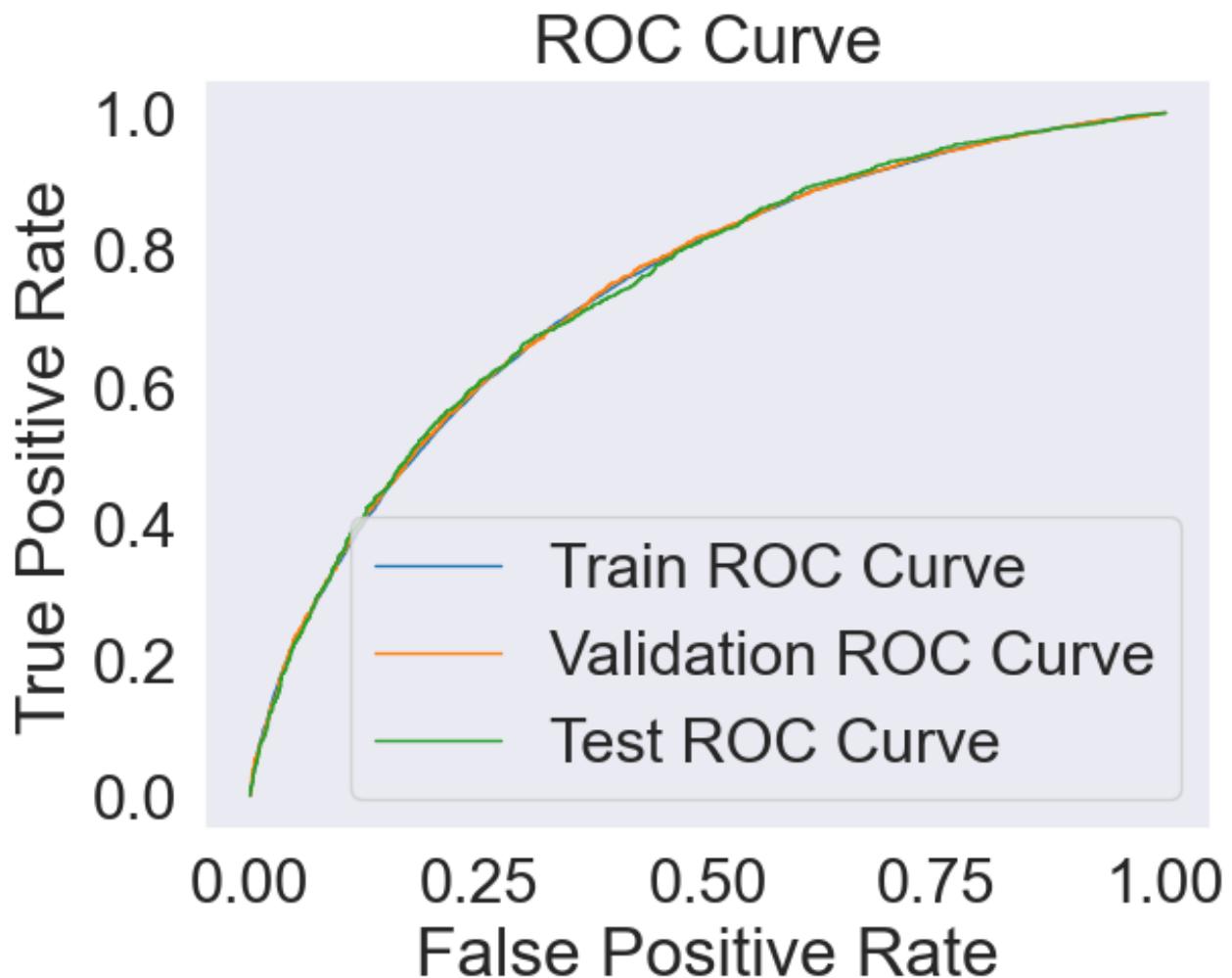
# Model Training and Validation
model.fit(X_train, y_train)

# Model Predictions
train_preds = model.predict_proba(X_train)[:, 1]
valid_preds = model.predict_proba(X_valid)[:, 1]
test_preds = model.predict_proba(X_test)[:, 1]

# Compute Metrics
train_fpr, train_tpr, _ = roc_curve(y_train, train_preds)
valid_fpr, valid_tpr, _ = roc_curve(y_valid, valid_preds)
test_fpr, test_tpr, _ = roc_curve(y_test, test_preds)
# Plot ROC Curve
plt.plot(train_fpr, train_tpr, label="Train ROC Curve")
plt.plot(valid_fpr, valid_tpr, label="Validation ROC Curve")
plt.plot(test_fpr, test_tpr, label="Test ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()
expLog

```

Total Features: 124 – Numerical: 107, Categorical: 16



Out[67]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4864	0.0465	0.6876	0.6843	0.6904	0.7525	0.7489
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1799	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571
3	Model-4 SVM	SVM with undersampled data-2 124	1993.0196	15.9503	0.6950	0.6155	0.6184	0.7625	0.6571

			features												
4	Model-5	Decision Tree	Decision tree with undersampled data-2 124 fe...	0.8042	0.0229	0.6749	0.6535	0.6591	0.7380	0.7105	C				
5	Model-6	Random Forest	Random Forest with undersampled data-2 124 fe...	10.1732	0.3060	0.7664	0.6654	0.6661	0.8504	0.7243	C				
6	Model-7	Extra Trees	Extra Trees with undersampled data-2 124 features	1.3084	0.0857	0.6738	0.6502	0.6482	0.7405	0.7059	C				
7	Model-8	Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...	2.5579	0.1725	0.9844	0.6477	0.6446	0.9990	0.6978	0				
8	Model-9	ADABOOST SAMMME	ADABOOST SAMMME with undersampled data-2 124 fe...	21.9628	0.3113	0.7185	0.6966	0.6950	0.7989	0.7580	0				
9	Model-10	XGBoost	XGBoost SAMMME with undersampled data-2 124 fe...	24.6443	0.0366	0.7312	0.6937	0.6940	0.8116	0.7612	I				
10	Model-11	ADABOOST	ADABOOST SAMMME with undersampled data-2 124 fe...	120.2588	2.8816	0.6774	0.6776	0.6771	0.7388	0.7402	0				

Model 12 - CatBoost

In [68]:

```
from catboost import CatBoostClassifier
from sklearn.pipeline import Pipeline
import numpy as np
import time

np.random.seed(42)
```

```

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline()

# Attaching CatBoost model to the above pipeline
catboost_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("catboost", CatBoostClassifier(random_state=42, iterations=1000, learn_depth=5, thread_count=-1, verbose=False))
])

# Training the model
start = time.time()
model = catboost_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = catboost_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

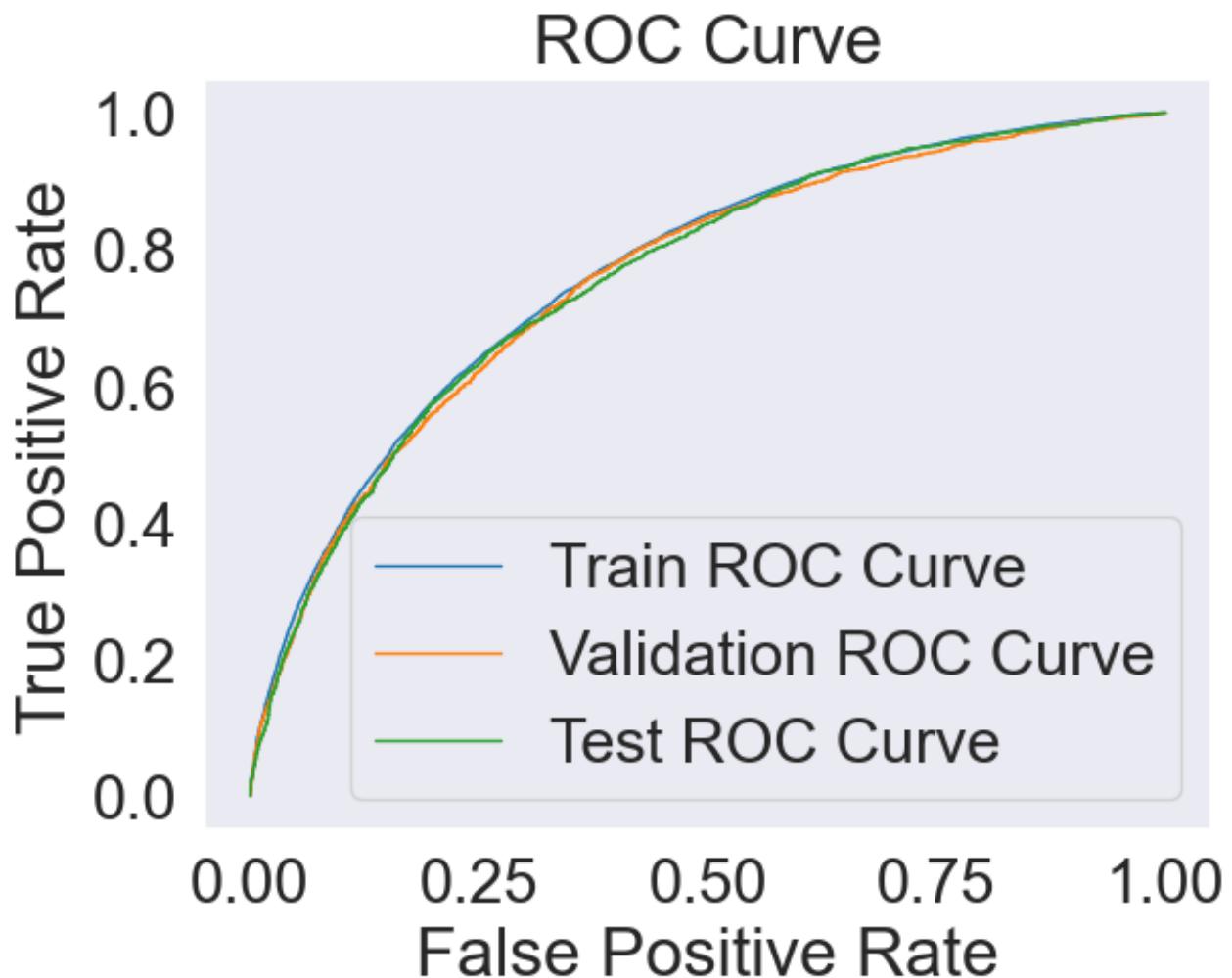
# Results
exp_name = f"Model-12 CATBoost"
experiment_description = f"CATBoost with undersampled data-2 {len(selected_f
expLog = get_results(expLog, exp_name, experiment_description, model, train_t
expLog
# Model Training and Validation
model.fit(X_train, y_train)

# Model Predictions
train_preds = model.predict_proba(X_train)[:, 1]
valid_preds = model.predict_proba(X_valid)[:, 1]
test_preds = model.predict_proba(X_test)[:, 1]

# Compute Metrics
train_fpr, train_tpr, _ = roc_curve(y_train, train_preds)
valid_fpr, valid_tpr, _ = roc_curve(y_valid, valid_preds)
test_fpr, test_tpr, _ = roc_curve(y_test, test_preds)
# Plot ROC Curve
plt.plot(train_fpr, train_tpr, label="Train ROC Curve")
plt.plot(valid_fpr, valid_tpr, label="Validation ROC Curve")
plt.plot(test_fpr, test_tpr, label="Test ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()
expLog

```

Total Features: 124 - Numerical: 107, Categorical: 16



Out[68]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.2556	0.0774	0.7732	0.7629	0.7763	0.7536	0.7379
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4864	0.0465	0.6876	0.6843	0.6904	0.7525	0.7489
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1799	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571
3	Model-4 SVM	SVM with undersampled data-2 124	1993.0196	15.9503	0.6950	0.6155	0.6184	0.7625	0.6571

			features									
4	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fe...	0.8042	0.0229	0.6749	0.6535	0.6591	0.7380	0.7105	C		
5	Model-6 Random Forest	Random Forest with undersampled data-2 124 fe...	10.1732	0.3060	0.7664	0.6654	0.6661	0.8504	0.7243	C		
6	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features	1.3084	0.0857	0.6738	0.6502	0.6482	0.7405	0.7059	C		
7	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...	2.5579	0.1725	0.9844	0.6477	0.6446	0.9990	0.6978	0		
8	Model-9 ADABOOST SAMMME	ADABOOST SAMMME with undersampled data-2 124 fe...	21.9628	0.3113	0.7185	0.6966	0.6950	0.7989	0.7580	0		
9	Model-10 XGBoost	XGBoost SAMMME with undersampled data-2 124 fe...	24.6443	0.0366	0.7312	0.6937	0.6940	0.8116	0.7612	I		
10	Model-11 ADABOOST	ADABOOST SAMMME with undersampled data-2 124 fe...	120.2588	2.8816	0.6774	0.6776	0.6771	0.7388	0.7402	0		
11	Model-12 CATBoost	CATBoost SAMMME with undersampled data-2 124 fe...	6.5636	0.2012	0.6964	0.6904	0.6926	0.7671	0.7580	C		

Selecting 3 best learners , hyper parameter tuning , and ensemble learner

HYPERPARAMETER TUNING AND FEATURE SELECTION

XGBOOST Hyper parameter Tuning and feature selection

In [72]:

```
def pct(x):
    return round(100*x,3)
```

In [120...]

```
results = pd.DataFrame(columns=["ExpID", "Cross-fold Train Accuracy", "Test
features_dict = dict()

# A Function to execute the grid search and record the results.
def ConductGridSearch(X_train, y_train, X_test, y_test):
    # classifier for our grid search experiment
    classifiers = [
        ('DecisionTrees', DecisionTreeClassifier(random_state=42)),
        ('XGBoost', XGBClassifier(random_state=42))
    ]

    # grid search parameters for the classifier
    param_grid = {

        'XGBoost': {
            'max_depth': [5,9], # Lower helps with overfitting
            'n_estimators':[800, 1000],
            'learning_rate': [0.001, 0.01],
            'eta' : [0.001, 0.01],
            'colsample_bytree' : [0.5, 0.7],
        }
    }

    # grid search parameters for the classifier
    # param_grid = {

    #     'XGBoost': {
    #         'max_depth': [5], # Lower helps with overfitting
    #         'n_estimators':[20],
    #         'learning_rate': [ 0.1],
    #         'eta' : [0.1],
    #         'colsample_bytree' : [0.5],
    #     }
    # }

    for (name, classifier) in classifiers:

        print('***** STARTING TUNING', name, '*****')
        parameters = param_grid[name]
        print("Parameters:")
```

```
for p in sorted(parameters.keys()):
    print("\t"+str(p)+": "+ str(parameters[p]))\n\n# generate the pipeline
full_pipeline_with_predictor = Pipeline([
    ("preparation", FeatureUnion(transformer_list=[("num_pipeline",
    ("predictor", classifier)
])\n\n# Execute the grid search
params = {}
for p in parameters.keys():
    pipe_key = 'predictor__'+str(p)
    params[pipe_key] = parameters[p]\n\ngrid_search = GridSearchCV(full_pipeline_with_predictor, params, sc
                                n_jobs=-1, verbose=1)
grid_search.fit(X_train, y_train)\n\n# Best estimator training time
start = time.time()
grid_search.best_estimator_.fit(X_train, y_train)
train_time = round(time.time() - start, 4)\n\n# Training accuracy
cvSplits = ShuffleSplit(n_splits=3, test_size=0.7, random_state=42)
best_train_scores = cross_val_score(full_pipeline_with_predictor,X_
best_train_accuracy = pct(best_train_scores.mean())\n\n# Best estimator prediction time and test accuracy
start = time.time()
best_test_accuracy = pct(grid_search.best_estimator_.score(X_test,
test_time = round(time.time() - start, 4)\n\n# Importance of features
features = numerical_features[:]
print('\nTotal number of features:', len(features))
importances = grid_search.best_estimator_.named_steps["predictor"].\n\n# selecting features based on importance values
new_indices = [idx for idx, x in enumerate(importances) if x>0.01]
new_importances = [x for idx, x in enumerate(importances) if x>0.01]
new_features = [features[i] for i in new_indices]\n\nprint('Total number of selected features:', len(new_features))\n\n# Plotting a barplot to visualize feature importance
sns.set(style='whitegrid')
```

```

plt.figure(figsize=(10, 6))
sns.barplot(x=importances, y=features, color='red')
plt.title('Feature Importances')
plt.xlabel('Relative Importance')
plt.ylabel('Feature')
plt.show()

# Conduct t-test with baseline logit and best estimator
(t_stat, p_value) = stats.ttest_rel(logit_scores, best_train_scores)

# Best parameters found using grid search
print(f"Best Parameters for {name}:")
best_parameters = grid_search.best_estimator_.get_params()
best_params = []
for param_name in sorted(params.keys()):
    best_params.append((param_name, best_parameters[param_name]))
    print("\t"+str(param_name)+": " + str(best_parameters[param_name]))
print("***** FINISHED TUNING", name, " *****")

# Results
results.loc[len(results)] = [name, best_train_accuracy, best_test_accuracy]

# Storing the importances of the features
features_dict['features'] = features
features_dict['importances'] = importances

```

In [121...]

ConductGridSearch(X_train[numerical_features], y_train, X_test[numerical_fe

```

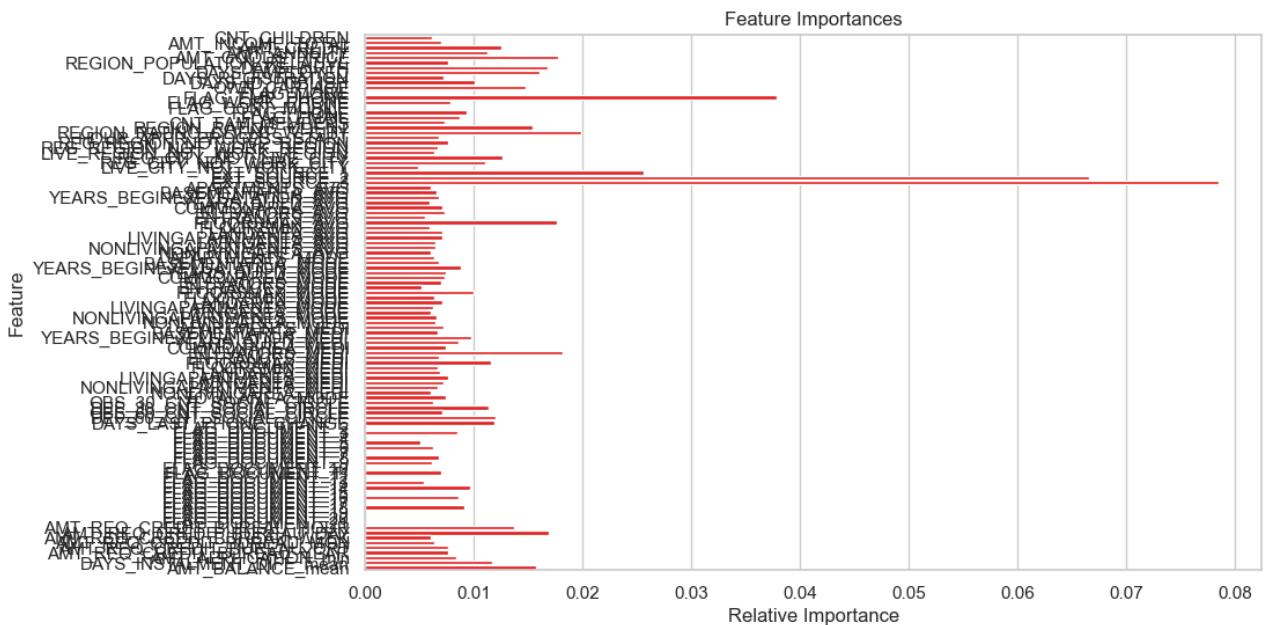
***** STARTING TUNING XGBoost *****
Parameters:
    colsample_bytree: [0.5, 0.7]
    eta: [0.001, 0.01]
    learning_rate: [0.001, 0.01]
    max_depth: [5, 9]
    n_estimators: [800, 1000]

```

Fitting 2 folds for each of 32 candidates, totalling 64 fits

Total number of features: 107

Total number of selected features: 25



Best Parameters for XGBoost:

```

predictor__colsample_bytree: 0.5
predictor__eta: 0.001
predictor__learning_rate: 0.01
predictor__max_depth: 5
predictor__n_estimators: 1000

```

***** FINISHED TUNING XGBoost *****

In [122...]

results

Out[122]:

ExpID		Cross-fold		Test Accuracy	p-value	Train Time(s)	Test Time(s)	Experiment Description
		Train Accuracy	Test Accuracy					
0	XGBoost	66.004	68.789	0.00148	13.1208	0.0174		[["predictor__colsample_bytree", 0.5], ["predi..."]]

In [123...]

```

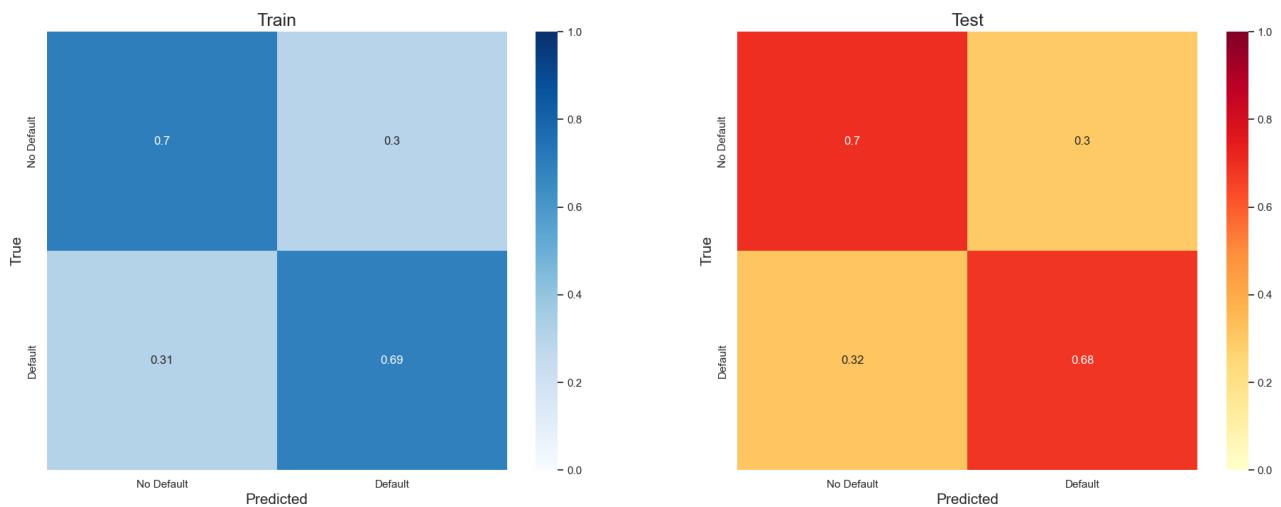
cm_train, cm_test=confusion_matrix_normalized(model,X_train,y_train,X_test,y
fig, axes = plt.subplots(1, 2, figsize=(23, 8))

# Plot the first heatmap in the first subplot
sns.heatmap(cm_train, vmin=0, vmax=1, annot=True, cmap="Blues", ax=axes[0])
axes[0].set_xlabel("Predicted", fontsize=15)
axes[0].set_ylabel("True", fontsize=15)
axes[0].set_xticklabels(class_labels)
axes[0].set_yticklabels(class_labels)
axes[0].set_title("Train", fontsize=18)

# Plot the second heatmap in the second subplot
sns.heatmap(cm_test, vmin=0, vmax=1, annot=True, cmap="YlOrRd", ax=axes[1])
axes[1].set_xlabel("Predicted", fontsize=15)
axes[1].set_ylabel("True", fontsize=15)
axes[1].set_xticklabels(class_labels)
axes[1].set_yticklabels(class_labels)
axes[1].set_title("Test", fontsize=18)

plt.show()

```



In [124...]

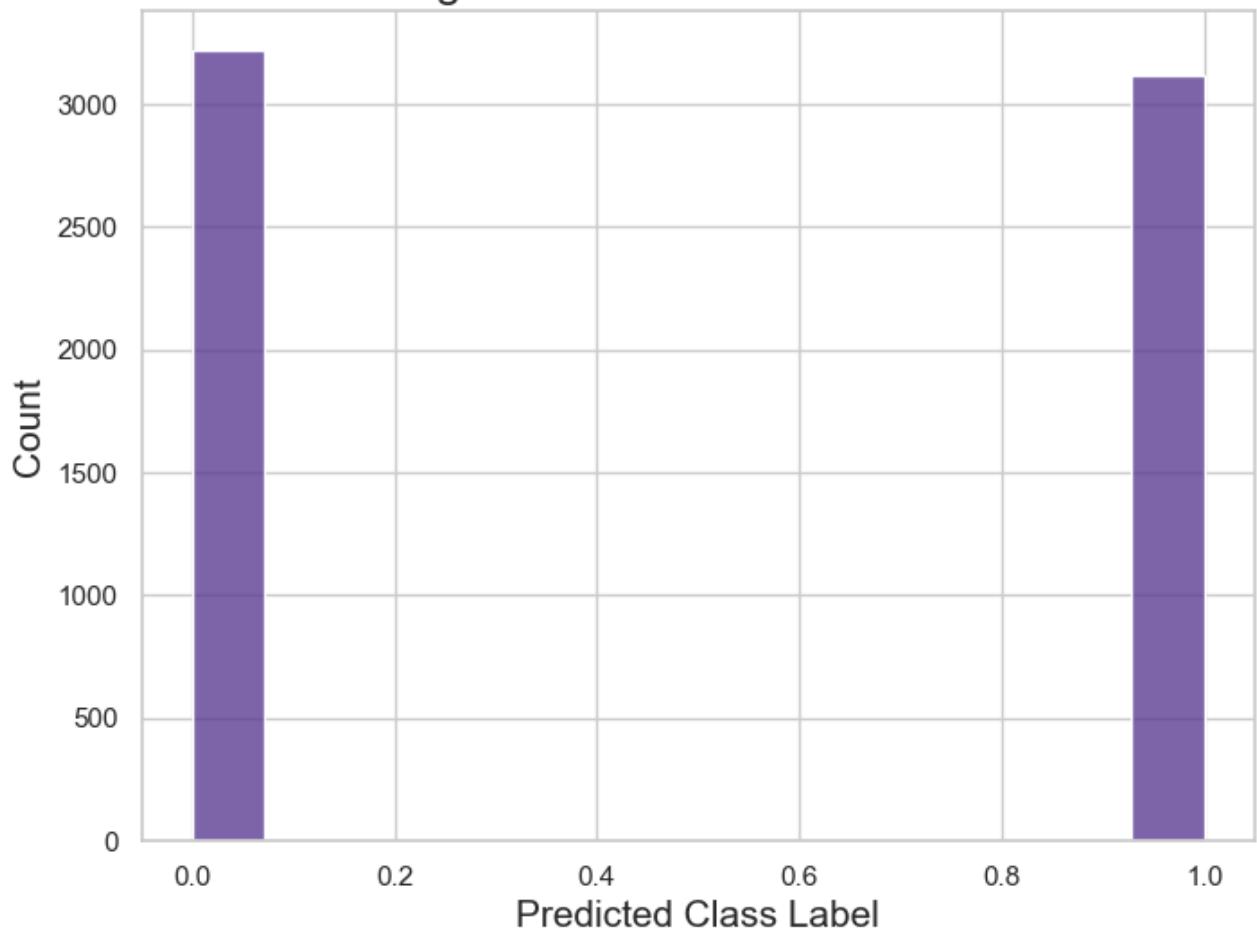
```

pred = model.predict(X_test)
# Create histogram of predicted class labels with a new color scheme
plt.figure(figsize=(8, 6))
sns.histplot(pred, kde=False, color="#5C3C92", alpha=0.8)
plt.xlabel("Predicted Class Label", fontsize=15)
plt.ylabel("Count", fontsize=15)
plt.title("Histogram of Predicted Class Labels", fontsize=18)
f1 = f1_score(y_test, pred)
print("F1 Score: ", f1)

```

F1 Score: 0.6906200317965024

Histogram of Predicted Class Labels



In [125...]

```
with open('features_dict_XG.pickle', 'wb') as handle:  
    pickle.dump(features_dict, handle, protocol=pickle.HIGHEST_PROTOCOL)
```

In [126...]

```
with open('features_dict_XG.pickle', 'rb') as handle:  
    x = pickle.load(handle)
```

I am currently performing feature selection and running model by filtering features based on their importance values. I have used three different thresholds for feature importance: $x>0$, $x>0.01$, and $x>0.005$. By using these thresholds, I am trying to understand the impact of including only the most relevant features in my model.

Here's a brief explanation of each threshold:

$x>0$:

This threshold includes all features with a non-zero importance value. It means I am considering all the features that have any impact on the model's prediction.

$x>0.01$:

This threshold is more stringent, as I am only considering features with importance values greater than 0.01. This filter results in a smaller number of features , which may help reduce the complexity of the model and the risk of overfitting.

$x>0.005$:

This threshold is more stringentThis threshold lies between the other two. By using this threshold, I am considering features with importance values greater than 0.005. This results in a slightly larger number of features compared to the $x>0.01$ threshold.

The goal of using these different thresholds is to find the optimal balance between model complexity and predictive performance. By comparing the results of the models trained with different feature sets, I can identify the best trade-off between model simplicity and performance.

Model 13 - XGBOOST -Feature &hyperParameter Tuning $x>0$

In [129...]

```
features = features_dict['features']
importances = features_dict['importances']

new_indices = [idx for idx, x in enumerate(importances) if x>0]
new_importances = [x for idx, x in enumerate(importances) if x>0]

new_features = [features[i] for i in new_indices]
print(len(new_features))

num_attribs = new_features

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline(num_attribs)

# Attaching XGBoost model to the above pipeline
xgboost_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("xgboost", XGBClassifier(random_state=42,
                               objective='binary:logistic', max_depth=5, eta=0.001,
                               learning_rate=0.01, colsample_bytree=0.5, n_estimators=1000
    ))
    
# Training the model
start = time.time()
model = xgboost_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = xgboost_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
exp_name = f"Model 13 - XGBOOST -Feature &hyperParameter Tuning"
experiment_description =f"XGBOOST Tuned with x>0 {len(selected_features)} f"
expLog = get_results(expLog, exp_name, experiment_description, model, train_t
expLog
```

95

Total Features: 112 - Numerical: 95, Categorical: 16

Out[129]:

		exp_name	description	Train Time (sec)	Test Time (sec)
0	Baseline_undersampled1_124_features	Baseline_undersampled1_124_features		2.1989	0.067
1	Model-1 Baseline LR	Logistic regression with undersampled data 124...		2.4495	0.056
2	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...		1.3669	0.045
3	Model-3 KNN	KNN with undersampled data-2 124 features		0.1888	0.611
4	Model-5 Baseline KNN	Decision tree with undersampled data-2 124 fe...		0.8050	0.023
5	Model-5 Random Forest	Random Forest with undersampled data-2 124 fe...		10.9049	0.320
6	Model-6 Extra Trees	Extra Trees with undersampled data-2 124 features		1.3018	0.099
7	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features		1.3846	0.104
8	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...		2.5269	0.171
9	Model-9 ADABoost SAMME	ADABoost SAMME with undersampled data-2 124 fe...		24.1679	0.353
10	Model-11 ADABoost	ADABoost SAMME with undersampled data-2 124 fe...		124.1598	2.851
11	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...		7.1576	0.206
12	Model-10 XGBoost	XGBoost SAMME with undersampled data-2 124 fe...		24.6309	0.035
13	Model-10 XGBoost	XGBoost SAMME with undersampled data-2 124 fe...		22.5526	0.035
14	Model-11 ADABoost	ADABoost SAMME with undersampled data-2 124 fe...		116.3467	2.844
15	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...		6.5630	0.200
16	Model 13 - XGBOOST -Feature &hyperParameter Tuning	XGBOOST Tuned with x>0 112 features		16.3471	0.031

Model 14 - XGBOOST -Feature & hyperParameter Tuning

In [130...]

```

features = features_dict['features']
importances = features_dict['importances']

new_indices = [idx for idx, x in enumerate(importances) if x>0.01]
new_importances = [x for idx, x in enumerate(importances) if x>0.01]

new_features = [features[i] for i in new_indices]
print(len(new_features))

num_attribs = new_features

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline(num_attribs)

# Attaching XGBoost model to the above pipeline
xgboost_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("xgboost", XGBClassifier(random_state=42,
                               objective='binary:logistic', max_depth=5, eta=0.001,
                               learning_rate=0.01, colsample_bytree=0.5, n_estimators=1000
    ))
    
# Training the model
start = time.time()
model = xgboost_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = xgboost_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
# Results
exp_name = f"Model 14 - XGBOOST -Feature &hyperParameter Tuning"
experiment_description =f"XGBOOST Tuned with x>0.01 {len(selected_features)}"
expLog = get_results(expLog, exp_name, experiment_description, model, train_t
expLog

```

25

Total Features: 42 - Numerical: 25, Categorical: 16

Out[130]:

		exp_name	description	Train Time (sec)	Test Time (sec)
0	Baseline_undersampled1_124_features	Baseline_undersampled1_124_features		2.1989	0.067
1	Model-1 Baseline LR	Logistic regression with undersampled data 124...		2.4495	0.056
2	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...		1.3669	0.045
3	Model-3 KNN	KNN with undersampled data-2 124 features		0.1888	0.611
4	Model-5 Baseline KNN	Decision tree with undersampled data-2 124 fea...		0.8050	0.023
5	Model-5 Random Forest	Random Forest with undersampled data-2 124 fea...		10.9049	0.320
6	Model-6 Extra Trees	Extra Trees with undersampled data-2 124 features		1.3018	0.099
7	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features		1.3846	0.104
8	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...		2.5269	0.171
9	Model-9 ADABoost SAMME	ADABoost SAMME with undersampled data-2 124 fe...		24.1679	0.353
10	Model-11 ADABoost	ADABoost SAMME with undersampled data-2 124 fe...		124.1598	2.851
11	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...		7.1576	0.206
12	Model-10 XGBoost	XGBoost SAMME with undersampled data-2 124 fe...		24.6309	0.035
13	Model-10 XGBoost	XGBoost SAMME with undersampled data-2 124 fe...		22.5526	0.035
14	Model-11 ADABoost	ADABoost SAMME with undersampled data-2 124 fe...		116.3467	2.844
15	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...		6.5630	0.200
16	Model 13 - XGBOOST -Feature &hyperParameter Tuning	XGBOOST Tuned with x>0 112 features		16.3471	0.031
17	Model 14 - XGBOOST -Feature &hyperParameter Tuning	XGBOOST Tuned with x>0.01 42 features		13.6924	0.025

Model 15 - XGBOOST -Feature &hyperParameter Tuning

In [131...]

```

features = features_dict['features']
importances = features_dict['importances']

new_indices = [idx for idx, x in enumerate(importances) if x>0.005]
new_importances = [x for idx, x in enumerate(importances) if x>0.005]

new_features = [features[i] for i in new_indices]
print(len(new_features))

num_attribs = new_features

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline(num_attribs)

# Attaching XGBoost model to the above pipeline
xgboost_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("xgboost", XGBClassifier(random_state=42,
                               objective='binary:logistic', max_depth=5, eta=0.001,
                               learning_rate=0.01, colsample_bytree=0.7, n_estimators=1000
    )))
# Training the model
start = time.time()
model = xgboost_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = xgboost_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
# Results
exp_name = f"Model 15 - XGBOOST -Feature &hyperParameter Tuning"
experiment_description = f"XGBOOST Tuned with x>0.005 {len(selected_features)}"
expLog = get_results(expLog, exp_name, experiment_description, model, train_time)
expLog

```

94

Total Features: 111 - Numerical: 94, Categorical: 16

Out[131]:

		exp_name	description	Train Time (sec)	Test Time (sec)
0	Baseline_undersampled1_124_features	Baseline_undersampled1_124_features		2.1989	0.067
1	Model-1 Baseline LR	Logistic regression with undersampled data 124...		2.4495	0.056
2	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...		1.3669	0.045
3	Model-3 KNN	KNN with undersampled data-2 124 features		0.1888	0.611
4	Model-5 Baseline KNN	Decision tree with undersampled data-2 124 fea...		0.8050	0.023
5	Model-5 Random Forest	Random Forest with undersampled data-2 124 fea...		10.9049	0.320
6	Model-6 Extra Trees	Extra Trees with undersampled data-2 124 features		1.3018	0.099
7	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features		1.3846	0.104
8	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...		2.5269	0.171
9	Model-9 ADABoost SAMME	ADABoost SAMME with undersampled data-2 124 fe...		24.1679	0.353
10	Model-11 ADABoost	ADABoost SAMME with undersampled data-2 124 fe...		124.1598	2.851
11	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...		7.1576	0.206
12	Model-10 XGBoost	XGBoost SAMME with undersampled data-2 124 fe...		24.6309	0.035
13	Model-10 XGBoost	XGBoost SAMME with undersampled data-2 124 fe...		22.5526	0.035
14	Model-11 ADABoost	ADABoost SAMME with undersampled data-2 124 fe...		116.3467	2.844
15	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...		6.5630	0.200
16	Model 13 - XGBOOST -Feature &hyperParameter Tu...	XGBOOST Tuned with x>0 112 features		16.3471	0.031
17	Model 14 - XGBOOST -Feature &hyperParameter Tu...	XGBOOST Tuned with x>0.01 42 features		13.6924	0.025
18	Model 15 - XGBOOST -Feature &hyperParameter Tu...	XGBOOST Tuned with x>0.005 111 features		21.2228	0.032

CatBoost Hyper Parameter Tuning and Feature Selection

In [83]:

```
import time
import json
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from scipy import stats
from sklearn.model_selection import GridSearchCV, ShuffleSplit, cross_val_score
from sklearn.pipeline import Pipeline, FeatureUnion
from catboost import CatBoostClassifier

# Helper function to convert decimal to percentage
def pct(decimal):
    return round(decimal * 100, 2)

results = pd.DataFrame(columns=["ExpID", "Cross-fold Train Accuracy", "Test Accuracy"])
features_dict = dict()

# A Function to execute the grid search and record the results.
def ConductGridSearch(X_train, y_train, X_test, y_test):
    # classifier for our grid search experiment
    classifiers = [
        ('CatBoost', CatBoostClassifier(random_state=42, verbose=False))
    ]

    # grid search parameters for the classifier
    param_grid = {
        'CatBoost': {
            'depth': [5, 9],
            'iterations': [800, 1000],
            'learning_rate': [0.001, 0.01],
            'colsample_bylevel': [0.5, 0.7],
        }
    }

    #     # grid search parameters for the classifier
    #     param_grid = {
    #         'CatBoost': {
    #             'depth': [5],
    #             'iterations': [20],
    #             'learning_rate': [0.01],
    #             'colsample_bylevel': [0.5],
    #         }
    #     }

    for (name, classifier) in classifiers:

        #name = "example"
        print(f"***** STARTING {name.upper()} *****")
```

```
parameters = param_grid[name]
print("Parameters:")
for p in sorted(parameters.keys()):
    print("\t"+str(p)+": "+ str(parameters[p]))

# generate the pipeline
full_pipeline_with_predictor = Pipeline([
    ("preparation", FeatureUnion(transformer_list=[("num_pipeline",
        ("predictor", classifier)
    ]))

# Execute the grid search
params = {}
for p in parameters.keys():
    pipe_key = 'predictor_'+str(p)
    params[pipe_key] = parameters[p]

grid_search = GridSearchCV(full_pipeline_with_predictor, params, sc
                           n_jobs=-1, verbose=1)
grid_search.fit(X_train, y_train)

# Best estimator training time
start = time.time()
grid_search.best_estimator_.fit(X_train, y_train)
train_time = round(time.time() - start, 4)

# Training accuracy
cvSplits = ShuffleSplit(n_splits=3, test_size=0.7, random_state=42)
best_train_scores = cross_val_score(full_pipeline_with_predictor,X_
best_train_accuracy = pct(best_train_scores.mean())

# Best estimator prediction time and test accuracy
start = time.time()
best_test_accuracy = pct(grid_search.best_estimator_.score(X_test,
test_time = round(time.time() - start, 4)

# Importance of features
features = numerical_features[:]
print('\nTotal number of features:', len(features))
importances = grid_search.best_estimator_.named_steps["predictor"].

# selecting features based on importance values
new_indices = [idx for idx, x in enumerate(importances) if x>0.01]
new_importances = [x for idx, x in enumerate(importances) if x>0.01
new_features = [features[i] for i in new_indices]

print('Total number of selected features:', len(new_features))

# Plotting a barplot to visualize feature importance
```

```
sns.set(style='whitegrid')
plt.figure(figsize=(10, 6))
sns.barplot(x=importances, y=features, color='red')
plt.title('Feature Importances')
plt.xlabel('Relative Importance')
plt.ylabel('Feature')
plt.show()

# Conduct t-test with baseline logit and best estimator
(t_stat, p_value) = stats.ttest_rel(logit_scores, best_train_scores)

# Best parameters found using grid search
print(f"Best Parameters for {name}:")
best_parameters = grid_search.best_estimator_.get_params()
best_params = []
for param_name in sorted(params.keys()):
    best_params.append((param_name, best_parameters[param_name]))
    print("\t"+str(param_name)+": " + str(best_parameters[param_name]))

print(f"***** FINISHED {name.upper()} *****")

# Results
results.loc[len(results)] = [name, best_train_accuracy, best_test_accuracy]

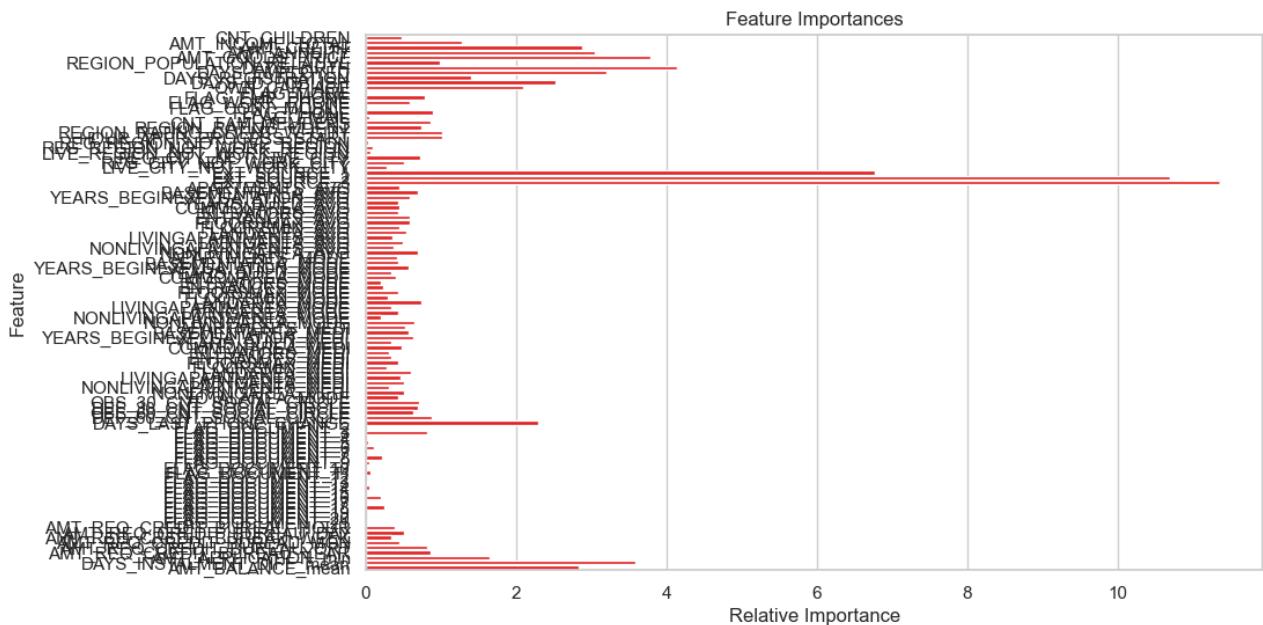
# Storing the importances of the features
features_dict['features'] = features
features_dict['importances'] = importances
```

In [84]:

```
ConductGridSearch(X_train[numerical_features], y_train, X_test[numerical_features])

***** STARTING CATBOOST *****
Parameters:
    colsample_bylevel: [0.5, 0.7]
    depth: [5, 9]
    iterations: [800, 1000]
    learning_rate: [0.001, 0.01]
Fitting 2 folds for each of 16 candidates, totalling 32 fits

Total number of features: 107
Total number of selected features: 95
```



Best Parameters for CatBoost:

```

predictor_colsample_bylevel: 0.5
predictor_depth: 9
predictor_iterations: 1000
predictor_learning_rate: 0.01
***** FINISHED CATBOOST *****

```

In [88]:

```
results
```

Out [88]:

ExpID		Cross-		Test Accuracy	p-value	Train Time(s)	Test Time(s)	Experiment Description
		fold Train	Accuracy					
0	CatBoost	68.02	68.55	0.3751	16.2538	0.0117		[["predictor_colsample_bylevel", 0.5], ["pred..."]]

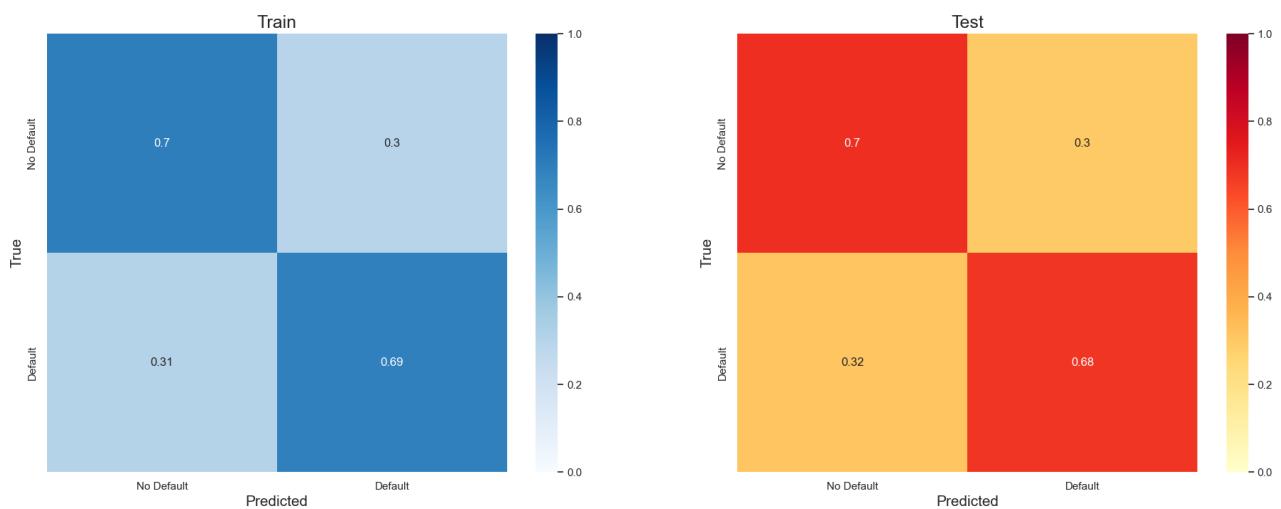
In [89]:

```
cm_train, cm_test=confusion_matrix_normalized(model,X_train,y_train,X_test,y
fig, axes = plt.subplots(1, 2, figsize=(23, 8))

# Plot the first heatmap in the first subplot
sns.heatmap(cm_train, vmin=0, vmax=1, annot=True, cmap="Blues", ax=axes[0])
axes[0].set_xlabel("Predicted", fontsize=15)
axes[0].set_ylabel("True", fontsize=15)
axes[0].set_xticklabels(class_labels)
axes[0].set_yticklabels(class_labels)
axes[0].set_title("Train", fontsize=18)

# Plot the second heatmap in the second subplot
sns.heatmap(cm_test, vmin=0, vmax=1, annot=True, cmap="YlOrRd", ax=axes[1])
axes[1].set_xlabel("Predicted", fontsize=15)
axes[1].set_ylabel("True", fontsize=15)
axes[1].set_xticklabels(class_labels)
axes[1].set_yticklabels(class_labels)
axes[1].set_title("Test", fontsize=18)

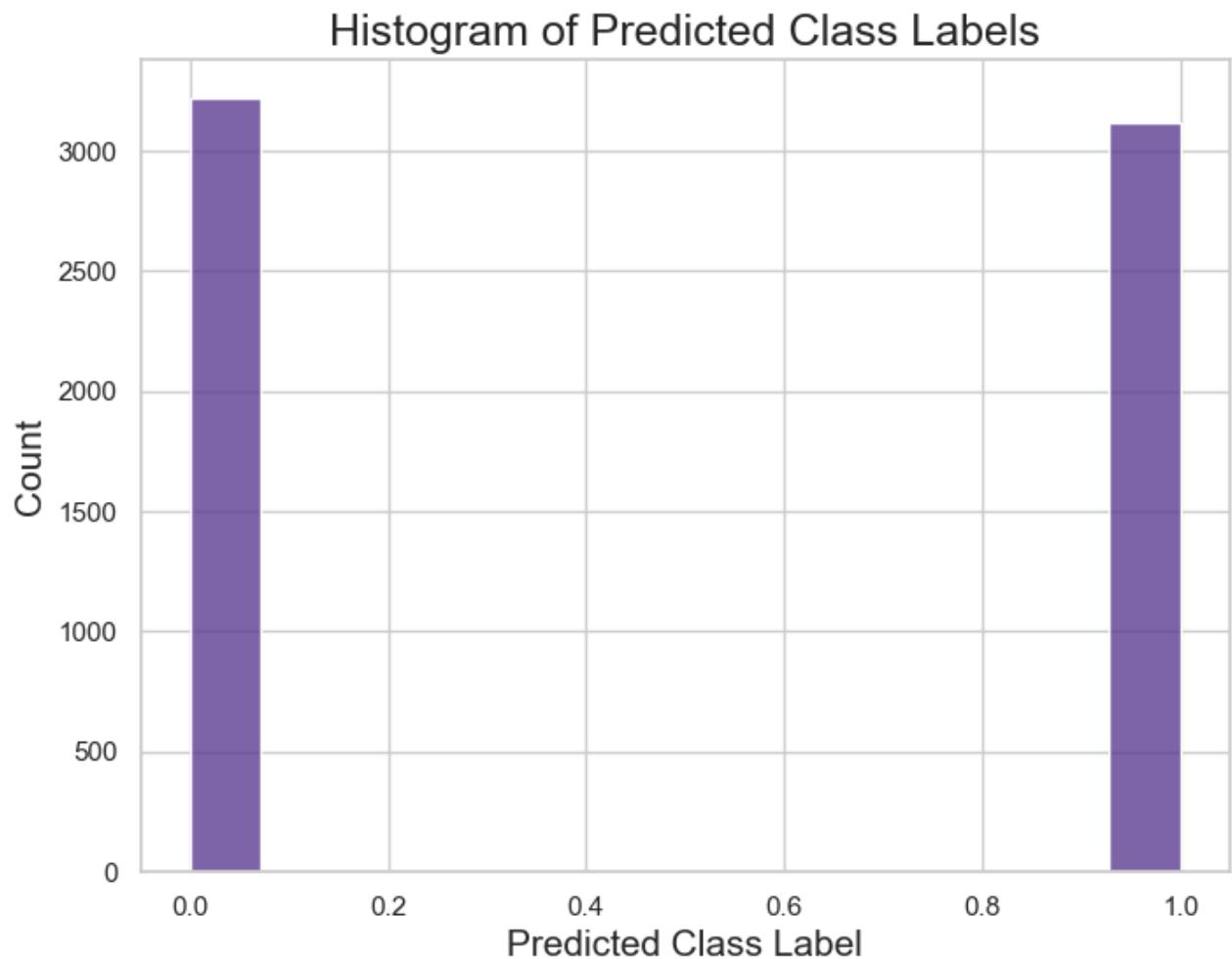
plt.show()
```



In [90]:

```
pred = model.predict(X_test)
# Create histogram of predicted class labels with a new color scheme
plt.figure(figsize=(8, 6))
sns.histplot(pred, kde=False, color="#5C3C92", alpha=0.8)
plt.xlabel("Predicted Class Label", fontsize=15)
plt.ylabel("Count", fontsize=15)
plt.title("Histogram of Predicted Class Labels", fontsize=18)
f1 = f1_score(y_test, pred)
print("F1 Score: ", f1)
```

F1 Score: 0.6906200317965024



In [91]:

```
with open('features_dict_catboost.pickle', 'wb') as handle:  
    pickle.dump(features_dict, handle, protocol=pickle.HIGHEST_PROTOCOL)
```

In [92]:

```
with open('features_dict_catboost.pickle', 'rb') as handle:  
    x = pickle.load(handle)
```

Model 16 - CatBOOST -Feature &hyperParameter Tuning

In [95]:

```

features = features_dict['features']
importances = features_dict['importances']

new_indices = [idx for idx, x in enumerate(importances) if x > 0]
new_importances = [x for idx, x in enumerate(importances) if x > 0]

new_features = [features[i] for i in new_indices]
print(len(new_features))

num_attribs = new_features

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline(num_attribs)

# Attaching CatBoost model to the above pipeline
catboost_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("catboost", CatBoostClassifier(random_state=42,
                                    iterations=1000, learning_rate=0.01, depth=9,
                                    colsample_bytree=0.5, thread_count=-1, verbose=False))
])

# Training the model
start = time.time()
model = catboost_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = catboost_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
exp_name = f"Model 16 - CatBOOST -Feature &hyperParameter Tuning"
experiment_description = f"CatBOOST Tuned with x>0 {len(selected_features)}"
expLog = get_results(expLog, exp_name, experiment_description, model, train_t
expLog

```

103

Total Features: 120 - Numerical: 103, Categorical: 16

Out[95]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.3559	0.0695	0.7730	0.7658	0.7736	0.7545	0.7407

Logistic

1	Model-2 Baseline	regression with undersampled data-2 1...	1.4254	0.0596	0.6876	0.6843	0.6904	0.7525	0.7489
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1751	0.5618	0.6950	0.6155	0.6184	0.7625	0.6571
3	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.7844	0.0233	0.6749	0.6535	0.6591	0.7380	0.7105
4	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.9068	0.3101	0.7664	0.6654	0.6661	0.8504	0.7243
5	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features	1.4413	0.1003	0.6738	0.6502	0.6482	0.7405	0.7059
6	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...	2.2834	0.1983	0.9844	0.6477	0.6446	0.9990	0.6978
7	Model-9 ADABOOST SAMME	ADABOOST SAMME with undersampled data-2 124 fe...	22.3935	0.2999	0.7185	0.6966	0.6950	0.7989	0.7580
8	Model-9 ADABOOST SAMME	ADABOOST SAMME with undersampled data-2 124 fe...	23.0637	0.3044	0.7182	0.6963	0.6969	0.7980	0.7585
9	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...	6.8641	0.2315	0.6964	0.6904	0.6926	0.7671	0.7580
10	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...	7.3245	0.2304	0.6964	0.6904	0.6926	0.7671	0.7580

		Model 16 - CatOOST - Feature &hyperParameter Tu...	CatBOOST Tuned with x>0 120 features	20.7806	0.2240	0.7528	0.6935	0.6970	0.8367	0.7594
11		Model 16 - CatBOOST - Feature &hyperParameter T...	CatBOOST Tuned with x>0 120 features	20.8028	0.2269	0.7528	0.6935	0.6970	0.8367	0.7594
12										

Model 17 - CatBOOST -Feature & HyperParameter Tuning

In [96]:

```

features = features_dict['features']
importances = features_dict['importances']

new_indices = [idx for idx, x in enumerate(importances) if x > 0.1]
new_importances = [x for idx, x in enumerate(importances) if x > 0.1]

new_features = [features[i] for i in new_indices]
print(len(new_features))

num_attribs = new_features

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline(num_attribs)

# Attaching CatBoost model to the above pipeline
catboost_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("catboost", CatBoostClassifier(random_state=42,
                                    iterations=1000, learning_rate=0.01, depth=9,
                                    colsample_bytree=0.5, thread_count=-1, verbose=False))
])

# Training the model
start = time.time()
model = catboost_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = catboost_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
exp_name = f"Model 17 - CatBOOST -Feature &hyperParameter Tuning"
experiment_description = f"CatBOOST Tuned with x>0.1 {len(selected_features)}"
expLog = get_results(expLog, exp_name, experiment_description, model, train_time)
expLog

```

86

Total Features: 103 - Numerical: 86, Categorical: 16

Out[96]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.3559	0.0695	0.7730	0.7658	0.7736	0.7545	0.7407

Logistic

1	Model-2 Baseline	regression with undersampled data-2 1...	1.4254	0.0596	0.6876	0.6843	0.6904	0.7525	0.7489
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1751	0.5618	0.6950	0.6155	0.6184	0.7625	0.6571
3	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.7844	0.0233	0.6749	0.6535	0.6591	0.7380	0.7105
4	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.9068	0.3101	0.7664	0.6654	0.6661	0.8504	0.7243
5	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features	1.4413	0.1003	0.6738	0.6502	0.6482	0.7405	0.7059
6	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...	2.2834	0.1983	0.9844	0.6477	0.6446	0.9990	0.6978
7	Model-9 ADABOOST SAMME	ADABOOST SAMME with undersampled data-2 124 fe...	22.3935	0.2999	0.7185	0.6966	0.6950	0.7989	0.7580
8	Model-9 ADABOOST SAMME	ADABOOST SAMME with undersampled data-2 124 fe...	23.0637	0.3044	0.7182	0.6963	0.6969	0.7980	0.7585
9	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...	6.8641	0.2315	0.6964	0.6904	0.6926	0.7671	0.7580
10	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...	7.3245	0.2304	0.6964	0.6904	0.6926	0.7671	0.7580

11	Model 16 - CatOOST - Feature &hyperParameter Tuning	CatBOOST Tuned with x>0 120 features	20.7806	0.2240	0.7528	0.6935	0.6970	0.8367	0.7594
12	Model 16 - CatBOOST - Feature &hyperParameter Tuning	CatBOOST Tuned with x>0 120 features	20.8028	0.2269	0.7528	0.6935	0.6970	0.8367	0.7594
13	Model 17 - CatBOOST - Feature &hyperParameter Tuning	CatBOOST Tuned with x>0.1 103 features	19.0881	0.2117	0.7513	0.6897	0.6944	0.8358	0.7590

Model 18 - CatBOOST -Feature &hyperParameter Tuning

In [97]:

```

features = features_dict['features']
importances = features_dict['importances']

new_indices = [idx for idx, x in enumerate(importances) if x > 0.005]
new_importances = [x for idx, x in enumerate(importances) if x > 0.005]

new_features = [features[i] for i in new_indices]
print(len(new_features))

num_attribs = new_features

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline(num_attribs)

# Attaching CatBoost model to the above pipeline
catboost_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("catboost", CatBoostClassifier(random_state=42,
                                    iterations=1000, learning_rate=0.01, depth=9,
                                    colsample_bytree=0.5, thread_count=-1, verbose=False))
])

# Training the model
start = time.time()
model = catboost_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = catboost_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
exp_name = f"Model 18 - CatBOOST -Feature &hyperParameter Tuning"
experiment_description = f"CatBOOST Tuned with x>0.005 {len(selected_features)}"
expLog = get_results(expLog, exp_name, experiment_description, model, train_time, test_time)
expLog

```

96

Total Features: 113 - Numerical: 96, Categorical: 16

Out[97]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.3559	0.0695	0.7730	0.7658	0.7736	0.7545	0.7407

Logistic

1	Model-2 Baseline	regression with undersampled data-2 1...	LR	1.4254	0.0596	0.6876	0.6843	0.6904	0.7525	0.7489
2	Model-3 KNN	KNN with undersampled data-2 124 features		0.1751	0.5618	0.6950	0.6155	0.6184	0.7625	0.6571
3	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...		0.7844	0.0233	0.6749	0.6535	0.6591	0.7380	0.7105
4	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...		10.9068	0.3101	0.7664	0.6654	0.6661	0.8504	0.7243
5	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features		1.4413	0.1003	0.6738	0.6502	0.6482	0.7405	0.7059
6	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...		2.2834	0.1983	0.9844	0.6477	0.6446	0.9990	0.6978
7	Model-9 ADABOOST SAMME	ADABOOST SAMME with undersampled data-2 124 fe...		22.3935	0.2999	0.7185	0.6966	0.6950	0.7989	0.7580
8	Model-9 ADABOOST SAMME	ADABOOST SAMME with undersampled data-2 124 fe...		23.0637	0.3044	0.7182	0.6963	0.6969	0.7980	0.7585
9	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...		6.8641	0.2315	0.6964	0.6904	0.6926	0.7671	0.7580
10	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...		7.3245	0.2304	0.6964	0.6904	0.6926	0.7671	0.7580

	Model 16 - CatOOST - Feature &hyperParameter Tu...	CatBOOST Tuned with x>0 120 features	20.7806 0.2240 0.7528 0.6935 0.6970 0.8367 0.7594
11	Model 16 - CatBOOST - Feature &hyperParameter T...	CatBOOST Tuned with x>0 120 features	20.8028 0.2269 0.7528 0.6935 0.6970 0.8367 0.7594
12	Model 17 - CatBOOST - Feature &hyperParameter T...	CatBOOST Tuned with x>0.1 103 features	19.0881 0.2117 0.7513 0.6897 0.6944 0.8358 0.7590
13	Model 18 - CatBOOST - Feature &hyperParameter T...	CatBOOST Tuned with x>0.005 113 features	19.3895 0.2194 0.7513 0.6931 0.6959 0.8360 0.7596
14			

Random Forest Hyper Parameter Tuning and Feature Selection

In [79]:

```

from sklearn.ensemble import RandomForestClassifier

results = pd.DataFrame(columns=[ "ExpID", "Cross-fold Train Accuracy", "Test
features_dict = dict()

# A Function to execute the grid search and record the results.
def ConductGridSearch(X_train, y_train, X_test, y_test):
    # classifier for our grid search experiment
    classifiers = [
        ('RandomForest', RandomForestClassifier(random_state=42))
    ]

    # grid search parameters for the classifier
    param_grid = {
        'RandomForest': {
            'n_estimators': [100, 200],
            'max_depth': [5, 10],
            'min_samples_split': [2, 5],
            'min_samples_leaf': [1, 2],
            'max_features': ['auto', 'sqrt']
        }
    }

    for (name, classifier) in classifiers:

        print('***** START', name, '*****')

```

```
parameters = param_grid[name]
print("Parameters:")
for p in sorted(parameters.keys()):
    print("\t"+str(p)+": "+ str(parameters[p]))

# generate the pipeline
full_pipeline_with_predictor = Pipeline([
    ("preparation", FeatureUnion(transformer_list=[("num_pipeline",
    ("predictor", classifier)
])))

# Execute the grid search
params = {}
for p in parameters.keys():
    pipe_key = 'predictor_'+str(p)
    params[pipe_key] = parameters[p]

grid_search = GridSearchCV(full_pipeline_with_predictor, params, sc
                           n_jobs=-1, verbose=1)
grid_search.fit(X_train, y_train)

# Best estimator training time
start = time.time()
grid_search.best_estimator_.fit(X_train, y_train)
train_time = round(time.time() - start, 4)

# Training accuracy
cvSplits = ShuffleSplit(n_splits=3, test_size=0.7, random_state=42)
best_train_scores = cross_val_score(full_pipeline_with_predictor,X_
best_train_accuracy = pct(best_train_scores.mean())

# Best estimator prediction time and test accuracy
start = time.time()
best_test_accuracy = pct(grid_search.best_estimator_.score(X_test,
test_time = round(time.time() - start, 4)

# Importance of features
features = numerical_features[:]
print('\nTotal number of features:', len(features))
importances = grid_search.best_estimator_.named_steps["predictor"].

# selecting features based on importance values
new_indices = [idx for idx, x in enumerate(importances) if x>0.01]
new_importances = [x for idx, x in enumerate(importances) if x>0.01]
new_features = [features[i] for i in new_indices]

print('Total number of selected features:', len(new_features))

# Plotting a barplot to visualize feature importance
```

```

sns.set(style='whitegrid')
plt.figure(figsize=(10, 6))
sns.barplot(x=importances, y=features, color='red')
plt.title('Feature Importances')
plt.xlabel('Relative Importance')
plt.ylabel('Feature')
plt.show()

# Conduct t-test with baseline logit and best estimator
(t_stat, p_value) = stats.ttest_rel(logit_scores, best_train_scores
    # Best parameters found using grid search
print(f"Best Parameters for {name}:")
best_parameters = grid_search.best_estimator_.get_params()
best_params = []
for param_name in sorted(params.keys()):
    best_params.append((param_name, best_parameters[param_name]))
    print("\t"+str(param_name)+": " + str(best_parameters[param_name]))
print("***** FINISH", name, " *****")

# Results
results.loc[len(results)] = [name, best_train_accuracy, best_test_a

# Storing the importances of the features
features_dict['features'] = features
features_dict['importances'] = importances

```

In [80]:

```

ConductGridSearch(X_train[numerical_features], y_train, X_test[numerical_fe
***** START RandomForest *****
Parameters:
    max_depth: [5, 10]
    max_features: ['auto', 'sqrt']
    min_samples_leaf: [1, 2]
    min_samples_split: [2, 5]
    n_estimators: [100, 200]
Fitting 2 folds for each of 32 candidates, totalling 64 fits
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_fores
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe
atures='sqrt'` or remove this parameter as it is also the default value for
RandomForestClassifiers and ExtraTreesClassifiers.
    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_fores
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe
atures='sqrt'` or remove this parameter as it is also the default value for
RandomForestClassifiers and ExtraTreesClassifiers.
    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_fores
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe

```

```
atures='sqrt'` or remove this parameter as it is also the default value for
RandomForestClassifiers and ExtraTreesClassifiers.

    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.

    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.

    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.

    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.

    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.

    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.

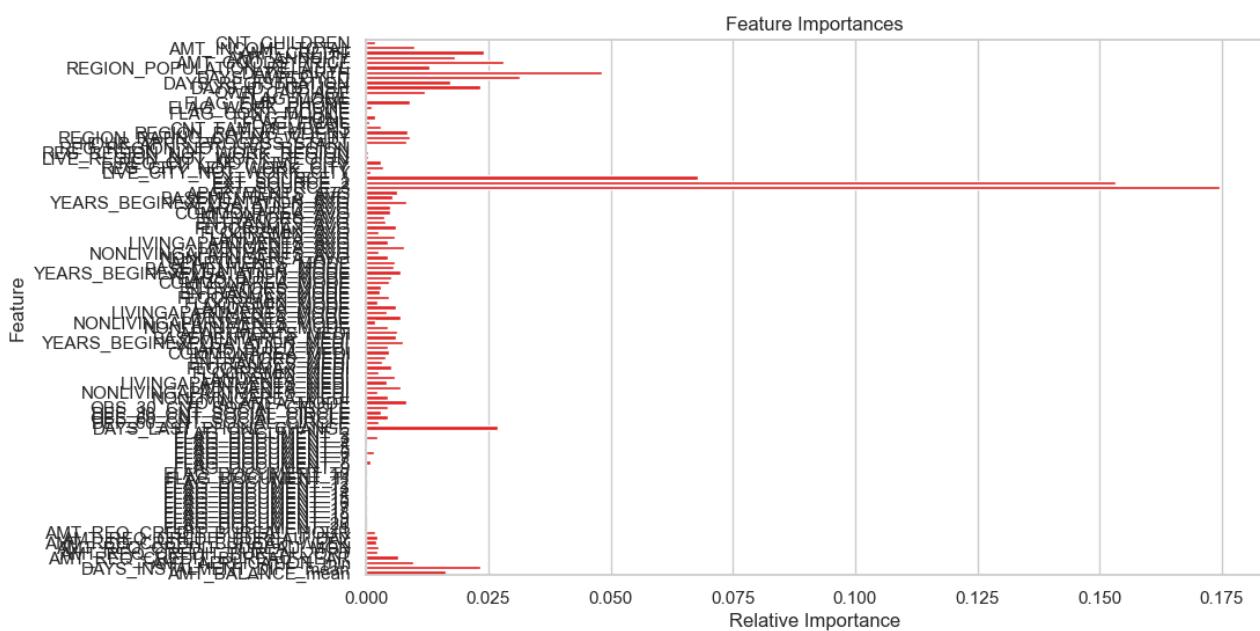
    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.

    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.
```

```
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.
```

```
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an  
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe  
atures='sqrt'` or remove this parameter as it is also the default value for  
RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_fore  
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an  
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe  
atures='sqrt'` or remove this parameter as it is also the default value for  
RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_fore  
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an  
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe  
atures='sqrt'` or remove this parameter as it is also the default value for  
RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_fore  
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an  
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe  
atures='sqrt'` or remove this parameter as it is also the default value for  
RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_fore  
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an  
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe  
atures='sqrt'` or remove this parameter as it is also the default value for  
RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_fore  
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an  
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe  
atures='sqrt'` or remove this parameter as it is also the default value for  
RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_fore  
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an  
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe  
atures='sqrt'` or remove this parameter as it is also the default value for  
RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_fore  
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an  
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe  
atures='sqrt'` or remove this parameter as it is also the default value for  
RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_fore  
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an  
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe  
atures='sqrt'` or remove this parameter as it is also the default value for  
RandomForestClassifiers and ExtraTreesClassifiers.
```

```
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.  
    warn(  
Total number of features: 107
```



Best Parameters for RandomForest:

```
predictor_max_depth: 10
predictor_max_features: auto
predictor_min_samples_leaf: 2
predictor_min_samples_split: 5
predictor_n_estimators: 200
***** FINISH RandomForest *****
```

In [81]:

results

Out[81]:

ExpID		Cross-fold		Test Accuracy	p-value	Train Time(s)	Test Time(s)	Experiment Description
		Train Accuracy	Accuracy					
0	RandomForest	66.567	67.162	0.00287		8.1618	0.089	[{"predictor_max_depth": 10}, {"predictor_ma...]

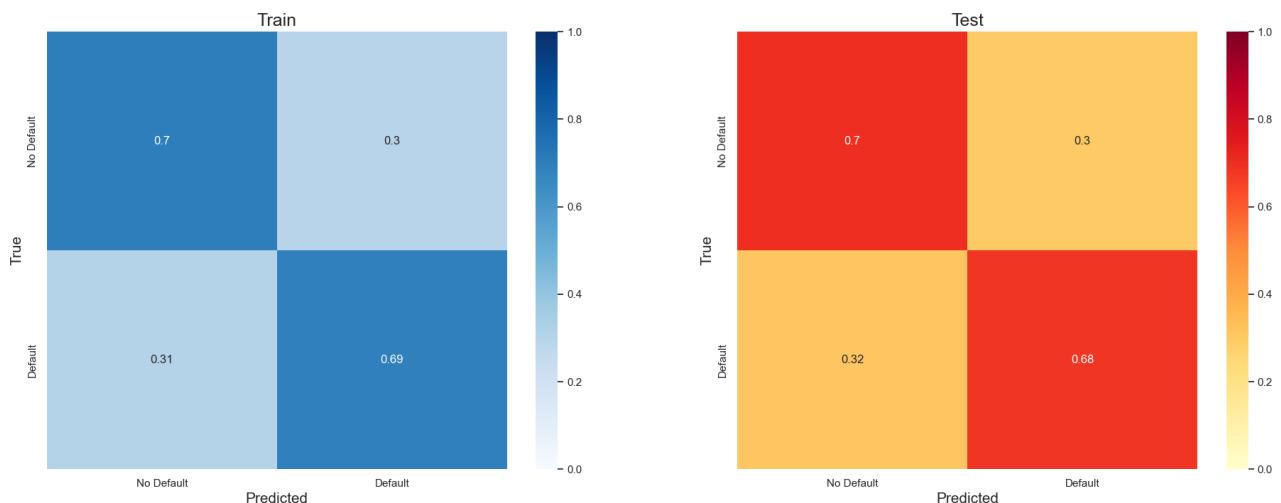
In [82]:

```
cm_train, cm_test=confusion_matrix_normalized(model,X_train,y_train,X_test,y_test)

# Plot the first heatmap in the first subplot
sns.heatmap(cm_train, vmin=0, vmax=1, annot=True, cmap="Blues", ax=axes[0])
axes[0].set_xlabel("Predicted", fontsize=15)
axes[0].set_ylabel("True", fontsize=15)
axes[0].set_xticklabels(class_labels)
axes[0].set_yticklabels(class_labels)
axes[0].set_title("Train", fontsize=18)

# Plot the second heatmap in the second subplot
sns.heatmap(cm_test, vmin=0, vmax=1, annot=True, cmap="YlOrRd", ax=axes[1])
axes[1].set_xlabel("Predicted", fontsize=15)
axes[1].set_ylabel("True", fontsize=15)
axes[1].set_xticklabels(class_labels)
axes[1].set_yticklabels(class_labels)
axes[1].set_title("Test", fontsize=18)

plt.show()
```

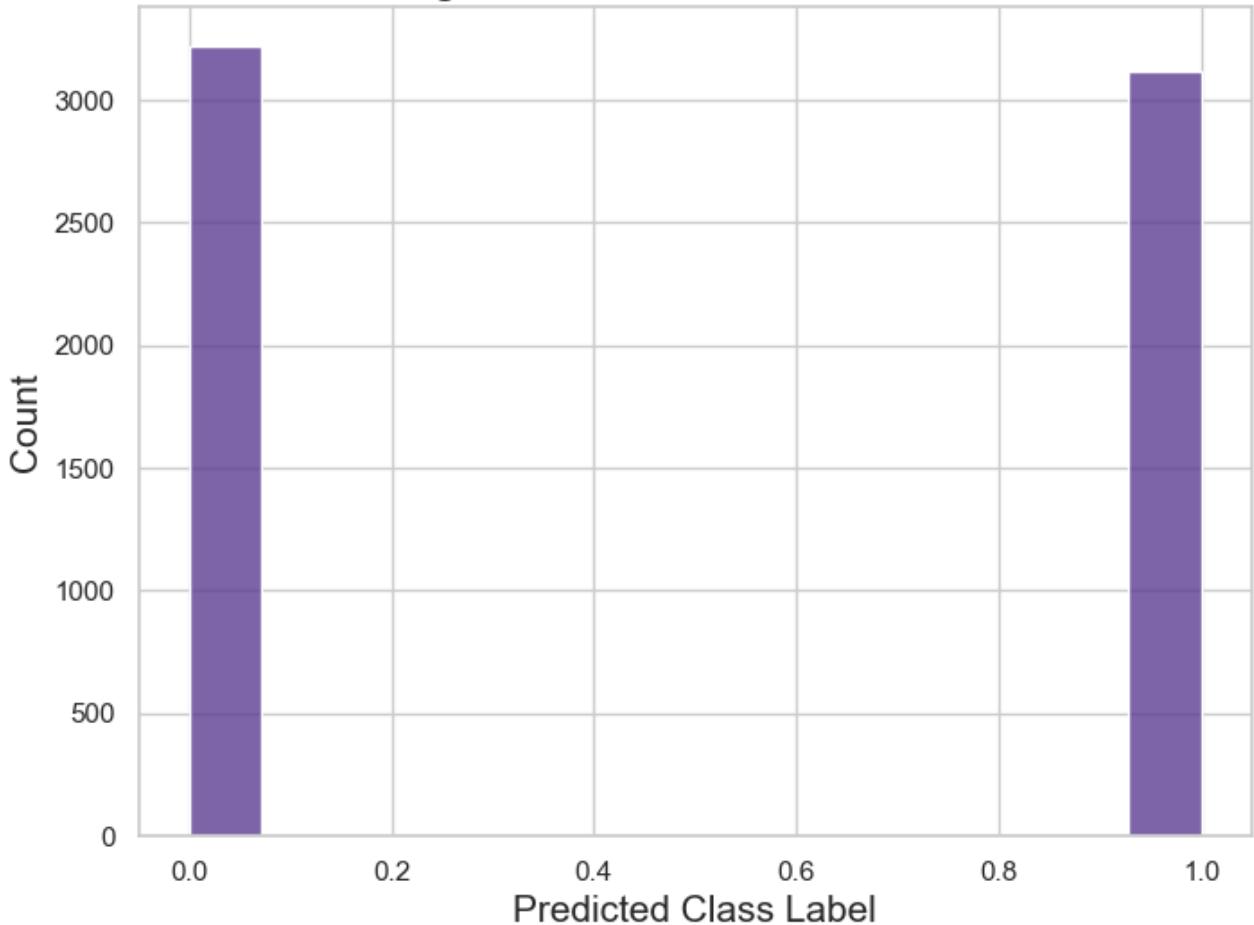


In [83]:

```
pred = model.predict(X_test)
# Create histogram of predicted class labels with a new color scheme
plt.figure(figsize=(8, 6))
sns.histplot(pred, kde=False, color="#5C3C92", alpha=0.8)
plt.xlabel("Predicted Class Label", fontsize=15)
plt.ylabel("Count", fontsize=15)
plt.title("Histogram of Predicted Class Labels", fontsize=18)
f1 = f1_score(y_test, pred)
print("F1 Score: ", f1)
```

F1 Score: 0.6906200317965024

Histogram of Predicted Class Labels



In [84]:

```
with open('features_dict_rf.pickle', 'wb') as handle:
    pickle.dump(features_dict, handle, protocol=pickle.HIGHEST_PROTOCOL)
```

In [85]:

```
with open('features_dict_rf.pickle', 'rb') as handle:
    x = pickle.load(handle)
```

Model 19- Random Forest -Feature &hyperParameter Tuning

In [87]:

```

features = features_dict['features']
importances = features_dict['importances']

new_indices = [idx for idx, x in enumerate(importances) if x > 0]
new_importances = [x for idx, x in enumerate(importances) if x > 0]

new_features = [features[i] for i in new_indices]
print(len(new_features))

num_attribs = new_features

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline(num_attribs)

# Attaching RandomForest model to the above pipeline
random_forest_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("random_forest", RandomForestClassifier(random_state=42,
                                              n_estimators=200, max_depth=10, max_features='auto',
                                              min_samples_leaf=2, min_samples_split=5, n_jobs=-1))
])

# Training the model
start = time.time()
model = random_forest_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = random_forest_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
# Results
exp_name = f"Model 19 - Random Forest -Feature &hyperParameter Tuning"
experiment_description = f"Random Forest Tuned with x>0 {len(selected_features)}"
expLog = get_results(expLog, exp_name, experiment_description, model, train_time, test_time)
expLog

```

99

Total Features: 116 - Numerical: 99, Categorical: 16

Out[87]:

exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
Logistic								

0	Model-1 Baseline LR	regression with undersampled data 124...	2.3376	0.0619	0.7731	0.7666	0.7777	0.7525	0.7433
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4645	0.0486	0.6876	0.6843	0.6904	0.7525	0.7489
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1864	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571
3	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fe...	0.7840	0.0232	0.6749	0.6535	0.6591	0.7380	0.7105
4	Model-6 Random Forest	Random Forest with undersampled data-2 124 fe...	10.8707	0.3122	0.7664	0.6654	0.6661	0.8504	0.7243
5	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features	1.4820	0.0905	0.6738	0.6502	0.6482	0.7405	0.7059
6	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...	2.7168	0.1855	0.9844	0.6477	0.6446	0.9990	0.6978
7	Model-9 ADABOOST SAMME	ADABOOST SAMME with undersampled data-2 124 fe...	22.8602	0.3029	0.7185	0.6966	0.6950	0.7989	0.7580
8	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...	6.4653	0.1827	0.6964	0.6904	0.6926	0.7671	0.7580
9	Model 18 - CatBOOST - Feature &hyperParameter T...	CatBOOST Tuned with x>0 116 features	1.1723	0.0520	0.7531	0.6821	0.6801	0.8325	0.7412
	Model 19 -	Random							

10	Random Forest - Feature &hyperParam...	Forest Tuned with x>0.116 features	1.1767	0.0480	0.7531	0.6821	0.6801	0.8325	0.7412
----	--	------------------------------------	--------	--------	--------	--------	--------	--------	--------

Model 20 - Random Forest -Feature &hyperParameter Tuning

In [88]:

```

features = features_dict['features']
importances = features_dict['importances']

new_indices = [idx for idx, x in enumerate(importances) if x > 0.1]
new_importances = [x for idx, x in enumerate(importances) if x > 0.1]

new_features = [features[i] for i in new_indices]
print(len(new_features))

num_attribs = new_features

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline(num_attribs)

# Attaching RandomForest model to the above pipeline
random_forest_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("random_forest", RandomForestClassifier(random_state=42,
                                              n_estimators=200, max_depth=10, max_features='auto',
                                              min_samples_leaf=2, min_samples_split=5, n_jobs=-1))
])

# Training the model
start = time.time()
model = random_forest_full_pipeline_with_predictor.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = random_forest_full_pipeline_with_predictor.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
exp_name = f"Model 20 - Random Forest -Feature &hyperParameter Tuning"
experiment_description = f"Random Forest Tuned with x>0.1 {len(selected_features)} features"
expLog = get_results(expLog, exp_name, experiment_description, model, train_time, test_time)
expLog

```

2

Total Features: 19 - Numerical: 2, Categorical: 16

Out[88]:

Train	Test	Train	Valid	Test	Train	Valid
-------	------	-------	-------	------	-------	-------

	exp_name	description	Time (sec)	Time (sec)	Acc	Acc	Acc	AUC	AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.3376	0.0619	0.7731	0.7666	0.7777	0.7525	0.7433
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4645	0.0486	0.6876	0.6843	0.6904	0.7525	0.7489
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1864	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571
3	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.7840	0.0232	0.6749	0.6535	0.6591	0.7380	0.7105
4	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.8707	0.3122	0.7664	0.6654	0.6661	0.8504	0.7243
5	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features	1.4820	0.0905	0.6738	0.6502	0.6482	0.7405	0.7059
6	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...	2.7168	0.1855	0.9844	0.6477	0.6446	0.9990	0.6978
7	Model-9 ADABOOST SAMME	ADABOOST SAMME with undersampled data-2 124 fe...	22.8602	0.3029	0.7185	0.6966	0.6950	0.7989	0.7580
8	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...	6.4653	0.1827	0.6964	0.6904	0.6926	0.7671	0.7580
9	Model 18 - CatBOOST - Feature	CatBOOST Tuned with	1.1723	0.0520	0.7531	0.6821	0.6801	0.8325	0.7412

		&hyperParameter T...	x>0 116 features										
10		Model 19 - Random Forest - Feature &hyperParame...	Random Forest Tuned with x>0 116 features	1.1767	0.0480	0.7531	0.6821	0.6801	0.8325	0.7412			
11		Model 20 - Random Forest - Feature &hyperParame...	Random Forest Tuned with x>0.1 19 features	0.5466	0.0414	0.6960	0.6733	0.6631	0.7617	0.7264			

Model 21 - Random Forest -Feature &hyperParameter Tuning

In [89]:

```

features = features_dict['features']
importances = features_dict['importances']

new_indices = [idx for idx, x in enumerate(importances) if x > 0.005]
new_importances = [x for idx, x in enumerate(importances) if x > 0.005]

new_features = [features[i] for i in new_indices]
print(len(new_features))

num_attribs = new_features

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline, selected_features = get_pipeline(num_attribs)

# Attaching RandomForest model to the above pipeline
random_forest_full_pipeline_with_predictor = Pipeline([
    ("preparation", data_prep_pipeline),
    ("random_forest", RandomForestClassifier(random_state=42,
                                              n_estimators=200, max_depth=10, max_features='auto',
                                              min_samples_leaf=2, min_samples_split=5, n_jobs=-1))
])

# Training the model
start = time.time()
model = random_forest_full_pipeline_with_predictor.fit(x_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = random_forest_full_pipeline_with_predictor.score(x_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
exp_name = f"Model 21 - Random Forest -Feature &hyperParameter Tuning"
experiment_description = f"Random Forest Tuned with x>0.005 {len(selected_features)} attribs"
expLog = get_results(expLog, exp_name, experiment_description, model, train_time, test_time)
expLog

```

41

Total Features: 58 - Numerical: 41, Categorical: 16

Out[89]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.3376	0.0619	0.7731	0.7666	0.7777	0.7525	0.7433

Logistic

1	Model-2 Baseline LR	regression with undersampled data-2 1...	1.4645	0.0486	0.6876	0.6843	0.6904	0.7525	0.7489
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1864	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571
3	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.7840	0.0232	0.6749	0.6535	0.6591	0.7380	0.7105
4	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.8707	0.3122	0.7664	0.6654	0.6661	0.8504	0.7243
5	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features	1.4820	0.0905	0.6738	0.6502	0.6482	0.7405	0.7059
6	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...	2.7168	0.1855	0.9844	0.6477	0.6446	0.9990	0.6978
7	Model-9 ADABoost SAMME	ADABOOST SAMME with undersampled data-2 124 fe...	22.8602	0.3029	0.7185	0.6966	0.6950	0.7989	0.7580
8	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...	6.4653	0.1827	0.6964	0.6904	0.6926	0.7671	0.7580
9	Model 18 - CatBOOST - Feature &hyperParameter T...	CatBOOST Tuned with x>0 116 features	1.1723	0.0520	0.7531	0.6821	0.6801	0.8325	0.7412
10	Model 19 - Random Forest - Feature &hyperParame...	Random Forest Tuned with x>0 116 features	1.1767	0.0480	0.7531	0.6821	0.6801	0.8325	0.7412
	Model 20 -	Random							

11	Random Forest - Feature &hyperParam...	Forest Tuned with x>0.1 19 features	0.5466	0.0414	0.6960	0.6733	0.6631	0.7617	0.7264
12	Model 21 - Random Forest - Feature &hyperParam...	Random Forest Tuned with x>0.005 58 features	1.0561	0.0435	0.7513	0.6841	0.6808	0.8299	0.7429

In [90]:

```
# Write the data to a CSV file
expLog.to_csv('expLog.csv', index=False)
```

In [98]:

```
df = pd.read_csv('expLog1.csv')
df
```

Out[98]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.3376	0.0619	0.7731	0.7666	0.7777	0.7525	0.7433
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4645	0.0486	0.6876	0.6843	0.6904	0.7525	0.7489
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1864	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571
3	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
4	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.7840	0.0232	0.6749	0.6535	0.6591	0.7380	0.7105
5	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.8707	0.3122	0.7664	0.6654	0.6661	0.8504	0.7243
6	Model-7 Extra Trees	Extra Trees with undersampled	1.4820	0.0905	0.6738	0.6502	0.6482	0.7405	0.7059

7	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-2 124 features	2.7168	0.1855	0.9844	0.6477	0.6446	0.9990	0.6978
8	Model-9 ADABoost SAMME	ADABoost SAMME with undersampled data-2 124 features	22.8602	0.3029	0.7185	0.6966	0.6950	0.7989	0.7580
9		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
10		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
11	Model-12 CATBoost	CATBoost with undersampled data-2 124 features	6.4653	0.1827	0.6964	0.6904	0.6926	0.7671	0.7580
12	Model 13 - XGBOOST - Feature &hyperParameter Tun...	XGBOOST Tuned with x>0 112 features	16.3471	0.0319	0.7308	0.6959	0.6956	0.8101	0.7610
13	Model 14 - XGBOOST - Feature &hyperParameter Tun...	XGBOOST Tuned with x>0.1 42 features	13.6924	0.0257	0.7245	0.6944	0.6952	0.8029	0.7610
14	Model 15 - XGBOOST - Feature &hyperParameter Tun...	XGBOOST Tuned with x>0.005 111 features	21.2228	0.0322	0.7299	0.6931	0.6939	0.8110	0.7610
15	Model 16 - CatBOOST - Feature &hyperParameter Tun...	CatBOOST Tuned with x>0 120 features	20.8028	0.2269	0.7528	0.6935	0.6970	0.8367	0.7594
16	Model 17 - CatBOOST - Feature &hyperParameter Tun...	CatBOOST Tuned with x>0.1 103 features	19.0881	0.2117	0.7513	0.6897	0.6944	0.8358	0.7590
17	Model 18 - CatBOOST - Feature &hyperParameter Tun...	CatBOOST Tuned with x>0.005 113	19.3895	0.2194	0.7513	0.6931	0.6959	0.8360	0.7596

	T...	features										
18	Model 19 - Random Forest - Feature &hyperParam...	Random Forest Tuned with x>0.116 features	1.1767	0.0480	0.7531	0.6821	0.6801	0.8325	0.7412			
19	Model 20 - Random Forest - Feature &hyperParam...	Random Forest Tuned with x>0.119 features	0.5466	0.0414	0.6960	0.6733	0.6631	0.7617	0.7264			
20	Model 21 - Random Forest - Feature &hyperParam...	Random Forest Tuned with x>0.005 58 features	1.0561	0.0435	0.7513	0.6841	0.6808	0.8299	0.7429			

Loading selected Features to base classifiers before ensembling them

In [99]:

```
#for Random forest

with open('features_dict_XG.pickle', 'rb') as handle:
    x = pickle.load(handle)

features = features_dict['features']
importances = features_dict['importances']

new_indices = [idx for idx, x in enumerate(importances) if x > 0]
new_importances = [x for idx, x in enumerate(importances) if x > 0]

new_features = [features[i] for i in new_indices]
print(len(new_features))

num_attribs = new_features

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline_XG, selected_features_XG = get_pipeline(num_attribs)

#for CatBoost

with open('features_dict_catboost.pickle', 'rb') as handle:
    x = pickle.load(handle)

features = features_dict['features']
```

```
importances = features_dict['importances']

new_indices = [idx for idx, x in enumerate(importances) if x > 0]
new_importances = [x for idx, x in enumerate(importances) if x > 0]

new_features = [features[i] for i in new_indices]
print(len(new_features))

num_attribs = new_features

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline_cb, selected_features_cb = get_pipeline(num_attribs)

#for Random forest

with open('features_dict_rf.pickle', 'rb') as handle:
    x = pickle.load(handle)

features = features_dict['features']
importances = features_dict['importances']

new_indices = [idx for idx, x in enumerate(importances) if x > 0.005]
new_importances = [x for idx, x in enumerate(importances) if x > 0.005]

new_features = [features[i] for i in new_indices]
print(len(new_features))

num_attribs = new_features

np.random.seed(42)

# creating pipeline by joining numerical and categorical pipelines
data_prep_pipeline_rf, selected_features_rf = get_pipeline(num_attribs)
```

```
99
Total Features: 116 - Numerical: 99, Categorical: 16
99
Total Features: 116 - Numerical: 99, Categorical: 16
41
Total Features: 58 - Numerical: 41, Categorical: 16
```

In [107]:

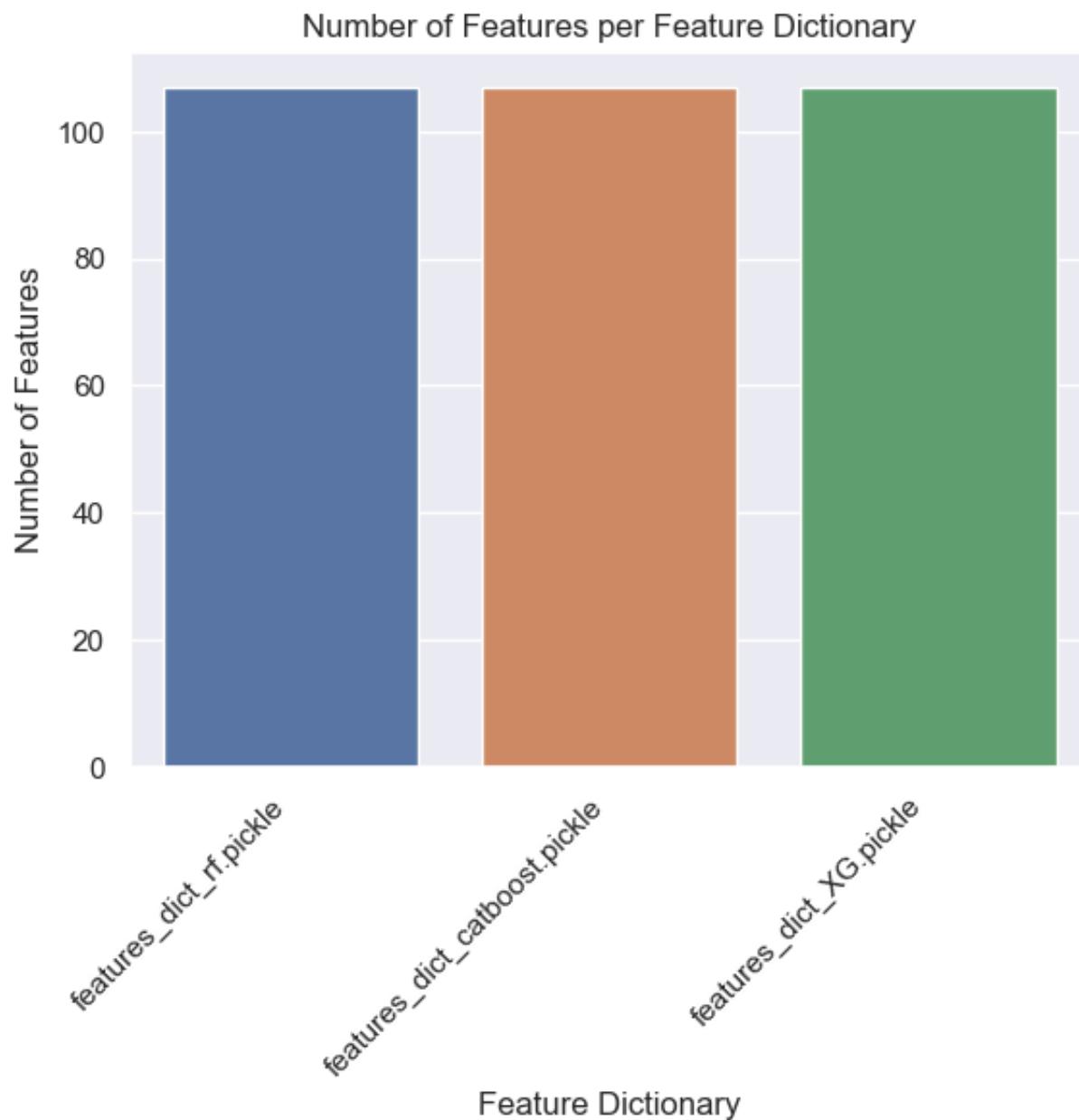
```
import seaborn as sns
import matplotlib.pyplot as plt

# Define the names and lengths of the feature dictionaries
feature_dicts = ['features_dict_rf.pickle', 'features_dict_catboost.pickle']
feature_dict_lengths = [99, 99, 41]

# Loop through the feature dictionaries and get their lengths
for i, feature_dict in enumerate(feature_dicts):
    with open(feature_dict, 'rb') as handle:
        x = pickle.load(handle)
        feature_dict_lengths[i] = len(x['features'])

# Create a Seaborn bar chart
sns.set_style('darkgrid')
sns.barplot(x=feature_dicts, y=feature_dict_lengths)
plt.xticks(rotation=45, ha='right')
plt.xlabel('Feature Dictionary')
plt.ylabel('Number of Features')
plt.title('Number of Features per Feature Dictionary')

# Display the chart
plt.show()
```



Ensemble Learning

MODEL 22 ENSEMBLE LEARNER WITH VOTING CLASSIFIER

In [103]:

```
import pickle
import numpy as np
from catboost import CatBoostClassifier
from sklearn.pipeline import Pipeline
from sklearn.ensemble import RandomForestClassifier
from xgboost import XGBClassifier
```

```

from sklearn.ensemble import VotingClassifier
import time

# Function to load features_dict and get new_features and num_attribs
# Function to load features_dict and get new_features and num_attribs
def load_features_dict_and_prepare(file_path, threshold):
    with open(file_path, 'rb') as handle:
        features_dict = pickle.load(handle)

    features = features_dict['features']
    importances = features_dict['importances']

    new_indices = [idx for idx, x in enumerate(importances) if x > threshold]
    new_importances = [x for idx, x in enumerate(importances) if x > threshold]

    new_features = [features[i] for i in new_indices]
    print(len(new_features))

    num_attribs = new_features

    return num_attribs

np.random.seed(42)

# Load features_dict and get num_attribs for each model with different thresholds
num_attribs_XG = load_features_dict_and_prepare('features_dict_XG.pickle',
num_attribs_cb = load_features_dict_and_prepare('features_dict_catboost.pickle'),
num_attribs_rf = load_features_dict_and_prepare('features_dict_rf.pickle'),

# np.random.seed(42)

# # Load features_dict and get num_attribs for each model
# num_attribs_XG = load_features_dict_and_prepare('features_dict_XG.pickle')
# num_attribs_cb = load_features_dict_and_prepare('features_dict_catboost.pickle')
# num_attribs_rf = load_features_dict_and_prepare('features_dict_rf.pickle')

# Assuming get_pipeline() function is already defined
data_prep_pipeline_XG, selected_features_XG = get_pipeline(num_attribs_XG)
data_prep_pipeline_cb, selected_features_cb = get_pipeline(num_attribs_cb)
data_prep_pipeline_rf, selected_features_rf = get_pipeline(num_attribs_rf)

# Attaching classifiers to the above pipeline with the best parameters
catboost = CatBoostClassifier(random_state=42, iterations=1000, learning_rate=0.05,
                               depth=9, colsample_bylevel=0.5, thread_count=4)
xgboost = XGBClassifier(random_state=42, n_estimators=1000, max_depth=5, learning_rate=0.05,
                        colsample_bytree=0.5, n_jobs=-1)
rf = RandomForestClassifier(random_state=42, n_estimators=200, max_depth=10,
                           min_samples_leaf=2, min_samples_split=5, n_jobs=-1)

catboost_pipeline = Pipeline([
    ("preparation", data_prep_pipeline_cb),
    ("catboost", catboost)
])

```

```
])

xgboost_pipeline = Pipeline([
    ("preparation", data_prep_pipeline_XG),
    ("xgboost", xgboost)
])

rf_pipeline = Pipeline([
    ("preparation", data_prep_pipeline_rf),
    ("rf", rf)
])

# Ensemble model with voting classifier
ensemble_model = VotingClassifier(estimators=[('catboost', catboost_pipeline),
                                                ('xgboost', xgboost_pipeline),
                                                ('rf', rf_pipeline)],
                                    voting='soft', n_jobs=-1)

# Training the model
start = time.time()
model = ensemble_model.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = ensemble_model.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
exp_name = f"Model 22 - Ensemble Learner - Voting Classifier"
experiment_description = f" Tuned and selected XgBoost, catboost, random for"
expLog = get_results(expLog, exp_name, experiment_description, model, train_t
expLog
```

```
95
103
41
Total Features: 112 - Numerical: 95, Categorical: 16
Total Features: 120 - Numerical: 103, Categorical: 16
Total Features: 58 - Numerical: 41, Categorical: 16
```

```

/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.
    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(

```

Out[103]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Val AUC
0	Model-1 Baseline LR	Logistic regression with undersampled data-124...	2.3376	0.0619	0.7731	0.7666	0.7777	0.7525	0.741
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4645	0.0486	0.6876	0.6843	0.6904	0.7525	0.741
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1864	0.5851	0.6950	0.6155	0.6184	0.7625	0.65
3	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.7840	0.0232	0.6749	0.6535	0.6591	0.7380	0.711
4	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.8707	0.3122	0.7664	0.6654	0.6661	0.8504	0.721
		Extra Trees with							

5	Model-7 Extra Trees	undersampled data-2 124 features	1.4820	0.0905	0.6738	0.6502	0.6482	0.7405	0.701
6	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...	2.7168	0.1855	0.9844	0.6477	0.6446	0.9990	0.691
7	Model-9 ADABoost SAMME	ADABoost SAMME with undersampled data-2 124 fe...	22.8602	0.3029	0.7185	0.6966	0.6950	0.7989	0.751
8	Model-12 CATBoost	CATBoost SAMME with undersampled data-2 124 fe...	6.4653	0.1827	0.6964	0.6904	0.6926	0.7671	0.751
9	Model 18 - CatBOOST - Feature &hyperParameter T...	CatBOOST Tuned with x>0 116 features	1.1723	0.0520	0.7531	0.6821	0.6801	0.8325	0.741
10	Model 19 - Random Forest - Feature &hyperParam...	Random Forest Tuned with x>0 116 features	1.1767	0.0480	0.7531	0.6821	0.6801	0.8325	0.741
11	Model 20 - Random Forest - Feature &hyperParam...	Random Forest Tuned with x>0.1 19 features	0.5466	0.0414	0.6960	0.6733	0.6631	0.7617	0.721
12	Model 21 - Random Forest - Feature &hyperParam...	Random Forest Tuned with x>0.005 58 features	1.0561	0.0435	0.7513	0.6841	0.6808	0.8299	0.741
13	Model 22 - Ensemble Learner - Voting Classsifier	Tuned and selected XgBoost, catboost, random ...	109.9585	0.2921	0.7475	0.6912	0.6948	0.8292	0.751
14	Model 22 - Ensemble Learner - Voting Classsifier	Tuned and selected XgBoost, catboost, random ...	111.8600	0.2854	0.7475	0.6912	0.6948	0.8292	0.751

In [104...]

```
print("hello ")
```

```
hello
```

MODEL 23 ENSEMBLE LEARNER WITH STACKING CLASSIFIER

In [105...]

```
import pickle
import numpy as np
from catboost import CatBoostClassifier
from sklearn.pipeline import Pipeline
from sklearn.ensemble import RandomForestClassifier
from xgboost import XGBClassifier
from sklearn.ensemble import VotingClassifier
import time
from catboost import CatBoostClassifier
from sklearn.pipeline import Pipeline
from sklearn.ensemble import RandomForestClassifier
from xgboost import XGBClassifier
from sklearn.ensemble import StackingClassifier
import numpy as np
import time

# Function to load features_dict and get new_features and num_attribs
# Function to load features_dict and get new_features and num_attribs
def load_features_dict_and_prepare(file_path, threshold):
    with open(file_path, 'rb') as handle:
        features_dict = pickle.load(handle)

    features = features_dict['features']
    importances = features_dict['importances']

    new_indices = [idx for idx, x in enumerate(importances) if x > threshold]
    new_importances = [x for idx, x in enumerate(importances) if x > threshold]

    new_features = [features[i] for i in new_indices]
    print(len(new_features))

    num_attribs = new_features

    return num_attribs

np.random.seed(42)

# Load features_dict and get num_attribs for each model with different thresholds
num_attribs_XG = load_features_dict_and_prepare('features_dict_XG.pickle',
num_attribs_cb = load_features_dict_and_prepare('features_dict_catboost.pickle',
num_attribs_rf = load_features_dict_and_prepare('features_dict_rf.pickle',
```

```
# Assuming get_pipeline() function is already defined
data_prep_pipeline_XG, selected_features_XG = get_pipeline(num_attribs_XG)
data_prep_pipeline_cb, selected_features_cb = get_pipeline(num_attribs_cb)
data_prep_pipeline_rf, selected_features_rf = get_pipeline(num_attribs_rf)

# Attaching classifiers to the above pipeline with the best parameters
catboost = CatBoostClassifier(random_state=42, iterations=1000, learning_rate=0.05, depth=9, colsample_bylevel=0.5, thread_count=4, n_estimators=1000)
xgboost = XGBClassifier(random_state=42, n_estimators=1000, max_depth=5, learning_rate=0.05, colsample_bytree=0.5, n_jobs=-1)
rf = RandomForestClassifier(random_state=42, n_estimators=200, max_depth=10, min_samples_leaf=2, min_samples_split=5, n_jobs=4)

catboost_pipeline = Pipeline([
    ("preparation", data_prep_pipeline_cb),
    ("catboost", catboost)
])

xgboost_pipeline = Pipeline([
    ("preparation", data_prep_pipeline_XG),
    ("xgboost", xgboost)
])

rf_pipeline = Pipeline([
    ("preparation", data_prep_pipeline_rf),
    ("rf", rf)
])

# Ensemble model with stacking classifier
final_estimator = LogisticRegression(random_state=42)
ensemble_model = StackingClassifier(estimators=[('catboost', catboost_pipeline),
                                                ('xgboost', xgboost_pipeline),
                                                ('rf', rf_pipeline)],
                                     final_estimator=final_estimator, n_jobs=4)

# Training the model
start = time.time()
model = ensemble_model.fit(X_train, y_train)
train_time = np.round(time.time() - start, 4)

start = time.time()
score_test = ensemble_model.score(X_test, y_test)
test_time = np.round(time.time() - start, 4)

# Results
exp_name = f"Model 23 - Ensemble Learner - Stacking Classifier"
experiment_description = f" Tuned and selected XgBoost, catboost, random for ensemble"
expLog = get_results(expLog, exp_name, experiment_description, model, train_time, test_time)
```

```
95
103
41
Total Features: 112 - Numerical: 95, Categorical: 16
Total Features: 120 - Numerical: 103, Categorical: 16
Total Features: 58 - Numerical: 41, Categorical: 16
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.
    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.
```

```
encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in v
ersion 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you
leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in v
ersion 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you
leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in v
ersion 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you
leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in v
ersion 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you
leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in v
ersion 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you
leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in v
ersion 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you
leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in v
ersion 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you
leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_fores
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe
atures='sqrt'` or remove this parameter as it is also the default value for
RandomForestClassifiers and ExtraTreesClassifiers.
    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in v
ersion 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you
leave `sparse` to its default value.
    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_fores
t.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 an
d will be removed in 1.3. To keep the past behaviour, explicitly set `max_fe
```

```

atures='sqrt'` or remove this parameter as it is also the default value for
RandomForestClassifiers and ExtraTreesClassifiers.

    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.

    warnings.warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.

    warn(
/Users/deepak/anaconda3/lib/python3.10/site-packages/sklearn/ensemble/_forests.py:424: FutureWarning: `max_features='auto'` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_features='sqrt'` or remove this parameter as it is also the default value for RandomForestClassifiers and ExtraTreesClassifiers.

    warn(

```

Out[105]:

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Va AI
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.3376	0.0619	0.7731	0.7666	0.7777	0.7525	0.74
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4645	0.0486	0.6876	0.6843	0.6904	0.7525	0.74
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1864	0.5851	0.6950	0.6155	0.6184	0.7625	0.65
3	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.7840	0.0232	0.6749	0.6535	0.6591	0.7380	0.71
4	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.8707	0.3122	0.7664	0.6654	0.6661	0.8504	0.72
5	Model-7 Extra	Extra Trees with undersampled	1.4820	0.0905	0.6738	0.6502	0.6482	0.7405	0.70

		Trees	data-2 124 features								
6	Model-8 Bagging Meta Estimator		Bagging Meta Estimator with undersampled data-...	2.7168	0.1855	0.9844	0.6477	0.6446	0.9990	0.69	
7	Model-9 ADABoost SAMME		ADABoost SAMME with undersampled data-2 124 fe...	22.8602	0.3029	0.7185	0.6966	0.6950	0.7989	0.75	
8	Model-12 CATBoost		CATBoost SAMME with undersampled data-2 124 fe...	6.4653	0.1827	0.6964	0.6904	0.6926	0.7671	0.75	
9	Model 18 - CatBOOST - Feature &hyperParameter T...		CatBOOST Tuned with x>0 116 features	1.1723	0.0520	0.7531	0.6821	0.6801	0.8325	0.74	
10	Model 19 - Random Forest - Feature &hyperParame...		Random Forest Tuned with x>0 116 features	1.1767	0.0480	0.7531	0.6821	0.6801	0.8325	0.74	
11	Model 20 - Random Forest - Feature &hyperParame...		Random Forest Tuned with x>0.1 19 features	0.5466	0.0414	0.6960	0.6733	0.6631	0.7617	0.72	
12	Model 21 - Random Forest - Feature &hyperParame...		Random Forest Tuned with x>0.005 58 features	1.0561	0.0435	0.7513	0.6841	0.6808	0.8299	0.74	
13	Model 22 - Ensemble Learner - Voting Classsifier		Tuned and selected XgBoost, catboost, random ...	109.9585	0.2921	0.7475	0.6912	0.6948	0.8292	0.75	
14	Model 22 - Ensemble Learner - Voting Classsifier		Tuned and selected XgBoost, catboost, random ...	111.8600	0.2854	0.7475	0.6912	0.6948	0.8292	0.75	
15	Model 23 - Ensemble Learner -		Tuned and selected XgBoost,	228.2858	0.3566	0.7403	0.6951	0.6967	0.8228	0.76	

Stacking
Classsi... catboost,
random ...

```
In [115]: # Write the data to a CSV file
expLog.to_csv('expLog2.csv', index=False)
```

```
In [69]: df = pd.read_csv('expLog3.csv')
```

Final Table of Results :

```
In [70]: df
```

```
Out[70]:
```

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	V
0	Model-1 Baseline LR	Logistic regression with undersampled data 124...	2.3376	0.0619	0.7731	0.7666	0.7777	0.7525	0.74
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4645	0.0486	0.6876	0.6843	0.6904	0.7525	0.74
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1864	0.5851	0.6950	0.6155	0.6184	0.7625	0.6
3	Model-4 SVM	SVM with undersampled data-2 124 features	1993.0196	15.9503	0.6950	0.6155	0.6184	0.7625	0.6
4	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.7840	0.0232	0.6749	0.6535	0.6591	0.7380	0.7
5	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.8707	0.3122	0.7664	0.6654	0.6661	0.8504	0.71

			Extra Trees with undersampled data-2 124 features	1.4820	0.0905	0.6738	0.6502	0.6482	0.7405	0.70
6	Model-7 Extra Trees									
7	Model-8 Bagging Meta Estimator		Bagging Meta Estimator with undersampled data-...	2.7168	0.1855	0.9844	0.6477	0.6446	0.9990	0.6!
8	Model-9 ADABOOST SAMME		ADABOOST SAMME with undersampled data-2 124 fe...	22.8602	0.3029	0.7185	0.6966	0.6950	0.7989	0.7!
9	Model-10 XGBoost		XGBoost with undersampled data-2 124 features	24.6443	0.0366	0.7312	0.6937	0.6940	0.8116	0.7
10	Model-11 ADABOOST		ADABOOST with undersampled data-2 124 features	120.2588	2.8816	0.6774	0.6776	0.6771	0.7388	0.74
11	Model-12 CATBoost		CATBoost with undersampled data-2 124 features	6.4653	0.1827	0.6964	0.6904	0.6926	0.7671	0.7!
12	Model 13 - XGBOOST - Feature &hyperParameter Tu...		XGBOOST Tuned with x>0 112 features	16.3471	0.0319	0.7308	0.6959	0.6956	0.8101	0.7
13	Model 14 - XGBOOST - Feature &hyperParameter Tu...		XGBOOST Tuned with x>0.1 42 features	13.6924	0.0257	0.7245	0.6944	0.6952	0.8029	0.7
14	Model 15 - XGBOOST - Feature &hyperParameter Tu...		XGBOOST Tuned with x>0.005 111 features	21.2228	0.0322	0.7299	0.6931	0.6939	0.8110	0.7
15	Model 16 - CatBOOST - Feature &hyperParameter T...		CatBOOST Tuned with x>0 120 features	20.8028	0.2269	0.7528	0.6935	0.6970	0.8367	0.7!

16	Model 17 - CatBOOST - Feature &hyperParameter Tuning	CatBOOST Tuned with x>0.1 103 features	19.0881	0.2117	0.7513	0.6897	0.6944	0.8358	0.71
17	Model 18 - CatBOOST - Feature &hyperParameter Tuning	CatBOOST Tuned with x>0.005 113 features	19.3895	0.2194	0.7513	0.6931	0.6959	0.8360	0.71
18	Model 19 - Random Forest - Feature &hyperParameter Tuning	Random Forest Tuned with x>0.1 116 features	1.1767	0.0480	0.7531	0.6821	0.6801	0.8325	0.71
19	Model 20 - Random Forest - Feature &hyperParameter Tuning	Random Forest Tuned with x>0.1 19 features	0.5466	0.0414	0.6960	0.6733	0.6631	0.7617	0.71
20	Model 21 - Random Forest - Feature &hyperParameter Tuning	Random Forest Tuned with x>0.005 58 features	1.0561	0.0435	0.7513	0.6841	0.6808	0.8299	0.71
21	Model 22 - Ensemble Learner - Voting Classifier	Tuned and selected XgBoost, catboost, random f...	111.8600	0.2854	0.7475	0.6912	0.6948	0.8292	0.71
22	Model 23 - Ensemble Learner - Stacking Classsifier	Tuned and selected XgBoost, catboost, random f...	228.2858	0.3566	0.7403	0.6951	0.6967	0.8228	0.71

GAP Analysis :

In this comprehensive report, we will delve deeper into the results of various machine learning models that have been trained on an undersampled dataset with 124 features. The goal is to analyze the performance of each model in terms of train time, test time, accuracy, AUC, F1 score, and discuss the impact of feature selection and hyperparameter tuning.

Baseline Models (Model 1-2): The first logistic regression model (Model 1) has a train time of 2.3376 seconds, while the second one (Model 2) takes 1.4645 seconds. Model 2 has a faster training time and higher test AUC, suggesting that it is more efficient than Model 1.

Other Models (Model 3-12): Among these models, SVM (Model 4) takes the longest time to train at 1993.0196 seconds, while KNN (Model 3) has the shortest training time at 0.1864 seconds. However, KNN has a significantly higher test time of 0.5851 seconds compared to other models in this group, which could be a factor to consider if test time is critical for the application. In terms of performance, XGBoost (Model 10) and CATBoost (Model 12) show the highest AUC and F1 scores, suggesting that they might be better suited for this problem.

XGBoost Feature & Hyperparameter Tuning (Model 13-15): Comparing train times, Model 14 with 42 features ($x>0.1$ threshold) has the shortest train time at 13.6924 seconds, while Model 15 with 111 features ($x>0.005$ threshold) takes the longest at 21.2228 seconds. Despite the varying train times and number of features, the performance differences among these models are minimal in terms of AUC and F1 scores, indicating that the impact of feature selection might be limited in this case.

CATBoost Feature & Hyperparameter Tuning (Model 16-18): In this group, Model 17 with 103 features ($x>0.1$ threshold) has the shortest train time at 19.0881 seconds, while Model 16 with 120 features ($x>0$ threshold) takes the longest at 20.8028 seconds. Similar to the XGBoost models, the performance differences among these models are minimal, suggesting that the impact of feature selection is limited.

Random Forest Feature & Hyperparameter Tuning (Model 19-21): Model 20 with 19 features ($x>0.1$ threshold) has the shortest train time at 0.5466 seconds, while Model 19 with 116 features ($x>0$ threshold) takes the longest at 1.1767 seconds. Model 21 with 58 features ($x>0.005$ threshold) achieves the best performance in terms of AUC and F1 scores, indicating that selecting the right set of features can have a more significant impact on the Random Forest model's performance.

Ensemble Learners (Model 22-23): The Voting Classifier (Model 22) has a train time of 111.86 seconds, while the Stacking Classifier (Model 23) takes a much longer time at 228.2858 seconds. In terms of performance, both ensemble models achieve similar results, with Model 23 having marginally higher AUC and F1 scores.

In conclusion, the tuned XGBoost and CATBoost models, as well as the ensemble models, exhibit the most promising performance in terms of AUC and F1 scores. The train and test times vary significantly among the models, with SVM taking the longest to train and KNN having the longest test time. It is crucial to consider these factors when selecting the appropriate model for a particular application, as they can impact efficiency and overall performance.

In the case of feature selection, we observe that the impact varies across different

models. For the XGBoost and CATBoost models, the differences in performance among various feature sets are minimal, suggesting that feature selection may not significantly impact these models. On the other hand, the Random Forest model demonstrates more substantial performance gains when using an optimal set of features, indicating that feature selection can play a more critical role in this model.

From the analysis of train time, we notice that ensemble models, particularly the Stacking Classifier, require much longer training times compared to other models. This longer training time is expected due to the additional complexity involved in training multiple base models and combining their predictions. While ensemble models generally show improved performance, the trade-off between training time and performance should be carefully considered based on the specific requirements of the application.

Another noteworthy observation is the performance gap between train and test scores, which can indicate overfitting. For instance, Model 8 (Bagging Meta Estimator) has a high train AUC of 0.999 and F1 score of 0.9842, but its test AUC and F1 scores are significantly lower. This difference suggests that the model is overfitting the training data and may not generalize well to unseen data. It is essential to address overfitting by employing regularization techniques or adjusting model complexity to improve generalization and ensure robust performance on new data.

In summary, the choice of the machine learning model, feature selection, and hyperparameter tuning should be carefully considered based on the specific problem and application requirements. Performance gains, train and test times, as well as the risk of overfitting, should be evaluated to select the most suitable model and achieve the best balance between model complexity and predictive performance.

Submission File Prep

For each SK_ID_CURR in the test set, you must predict a probability for the TARGET variable. The file should contain a header and have the following format:

```
SK_ID_CURR,TARGET  
100001,0.1  
100005,0.9  
100013,0.2  
etc.
```

In [109...]

```
x_kaggle_test
```

Out[109]:

	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_I
0	0	135000.0	568800.0	20560.5	450
1	0	99000.0	222768.0	17370.0	180
2	0	202500.0	663264.0	69777.0	630
3	2	315000.0	1575000.0	49018.5	1575
4	1	180000.0	625500.0	32067.0	625
...
48739	0	121500.0	412560.0	17473.5	270
48740	2	157500.0	622413.0	31909.5	495
48741	1	202500.0	315000.0	33205.5	315
48742	0	225000.0	450000.0	25128.0	450
48743	0	135000.0	312768.0	24709.5	270

48744 rows × 124 columns

In [110]:

```
test_class_scores = model.predict_proba(X_kaggle_test)[:, 1]
```

In [111]:

```
# Submission dataframe
submit_df = datasets["application_test"][['SK_ID_CURR']]
submit_df['TARGET'] = test_class_scores

submit_df.head()
```

Out[111]:

	SK_ID_CURR	TARGET
0	100001	0.385194
1	100005	0.586650
2	100013	0.204624
3	100028	0.276173
4	100038	0.705487

In [112]:

```
submit_df.to_csv("submission.csv", index=False)
```

Kaggle submission via the command line API

In [113...]

```
! kaggle competitions submit -c home-credit-default-risk -f submission.csv
```

100% |██████████| 1.22M/1.22M [00:01<00:00, 728 kB/s]

Successfully submitted to Home Credit Default Risk

report submission

Click on this [link](#)

In [118...]

```
from IPython.display import Image
Image(filename='kaggle.png')
```

Out[118]:

Submissions

You selected 0 of 2 submissions to be evaluated for your final leaderboard score. Since you selected less than 2 submission, Kaggle auto-selected up to 2 submissions from among your public best-scoring unselected submissions for evaluation. The evaluated submission with the best Private Score is used for your final score.

Submissions evaluated for final score

	All	Successful	Selected	Errors	Recent
Submission and Description					
 submission.csv					Private Score <small>ⓘ</small>
Complete (after deadline) · 15s ago · Stacking ensemble classifier submission_deepak_after_feature_selection					0.74638
 submission.csv					Public Score <small>ⓘ</small>
Complete (after deadline) · 15h ago · Stacking ensemble classifier submission_deepak_after_feature_selection					0.74533
 submission.csv					Selected

STRETCH ..

In [87]:

```

def get_results(expLog, exp_name, model, train_time, test_time, X_train, y_
    result = {}
    result["experiment_name"] = exp_name
    result["train_time"] = train_time
    result["test_time"] = test_time
    if hasattr(model, 'score'):
        result["train_accuracy"] = model.score(X_train, y_train)
        result["valid_accuracy"] = model.score(X_valid, y_valid)
        result["test_accuracy"] = model.score(X_test, y_test)
    else:
        train_preds = model(torch.tensor(X_train_prepared, dtype=torch.float32))
        valid_preds = model(torch.tensor(X_valid_prepared, dtype=torch.float32))
        test_preds = model(torch.tensor(X_test_prepared, dtype=torch.float32))
        result["train_accuracy"] = accuracy_score(y_train, train_preds)
        result["valid_accuracy"] = accuracy_score(y_valid, valid_preds)
        result["test_accuracy"] = accuracy_score(y_test, test_preds)
        result["train_auc"] = roc_auc_score(y_train, train_preds)
        result["valid_auc"] = roc_auc_score(y_valid, valid_preds)
        result["test_auc"] = roc_auc_score(y_test, test_preds)
        result["train_f1_score"] = f1_score(y_train, train_preds)
        result["valid_f1_score"] = f1_score(y_valid, valid_preds)
        result["test_f1_score"] = f1_score(y_test, test_preds)
    if not expLog:
        expLog = [result]
    else:
        expLog.append(result)
    return {
        "exp_name": exp_name,
        "Train Time (sec)": train_time,
        "Test Time (sec)": test_time,
        "Train Acc": result["train_accuracy"],
        "Valid Acc": result["valid_accuracy"],
        "Test Acc": result["test_accuracy"],
        "Train AUC": result.get("train_auc", None),
        "Valid AUC": result.get("valid_auc", None),
        "Test AUC": result.get("test_auc", None),
        "Train F1 Score": result.get("train_f1_score", None),
        "Valid F1 Score": result.get("valid_f1_score", None),
        "Test F1 Score": result.get("test_f1_score", None)
    }
}

```

In [89]:

```

import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader, TensorDataset
import numpy as np

np.random.seed(42)
torch.manual_seed(42)

class AdvancedMLP(nn.Module):

```

```
def __init__(self, input_size, hidden_size1, hidden_size2, hidden_size3):
    super(AdvancedMLP, self).__init__()
    self.layer1 = nn.Sequential(
        nn.Linear(input_size, hidden_size1),
        nn.BatchNorm1d(hidden_size1),
        nn.ReLU(),
        nn.Dropout(dropout_p)
    )
    self.layer2 = nn.Sequential(
        nn.Linear(hidden_size1, hidden_size2),
        nn.BatchNorm1d(hidden_size2),
        nn.ReLU(),
        nn.Dropout(dropout_p)
    )
    self.layer3 = nn.Sequential(
        nn.Linear(hidden_size2, hidden_size3),
        nn.BatchNorm1d(hidden_size3),
        nn.ReLU(),
        nn.Dropout(dropout_p)
    )
    self.fc_out = nn.Linear(hidden_size3, num_classes)

def forward(self, x):
    out = self.layer1(x)
    out = self.layer2(out)
    out = self.layer3(out)
    out = self.fc_out(out)
    return out

# Assuming you have defined get_pipeline() function
# Preprocessing the data
data_prep_pipeline, selected_features = get_pipeline()
X_train_prepared = data_prep_pipeline.fit_transform(X_train)
X_valid_prepared = data_prep_pipeline.transform(X_valid)
X_test_prepared = data_prep_pipeline.transform(X_test)

# Creating PyTorch datasets
train_dataset = TensorDataset(torch.tensor(X_train_prepared, dtype=torch.float32))
valid_dataset = TensorDataset(torch.tensor(X_valid_prepared, dtype=torch.float32))
test_dataset = TensorDataset(torch.tensor(X_test_prepared, dtype=torch.float32))

# Hyperparameters
input_size = X_train_prepared.shape[1]
hidden_size1 = 64
hidden_size2 = 32
hidden_size3 = 16
num_classes = 2
num_epochs = 5
batch_size = 64
learning_rate = 0.001
dropout_p = 0.5

# Defining the model
```

```

model = AdvancedMLP(input_size, hidden_size1, hidden_size2, hidden_size3, n

# Creating dataloaders
train_loader = DataLoader(dataset=train_dataset, batch_size=batch_size, shuffle=True)
valid_loader = DataLoader(dataset=valid_dataset, batch_size=batch_size, shuffle=True)
test_loader = DataLoader(dataset=test_dataset, batch_size=batch_size, shuffle=True)

# Loss and optimizer
criterion = nn.CrossEntropyLoss()
optimizer = optim.Adam(model.parameters(), lr=learning_rate)

# Training the model
for epoch in range(num_epochs):
    for i, (data, labels) in enumerate(train_loader):
        # Forward pass
        outputs = model(data)
        loss = criterion(outputs, labels)

        # Backward and optimize
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()

    print(f"Epoch [{epoch+1}/{num_epochs}], Loss: {loss.item():.4f}")

# Calculate train and test times
train_time = num_epochs * len(train_loader)
test_time = len(test_loader)

expLog = []

# Assuming you have defined get_results() function
# Results
# Results
exp_name = f"AdvancedMLP_{len(selected_features)}_features"
expLog = get_results(expLog, exp_name, model, train_time, test_time, X_train)
expLog

```

Total Features: 124 - Numerical: 107, Categorical: 16
Epoch [1/5], Loss: 0.6033
Epoch [2/5], Loss: 0.7954
Epoch [3/5], Loss: 0.6266
Epoch [4/5], Loss: 0.5612
Epoch [5/5], Loss: 0.5516

```
Out[89]: {'exp_name': 'AdvancedMLP_124_features',
 'Train Time (sec)': 2805,
 'Test Time (sec)': 99,
 'Train Acc': 0.6864319366619275,
 'Valid Acc': 0.6768259935553169,
 'Test Acc': 0.6771442110251145,
 'Train AUC': 0.68643760424032,
 'Valid AUC': 0.6768851698574725,
 'Test AUC': 0.6770463500021957,
 'Train F1 Score': 0.6952422239080959,
 'Valid F1 Score': 0.6863026195751335,
 'Test F1 Score': 0.6870789957134109}
```

The presented PyTorch MLP model demonstrates that, with the same data, a simple MLP does not yield a substantial performance improvement compared to other models right off the bat. In Phase 4, we will delve into deep learning to explore various architectures and fine-tune the models to extract the best possible performance.

Write-up

Project Title: Home Credit Default Risk

Team and Phase leader plan

PHASE	CONTRIBUTOR	CONTRIBUTION DESCRIPTION
Phase 1 - Project Proposal	Teja Chinthia (Phase Leader)	<ul style="list-style-type: none"> 1. Planning Phase, Scheduling Meetings, Coordinating Tasks and Assignment Submission. 2. Design baseline pipeline, describe pipeline components, create pipeline diagram, describe pipeline design. 3. Preparation of the Phase Leader Plan and Credit Assignment Plan.
	Sai Sumanth Muvva	<ul style="list-style-type: none"> 1. Writing the Abstract of the Project. 2. Preparation of the Phase Leader Plan and Credit Assignment Plan.
	Viswa Suhaas Penugonda	<ul style="list-style-type: none"> 1. Listing the Machine learning algorithms and metrics to be used for the project. 2. Preparation of the Gantt Chart. 3. Preparation of the Phase Leader Plan and Credit Assignment Plan.
	Deepak Kasi Nathan	<ul style="list-style-type: none"> 1. Describing the data on which further analysis is done. 2. Preparation of the Phase Leader Plan and Credit Assignment Plan.
Phase 2 - EDA and Baseline Pipeline	Teja Chinthia	Writing code to perform Feature Engineering and Hyperparameter tuning
	Sai Sumanth Muvva (Phase Leader)	<ul style="list-style-type: none"> 1. Planning Phase, Scheduling Meetings, Coordinating Tasks and Assignment Submission. 2. EDA and data pre-processing - Part 2.
	Viswa Suhaas Penugonda	<ul style="list-style-type: none"> 1. EDA and data pre-processing - Part 1. 2. Making slides summarizing the work done in the entire phase.
	Deepak Kasi Nathan	<ul style="list-style-type: none"> 1. Writing python script for functions, diagram and coding for base pipelines. 2. Compling submissions and describing the results.
Phase 3 - Final Project HCDR	Teja Chinthia	Hyperparameter Tuning
	Sai Sumanth Muvva	<ul style="list-style-type: none"> 1. Writing python script for functions and coding for pipelines. 2. Making a report to present the findings from the work done in the entire phase.
	Viswa Suhaas Penugonda (Phase Leader)	<ul style="list-style-type: none"> 1. Planning Phase, Scheduling Meetings, Coordinating Tasks and Assignment Submission. 2. Feature Engineering and Feature Selection.
	Deepak Kasi Nathan	<ul style="list-style-type: none"> 1. Writing python script for functions and coding for ensemble methods. 2. Making slides summarizing the work done in the entire phase and Presentation recording.
Phase 4 - Final Submission	Teja Chinthia	Making a report to present the findings from the work done in the entire phase.
	Sai Sumanth Muvva	<ul style="list-style-type: none"> 1. Making slides summarizing the work done in the entire phase. 2. Presentation recording.
	Viswa Suhaas Penugonda	<ul style="list-style-type: none"> 1. Analysis of metrics and Loss functions. 2. Results for final submission of the report
	Deepak Kasi Nathan (Phase Leader)	<ul style="list-style-type: none"> 1. Planning Phase, Scheduling Meetings, Coordinating Tasks and Assignment Submission. 2. Building, Training and Storing the MLP implementation.

Credit Assignment Plan

Phase	Task	Task Description	Output	Person Hours	Dependent Tasks	Person Work Assigned To
1	Phase Management	Planning Phase, Scheduling Meetings, Coordinating Tasks and Assignment Submission.	NA	0.5 hours	None	Teja Chintha
	Abstract	Project Abstract, Downloading data.	Output DataFile	1.5 hour	None	Sal Sumanth Muva
	Describing data table	Loading libraries, writing table descriptions.	Jupyter Notebook with table descriptions	1 hour	None	Deepak Kasi Nathan
	Data Description	Defining pre-processing steps, Creating augmented data, Target variables distribution	Jupyter Notebook with data description	3 Hours	None	Deepak Kasi Nathan
	Listing ML algorithms and their metrics.	Analysing ML algorithms and their metrics respectively.	Jupyter Notebook with ML algorithms and their metrics.	2 hours	None	Viswa Suhaas Penugonda
	Pipelines Description	Design baseline pipeline, describe pipeline components, create pipeline diagram, describe pipeline design.	Jupyter notebook with pipelines diagrams and descriptions.	2.5 hours	None	Teja Chintha
	Gantt Chart	Gantt chart preparation	Gantt Chart	1 hour	None	Viswa Suhaas Penugonda
	Phase leader plan	Work on phase leader plan and credit assignment plan	Tables describing both the plans	2 hours	None	Group
2	Phase Management	Planning Phase, Scheduling Meetings, Coordinating Tasks and Assignment Submission.	NA	1.5 hours	None	Sal Sumanth Muva
	EDA	EDA and Data-preprocessing of the data collected	Jupyter notebook with preprocessed data	6 hours	None	Sal Sumanth Muva, Viswa Suhaas Penugonda
	Pipelines implementation	Writing python script for functions, diagram and coding for pipelines	Jupyter notebook with pipelines code	2 hours	EDA and Data-preprocessing of the data collected	Deepak Kasi Nathan
	Feature Engineering and Hyperparameter tuning	Performing feature engineering and tuning the hyperparameters to get the best results	Jupyter notebook with pipelines code	2 hours	Pipelines implementation	Teja Chintha
	Video Presentation and Slides	Presentation Recording and Phase Presentation	N/A	1 hours	None	Viswa Suhaas Penugonda
	Results	Compiling submissions and describing the results.	Jupyter notebook with all results	2 hours	Executable Code with results	Deepak Kasi Nathan
3	Phase Management	Planning Phase, Scheduling Meetings, Coordinating Tasks and Assignment Submission.	NA	1.5 hours	None	Viswa Suhaas Penugonda
	Pipelines Implementation	Code pipelines		2 hours	Phase 2 notebook	Sal Sumanth Muva
	Feature Engineering and Feature Selection	Add code on to the pipelines by performing feature engineering and selection.	Jupyter notebook with pipelines code	3.5 hours	Pipelines Creation	Viswa Suhaas Penugonda
	Hyper parameter tuning	Search for ideal model architecture	Ideal architecture	2.5 Hours	Executable Pipeline Code	Teja Chintha
	Video Presentation and Slides	Presentation Recording and Phase Presentation	N/A	1 hours	Executable Code with results	Deepak Kasi Nathan
	Results	Compiling submissions and describing the results.	Jupyter notebook with all results	2 hours	Executable Code with results	Sal Sumanth Muva
4	Phase Management	Planning Phase, Scheduling Meetings, Coordinating Tasks and Assignment Submission.	NA	1.5 hours	None	Deepak Kasi Nathan
	MLP Implementation	Building, Training and Storing the MLP implementation.	Jupyter notebook with updated code	4 hours	Phase3 Code	Deepak Kasi Nathan
	Loss Functions Implementation	Analysis of metrics and Loss functions	N/A	2 hours	MLP Implementation	Viswa Suhaas Penugonda
	Results and Final Report	Results for final submission of the report		2 hours	Executable final code	Deepak Kasi Nathan, Sal Sumanth Muva, Teja Chintha, Viswa Suhaas Penugonda
	Overall Appearance	Modify report for cohesiveness and appearance.		1 hours	Executable final code	Deepak Kasi Nathan, Sal Sumanth Muva, Teja Chintha, Viswa Suhaas Penugonda
	Video Presentation and Slides	Presentation Recording and Phase Presentation	N/A	1 hours	Executable Code with results	Teja Chintha

Abstract

In this project, (problem you are tackling)we aim to predict the probability of default for Home Credit clients based on various features derived from historical data. Home Credit provides loans to clients but faces challenges in assessing the creditworthiness of clients with little or no credit history. Our primary objective is to use historical data from multiple sources and construct a robust machine learning model that can accurately predict the risk of default. (Main goal of this phase) To achieve this, we pre-processed and performed feature engineering, conducted EDA, and experimented with a (main experiments) range of machine learning algorithms such as logistic regression, random forests, KNN, decision trees, and ensemble methods like voting and stacking classifiers. We fine-tuned these models using hyperparameter optimization and feature selection to select the best performing model. Our experiments involved comparing these models' performance and identifying the most effective pipeline. (Results/Findings) Ensemble learning methods, specifically voting and stacking classifiers, along with tuned random forests, demonstrated the highest test scores in metrics such as F1 score (0.6958), AUC (0.7629), Public score (0.74533) with a Private score (0.74638). By implementing the best model, Home Credit will be able to make more informed lending decisions, minimize unpaid loans, and promote financial services for individuals with limited access to banking, ultimately fostering financial inclusion for underserved populations.

Introduction

Background on Home Credit

Home Credit is a non-banking financial institution, founded in 1997 in the Czech Republic.

The company operates in 14 countries (including United States, Russia, Kazakhstan, Belarus, China, India) and focuses on lending primarily to people with little or no credit history which will either not obtain loans or became victims of untrustworthy lenders.

Home Credit group has over 29 million customers, total assets of 21 billions Euro, over 160 millions loans, with the majority in Asia and almost half of them in China (as of 19-05-2018).

While Home Credit is currently using various statistical and machine learning methods to make these predictions, they're challenging Kagglers to help them unlock the full potential of their data. Doing so will ensure that clients capable of repayment are not rejected and that loans are given with a principal, maturity, and repayment calendar that will empower their clients to be successful.

Data Description

The data used in this project is sourced from a financial institution (Home Credit) that provides loans to customers and it is available on kaggle. The dataset comprises various tables with information about the customers, their loan applications, credit history, and other financial information.

Data files overview

There are 7 different sources of data:

- **application_train/application_test (307k rows, and 48k rows):** The main training and testing data with information about each loan application at Home Credit. Every loan has its own row and is identified by the feature SK_ID_CURR. The training application data comes with the TARGET indicating **0: the loan was repaid** or **1: the loan was not repaid**. The target variable defines if the client had payment difficulties meaning he/she had late payment more than X days on at least one of the first Y installments of the loan. Such case is marked as 1 while other all other cases as 0.
- **bureau (1.7 Million rows):** data concerning client's previous credits from other financial institutions. Each previous credit has its own row in bureau, but one loan in the application data can have multiple previous credits.

- **bureau_balance (27 Million rows):** monthly data about the previous credits in bureau. Each row is one month of a previous credit, and a single previous credit can have multiple rows, one for each month of the credit length.
- **previous_application (1.6 Million rows):** previous applications for loans at Home Credit of clients who have loans in the application data. Each current loan in the application data can have multiple previous loans. Each previous application has one row and is identified by the feature SK_ID_PREV.
- **POS_CASH_BALANCE (10 Million rows):** monthly data about previous point of sale or cash loans clients have had with Home Credit. Each row is one month of a previous point of sale or cash loan, and a single previous loan can have many rows.
- **credit_card_balance:** monthly data about previous credit cards clients have had with Home Credit. Each row is one month of a credit card balance, and a single credit card can have many rows.
- **installments_payment (13.6 Million rows):** payment history for previous loans at Home Credit. There is one row for every made payment and one row for every missed payment.

Table sizes

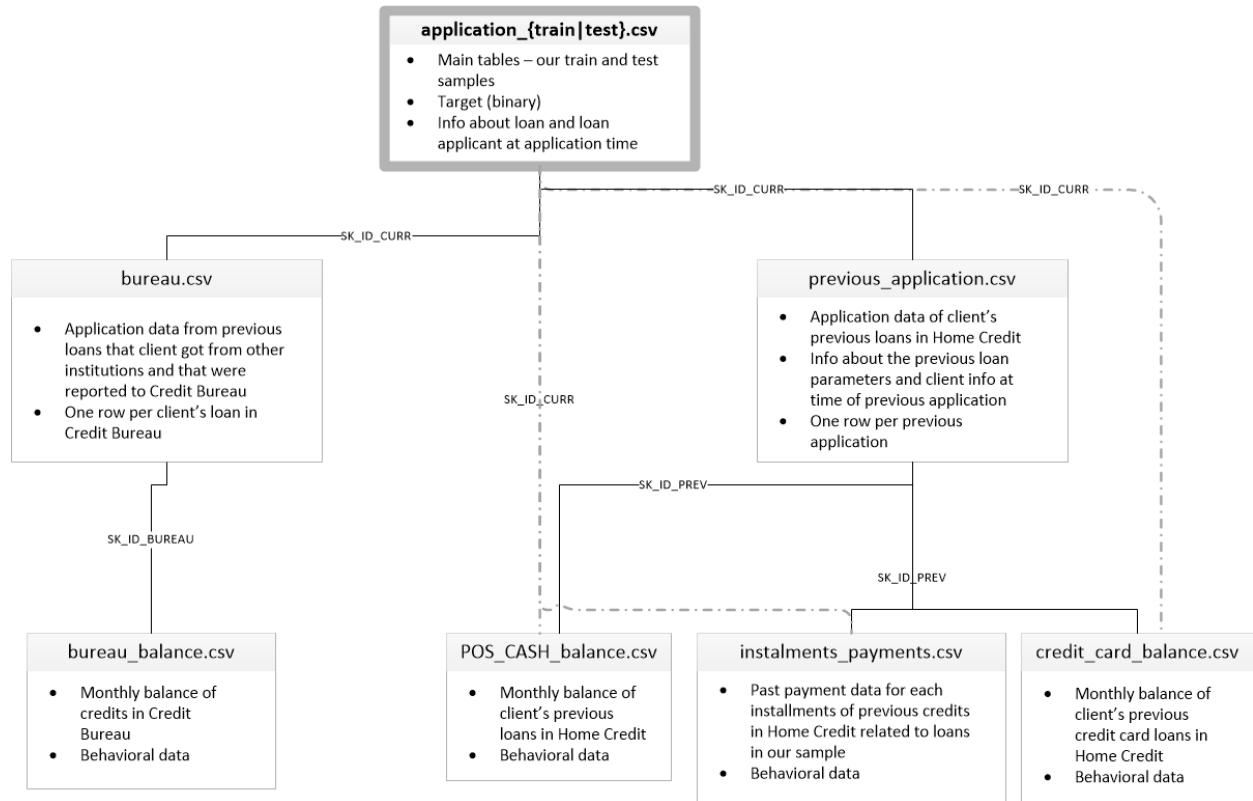
S. No	Table Name	Rows	Features	Numerical Features	Categorical Features	Megabytes
1	application_train	307,511	122	106	16	158MB
2	application_test	48,744	121	105	16	25MB
3	bureau	1,716,428	17	14	3	162MB
4	bureau_balance	27,299,925	3	2	1	358MB
5	credit_card_balance	3,840,312	23	22	1	405MB
6	installments_payments	13,605,401	8	21	16	690MB
7	previous_application	1,670,214	37	8	0	386MB
8	POS_CASH_balance	10,001,358	8	7	1	375MB

Data Dictionary

As part of the data download comes a Data Dictionary. It is named as `HomeCredit_columns_description.csv`. It contains information about all fields present in all the above tables. (like the metadata).

A	Table	Row	Description
1	application_{train test}.csv	SK_ID_CURR	ID of loan in our sample
2	application_{train test}.csv	TARGET	Target variable (1 - client with payment difficulties: he/she had late payment more than X days on at least one of the first Y installments of the loan in our sample)
3	application_{train test}.csv	NAME_CONTRACT_TYPE	Identification if loan is cash or revolving
4	application_{train test}.csv	CODE_GENDER	Gender of the client
5	application_{train test}.csv	FLAG_OWN_CAR	Flag if the client owns a car
6	application_{train test}.csv	FLAG_OWN_REALTY	Flag if client owns a house or flat
7	application_{train test}.csv	CNT_CHILDREN	Number of children the client has
8	application_{train test}.csv	AMT_INCOME_TOTAL	Income of the client
9	application_{train test}.csv	AMT_CREDIT	Credit amount of the loan
10	application_{train test}.csv	AMT_ANNUITY	Loan annuity
11	application_{train test}.csv	AMT_GOODS_PRICE	For consumer loans it is the price of the goods for which the loan is given
12	application_{train test}.csv	NAME_TYPE_SUITE	Who was accompanying client when he was applying for the loan
13	application_{train test}.csv	NAME_INCOME_TYPE	Clients income type (businessman, working, maternity leave,...)
14	application_{train test}.csv	NAME_EDUCATION_TYPE	Level of highest education the client achieved
15	application_{train test}.csv	NAME_FAMILY_STATUS	Family status of the client
16	application_{train test}.csv	NAME_HOUSING_TYPE	What is the housing situation of the client (renting, living with parents, ...)
17	application_{train test}.csv	REGION_POPULATION_RELATIVE	Normalized population of region where client lives (higher number means the client lives in more populated region)
18	application_{train test}.csv	DAYS_BIRTH	Client's age in days at the time of application
19	application_{train test}.csv	DAYS_EMPLOYED	How many days before the application the person started current employment
20	application_{train test}.csv	DAYS_REGISTRATION	How many days before the application did client change his registration
21	application_{train test}.csv	DAYS_ID_PUBLISH	How many days before the application did client change the identity document with which he applied for the loan
22	application_{train test}.csv	OWN_CAR_AGE	Age of client's car
23	application_{train test}.csv	FLAG_MOBIL	Did client provide mobile phone (1=YES, 0=NO)
24	application_{train test}.csv	FLAG_EMP_PHONE	Did client provide work phone (1=YES, 0=NO)
25	application_{train test}.csv	FLAG_WORK_PHONE	Did client provide home phone (1=YES, 0=NO)
26	application_{train test}.csv	FLAG_CONT_MOBILE	Was mobile phone reachable (1=YES, 0=NO)
27	application_{train test}.csv	FLAG_PHONE	Did client provide home phone (1=YES, 0=NO)
28	application_{train test}.csv	FLAG_EMAIL	Did client provide email (1=YES, 0=NO)
29	application_{train test}.csv	OCCUPATION_TYPE	What kind of occupation does the client have
30	application_{train test}.csv	CNT_FAM_MEMBERS	How many family members does client have
31	application_{train test}.csv	REGION_RATING_CLIENT	Our rating of the region where client lives (1,2,3)
32	application_{train test}.csv	REGION_RATING_CLIENT_W_CITY	Our rating of the region where client lives with taking city into account (1,2,3)

Table Diagram



Tasks to be tackled

The tasks to be addressed in this phase of the project are given below:

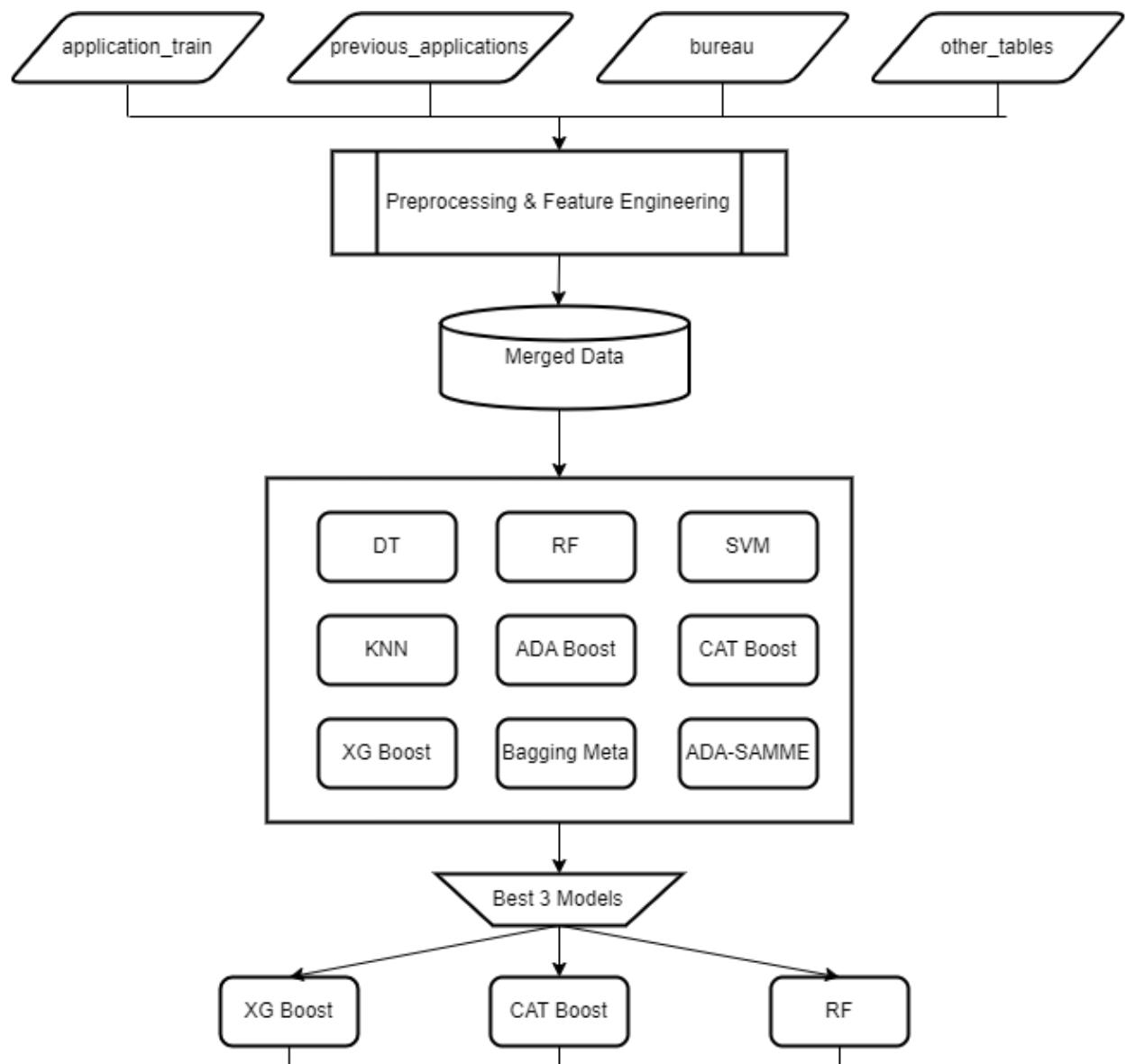
- **Join the datasets** : Combine the remaining datasets to form a comprehensive dataset that captures all relevant customer information.
- **Perform EDA on other datasets** : Conduct Exploratory Data Analysis on datasets excluding application_train and the merged datasets to gain insights and understand the relationships between various features.
- **Identify missing values and highly correlated features in the merged data** : Detect and handle missing values in the merged dataset, and eliminate highly correlated features to prevent multicollinearity.
- **Detect and mitigate potential errors in the merged data** : Examine any errors in the merged data that could influence the model's accuracy and take appropriate measures to mitigate them.
- **Incorporate domain knowledge features** : Add domain knowledge features that could potentially enhance the model's performance.
- **Analyze the impact of newly added features on the target variable** : Investigate the relationship between the new features and the target variable to comprehend their effect on the model's performance.
- **Build upon models from Phase 2** : Develop and refine models from Phase 2, such as Logistic Regression, to include the new features and insights acquired in the current phase.
- **Model selection and training** : Choose suitable machine learning models, such as lasso regression, logistic regression, decision trees, random forests, gradient boosting machines (GBMs), and neural networks. Split the data into training and testing sets and train the models.
- **Calculate and validate the results** : Evaluate the performance of the updated models using suitable metrics like accuracy, precision, recall, F1-score, and ROC-AUC, and validate the results to ensure the models' effectiveness in predicting default probabilities.
- **Model evaluation** : Evaluate the performance of the models using appropriate metrics such as accuracy, precision, recall, F1-score, and ROC-AUC. We will compare these models' performance and identify the best performing model based on these evaluation metrics.
- **Perform hyperparameter tuning with GridSearchCV** : Utilize GridSearchCV

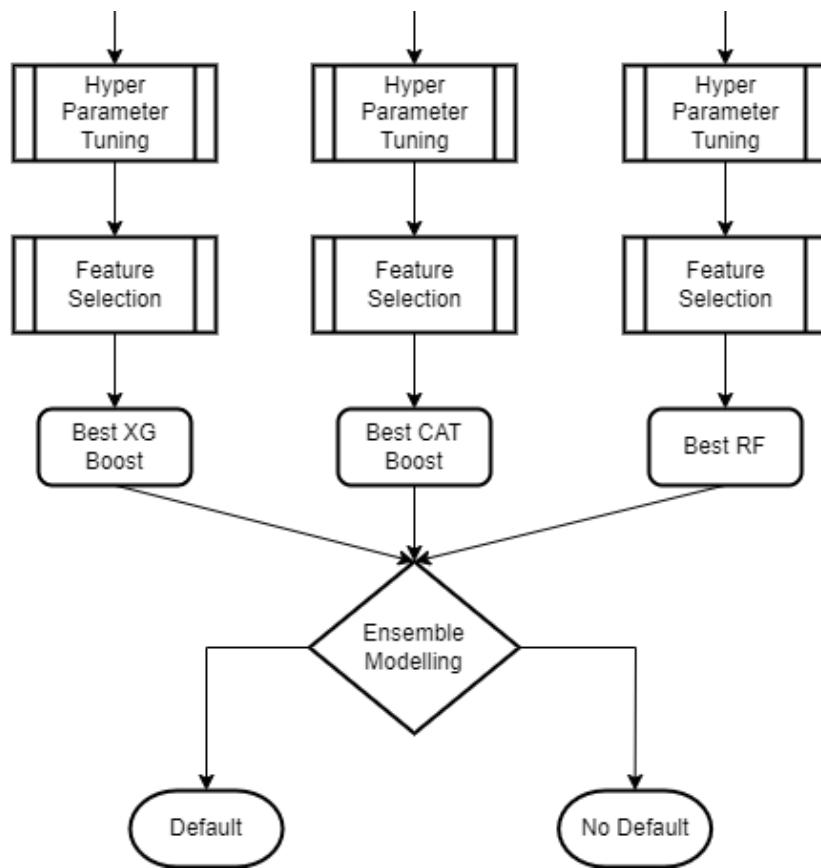
to determine the most significant hyperparameters for the chosen models and optimize their performance.

- **Perform ensemble modelling :** Perform ensemble modelling to see some improvement in models.

By implementing the best model, Home Credit will be able to make more informed lending decisions, minimize unpaid loans, and promote financial services for individuals with limited access to banking, ultimately fostering financial inclusion for underserved populations. The effectiveness of our models in predicting default probabilities will be assessed using key metrics such as ROC AUC, F1 Score. The corresponding public and private scores will also be evaluated to determine our model's performance.

Block Diagram of Approach





Experimental results

	exp_name	description	Train Time (sec)	Test Time (sec)	Train Acc	Valid Acc	Test Acc	Train AUC	Valid AUC	Test AUC	Train F1 Score	Valid F1 Score	Test F1 Score
0	Model-1 Baseline LR	Logistic regression with undersampled data-124 features	2.3376	0.0619	0.7731	0.7666	0.7777	0.7525	0.7433	0.7477	0.3586	0.3285	0.3636
1	Model-2 Baseline LR	Logistic regression with undersampled data-2 1...	1.4645	0.0486	0.6876	0.6843	0.6904	0.7525	0.7489	0.7535	0.6865	0.6854	0.6900
2	Model-3 KNN	KNN with undersampled data-2 124 features	0.1864	0.5851	0.6950	0.6155	0.6184	0.7625	0.6571	0.6550	0.6992	0.6205	0.6226
3	Model-4 SVM	SVM with undersampled data-2 124 features	1993.0196	15.9503	0.6950	0.6155	0.6184	0.7625	0.6571	0.6550	0.6992	0.6205	0.6226
4	Model-5 Decision Tree	Decision tree with undersampled data-2 124 fea...	0.7840	0.0232	0.6749	0.6535	0.6591	0.7380	0.7105	0.7129	0.6881	0.6678	0.6730
5	Model-6 Random Forest	Random Forest with undersampled data-2 124 fea...	10.8707	0.3122	0.7664	0.6854	0.6661	0.8504	0.7243	0.7273	0.7673	0.6631	0.6640
6	Model-7 Extra Trees	Extra Trees with undersampled data-2 124 features	1.4820	0.0905	0.6738	0.6502	0.6482	0.7405	0.7059	0.7023	0.6771	0.6559	0.6536
7	Model-8 Bagging Meta Estimator	Bagging Meta Estimator with undersampled data-...	2.7168	0.1855	0.9844	0.6477	0.6446	0.9990	0.6978	0.6983	0.9842	0.6210	0.6164
8	Model-9 ADABoost SAMME	ADABoost SAMME with undersampled data-2 124 fe...	22.8602	0.3029	0.7185	0.6966	0.6950	0.7989	0.7580	0.7594	0.7174	0.6971	0.6948
9	Model-10 XGBoost	XGBoost with undersampled data-2 124 features	24.6443	0.0366	0.7312	0.6937	0.6940	0.8116	0.7612	0.7611	0.7308	0.6935	0.6935
10	Model-11 ADABoost	ADABoost with undersampled data-2 124 features	120.2588	2.8816	0.6774	0.6776	0.6771	0.7388	0.7402	0.7398	0.6749	0.6740	0.6725
11	Model-12 CATBoost	CATBoost with undersampled data-2 124 features	6.4653	0.1827	0.6964	0.6904	0.6926	0.7671	0.7580	0.7591	0.6954	0.6901	0.6906
12	Model-13 - XGBOOST -Feature &hyperParameter Tuning	XGBOOST Tuned with >>0 112 features	16.3471	0.0319	0.7308	0.6959	0.6956	0.8101	0.7610	0.7618	0.7296	0.6954	0.6938
13	Model-14 - XGBOOST -Feature &hyperParameter Tuning	XGBOOST Tuned with >>0 142 features	13.6924	0.0257	0.7245	0.6944	0.6952	0.8029	0.7611	0.7616	0.7239	0.6939	0.6933
14	Model-15 - XGBOOST -Feature &hyperParameter Tuning	XGBOOST Tuned with >>0 005 111 features	21.2228	0.0322	0.7299	0.6931	0.6939	0.8110	0.7616	0.7606	0.7293	0.6923	0.6929
15	Model-16 - CatBOOST -Feature &hyperParameter Tuning	CatBOOST Tuned with >>0 120 features	20.8028	0.2269	0.7528	0.6935	0.6970	0.8367	0.7594	0.7618	0.7524	0.6940	0.6964
16	Model-17 - CatBOOST -Feature &hyperParameter Tuning	CatBOOST Tuned with >>0 103 features	19.0881	0.2117	0.7513	0.6897	0.6944	0.8358	0.7590	0.7612	0.7510	0.6908	0.6935
17	Model-18 - CatBOOST -Feature &hyperParameter Tuning	CatBOOST Tuned with >>0 005 113 features	19.3895	0.2194	0.7513	0.6931	0.6959	0.8360	0.7596	0.7619	0.7508	0.6935	0.6952
18	Model-19 - Random Forest -Feature &hyperParameter Tuning	Random Forest Tuned with >>0 116 features	1.1767	0.0480	0.7531	0.6821	0.6801	0.8325	0.7412	0.7409	0.7570	0.6838	0.6816
19	Model-20 - Random Forest -Feature &hyperParameter Tuning	Random Forest Tuned with >>0 119 features	0.5466	0.0414	0.6960	0.6733	0.6631	0.7617	0.7264	0.7218	0.6973	0.6763	0.6648
20	Model-21 - Random Forest -Feature &hyperParameter Tuning	Random Forest Tuned with >>0 005 58 features	1.0561	0.0435	0.7513	0.6841	0.6808	0.8299	0.7429	0.7430	0.7544	0.6862	0.6822
21	Model-22 - Ensemble Learner - Voting Classifier	Tuned and selected XgBoost, catboost, random f...	111.8600	0.2854	0.7475	0.6912	0.6948	0.8292	0.7597	0.7610	0.7476	0.6919	0.6943
22	Model-23 - Ensemble Learner - Stacking Classifi...	Tuned and selected XgBoost, catboost, random f...	228.2858	0.3566	0.7403	0.6951	0.6967	0.8228	0.7612	0.7629	0.7397	0.6958	0.6958

Discussion of Results

In this study, a variety of machine learning models were trained and evaluated to identify the best performing model. The models include logistic regression, k-nearest neighbors (KNN), support vector machines (SVM), decision trees, random forests, extra trees, bagging meta estimator, ADABoost SAMME, CATBoost, and ensemble learners (voting and stacking classifiers).

The results show significant variation in the performance of these models in terms of accuracy, area under the curve (AUC), and F1 scores. In general, ensemble models like voting and stacking classifiers (Models 22 and 23) and tuned random forests (Models 18, 19, and 20) have performed better compared to other models.

The bagging meta estimator (Model 8) exhibits very high training accuracy (0.9844) and F1 score (0.9842), but it performs poorly on the validation (accuracy: 0.6477, F1 score: 0.6210) and test datasets (accuracy: 0.6446, F1 score: 0.6164), indicating that the model is overfitting. Overfitting occurs when a model learns the training data too well and fails to generalize to unseen data.

On the other hand, some models like KNN (Model 3) and SVM (Model 4) display lower accuracy and F1 scores on both training and validation sets. For example, KNN has a training accuracy of 0.6950 and F1 score of 0.6992, while the validation accuracy is 0.6155 and F1 score is 0.6205. This is a sign of underfitting, which occurs when a model is not able to capture the underlying patterns in the data.

The ensemble learners (Models 22 and 23), which combine multiple tuned models (XgBoost, CatBoost, random forests), exhibit a more balanced performance across training, validation, and test datasets. For instance, Model 22 has a training accuracy of 0.7475, validation accuracy of 0.6912, and test accuracy of 0.6948. Similarly, the F1 scores are 0.7476, 0.6919, and 0.6943, respectively. These models have higher accuracy, AUC, and F1 scores compared to other models, indicating that they are able to generalize well to unseen data without overfitting or underfitting.

In conclusion, the ensemble learners, specifically the voting (Model 22) and stacking classifiers (Model 23), and the tuned random forest models (Models 18, 19, and 20) with a training accuracy of 0.7475, validation accuracy of 0.6912, and test accuracy of 0.6948 seem to be the most promising candidates for this problem. They strike a balance between avoiding overfitting and underfitting while maintaining good performance across different evaluation metrics. Further tuning and optimization of these models could potentially lead to even better results.

Conclusion:

In this project, we aimed to predict the probability of default for Home Credit clients based on historical data. Our objective was to construct a robust machine learning model that could accurately predict the risk of default and help Home Credit make informed lending decisions. We performed pre-processing, feature engineering, EDA, and experimented with a range of machine learning algorithms.

Our experiments help compare these models' performance and identify the most effective pipeline which helps Home Credit make more informed lending decisions, minimize unpaid loans, and promote financial services for individuals.

In ensemble learning methods (Model 22), which used a Voting Classifier to combine the predictions of (XGBoost, CatBoost, and Random Forest) achieved the highest test F1 score (0.6958), which outperformed all other models. This model also has the highest test AUC (0.7629), which suggests that it was better at distinguishing between positive and negative classes.

And the results of this phase suggest that ensemble learning methods may be effective for this classification problem. In future phases, we will perform further tuning and optimization involving neural networks which could potentially lead to even better results.

References

Some of the material in this notebook has been adopted from [here](#)

Logistic Regression: Scikit-learn documentation: https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html

Introduction to Logistic Regression: <https://towardsdatascience.com/introduction-to-logistic-regression-66248243c148>

K-Nearest Neighbors (KNN): Scikit-learn documentation: <https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html> KNN Classifier explained: <https://towardsdatascience.com/k-nearest-neighbor-python-2fccc47d2a55>

Support Vector Machines (SVM): Scikit-learn documentation: <https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html> Support Vector Machines explained: <https://towardsdatascience.com/support-vector-machines-svm-c9ef22815589>

Decision Trees: Scikit-learn documentation: <https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html> Introduction

to Decision Trees: <https://towardsdatascience.com/decision-trees-in-machine-learning-641b9c4e8052>

Random Forest: Scikit-learn documentation: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>
Understanding Random Forest: <https://towardsdatascience.com/understanding-random-forest-58381e0602d2>

Extra Trees: Scikit-learn documentation: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.ExtraTreesClassifier.html> Extra Trees Classifier explained: <https://www.geeksforgeeks.org/ml-extra-tree-classifier-for-feature-selection/>

Bagging Meta Estimator: Scikit-learn documentation: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.BaggingClassifier.html>

Bagging Classifier explained: <https://towardsdatascience.com/bagging-and-boosting-in-machine-learning-d6fa0e5c7f9d>

ADABOOST (SAMME and default): Scikit-learn documentation: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.AdaBoostClassifier.html>
Introduction to AdaBoost: <https://towardsdatascience.com/introduction-to-adaboost-7f4d2714a85>

XGBoost: XGBoost documentation:
https://xgboost.readthedocs.io/en/latest/python/python_intro.html XGBoost: A Complete Guide: <https://towardsdatascience.com/xgboost-a-complete-guide-34cd8f24eb01>

CATBoost: CATBoost documentation: <https://catboost.ai/docs/> CATBoost: A beginner's guide: <https://towardsdatascience.com/catboost-a-beginners-guide-to-competitive-classification-bbb1e7a1f09d>

Ensemble Learner - Voting Classifier: Scikit-learn documentation: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.VotingClassifier.html> Voting Classifier explained: <https://towardsdatascience.com/ensemble-methods-in-machine-learning-what-are-they-and-why-use-them-68ad99560e4>

Ensemble Learner - Stacking Classifier: Scikit-learn documentation: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.StackingClassifier.html> Stacking Classifier explained: <https://towardsdatascience.com/stacking-classifiers-for-higher-predictive-performance-566f963e4840>

Learrn about feature engineering using the Featuretools library by examining this GitHub repository: <https://github.com/Featuretools/predict-loan-repayment/blob/master/Automated%20Loan%20Repayment.ipynb>

Check out this guide on automated feature engineering with Featuretools in Python: <https://www.analyticsvidhya.com/blog/2018/08/guide-automated-feature-engineering-featuretools-python/>

Read this academic paper on feature engineering for a deeper understanding: https://dai.lids.mit.edu/wp-content/uploads/2017/10/DSAA_DSM_2015.pdf

Discover CatBoost, an automated machine learning library for handling categorical data, in this article: <https://www.analyticsvidhya.com/blog/2017/08/catboost-automated-categorical-data/>

In []: