

In [3]:

```
import pandas as pd
import matplotlib.pyplot as plt
import re
import time
import warnings
import numpy as np
from nltk.corpus import stopwords
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
warnings.filterwarnings("ignore")
import sys
import os
from tqdm import tqdm
import sqlite3
from sqlalchemy import create_engine # database connection
import csv
warnings.filterwarnings("ignore")
import datetime as dt
from nltk.corpus import stopwords
from sklearn.decomposition import TruncatedSVD
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.manifold import TSNE
import seaborn as sns
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import confusion_matrix
from sklearn.metrics.classification import accuracy_score, log_loss
from sklearn.feature_extraction.text import TfidfVectorizer
from collections import Counter
from scipy.sparse import hstack
from sklearn.multiclass import OneVsRestClassifier
from sklearn.svm import SVC
from collections import Counter, defaultdict
from sklearn.calibration import CalibratedClassifierCV
from sklearn.naive_bayes import MultinomialNB
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
import math
from sklearn.metrics import normalized_mutual_info_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import SGDClassifier
from mlxtend.classifier import StackingClassifier
from sklearn import model_selection
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import precision_recall_curve, auc, roc_curve

import cufflinks as cf
cf.go_offline()
```

In [4]:

```
#training data
df_train = pd.read_csv("application_train.csv")
df_train.head()
```

Out[4]:

	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REAL
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

In [5]:

```
#test data
df_test = pd.read_csv("application_test.csv")
df_test.head()
```

Out[5]:

	SK_ID_CURR	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REAL
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

Data Preparation:

Feature Engineering of Application data:

In [53]:

```
# Flag to represent when Total income is greater than Credit
df_train['INCOME_GT_CREDIT_FLAG'] = df_train['AMT_INCOME_TOTAL'] > df_train['AMT_CREDIT']
# Column to represent Credit Income Percent
df_train['CREDIT_INCOME_PERCENT'] = df_train['AMT_CREDIT'] / df_train['AMT_INCOME_TOTAL']
# Column to represent Annuity Income percent
df_train['ANNUITY_INCOME_PERCENT'] = df_train['AMT_ANNUITY'] / df_train['AMT_INCOME_TOTAL']
# Column to represent Credit Term
df_train['CREDIT_TERM'] = df_train['AMT_CREDIT'] / df_train['AMT_ANNUITY']
# Column to represent Days Employed percent in his Life
df_train['DAYS_EMPLOYED_PERCENT'] = df_train['DAYS_EMPLOYED'] / df_train['DAYS_BIRTH']
# Shape of Application data
print('The shape of application data:', df_train.shape)
```

The shape of application data: (307511, 129)

Create per person house features (living area per person, number of floors per person, etc.)

In [54]:

```

df_train['house_person'] = 1
df_train['house_person'].loc[df_train['NAME_HOUSING_TYPE']=='With parents'] +=2
df_train['house_person'].loc[(df_train['NAME_FAMILY_STATUS']=='Married')|(df_train['NAME_FAM
df_train['house_person'].loc[df_train['AGE']<55] += df_train['CNT_CHILDREN']

house = [f_ for f_ in df_train.columns if ('AVG' in f_) | ('MEDI' in f_) | ('MODE' in f_) &
for f_ in house:
    if df_train[f_].dtype != 'object':
        print (f_)
        df_train[f+'_PP'] = df_train[f_]/df_train['house_person']

```

APARTMENTS_AVG
 BASEMENTAREA_AVG
 YEARS_BEGINEXPLUATATION_AVG
 YEARS_BUILD_AVG
 COMMONAREA_AVG
 ELEVATORS_AVG
 ENTRANCES_AVG
 FLOORSMAX_AVG
 FLOORSMIN_AVG
 LANDAREA_AVG
 LIVINGAPARTMENTS_AVG
 LIVINGAREA_AVG
 NONLIVINGAPARTMENTS_AVG
 NONLIVINGAREA_AVG
 APARTMENTS_MODE
 BASEMENTAREA_MODE
 COMMONAREA_MODE
 ELEVATORS_MODE
 ENTRANCES_MODE
 FLOORSMAX_MODE
 FLOORSMIN_MODE
 LANDAREA_MODE
 LIVINGAPARTMENTS_MODE
 LIVINGAREA_MODE
 NONLIVINGAPARTMENTS_MODE
 NONLIVINGAREA_MODE
 APARTMENTS_MEDI
 BASEMENTAREA_MEDI
 YEARS_BEGINEXPLUATATION_MEDI
 YEARS_BUILD_MEDI
 COMMONAREA_MEDI
 ELEVATORS_MEDI
 ENTRANCES_MEDI
 FLOORSMAX_MEDI
 FLOORSMIN_MEDI
 LANDAREA_MEDI
 LIVINGAPARTMENTS_MEDI
 LIVINGAREA_MEDI
 NONLIVINGAPARTMENTS_MEDI
 NONLIVINGAREA_MEDI
 TOTALAREA_MODE

In [55]:

```
df_train.head()
```

Out[55]:

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

5 rows × 7 columns

Using Bureau Data:

In [56]:

```
print('Reading the data....', end='')
bureau = pd.read_csv('bureau.csv')
print('done!!!')
print('The shape of data:', bureau.shape)
print('First 5 rows of data:')
bureau.head()
```

Reading the data....done!!!

The shape of data: (1716428, 17)

First 5 rows of data:

Out[56]:

	SK_ID_CURR	SK_ID_BUREAU	CREDIT_ACTIVE	CREDIT_CURRENCY	DAYS_CREDIT	CREDIT_TERM
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	

Joining Bureau data to Application data:

In [57]:

```
# Combining numerical features
grp = bureau.drop(['SK_ID_BUREAU'], axis = 1).groupby(by=['SK_ID_CURR']).mean().reset_index()
grp.columns = ['BUREAU_'+column if column != 'SK_ID_CURR' else column for column in grp.columns]
application_bureau = df_train.merge(grp, on='SK_ID_CURR', how='left')
application_bureau.update(application_bureau[grp.columns].fillna(0))

# Combining categorical features
bureau_categorical = pd.get_dummies(bureau.select_dtypes('object'))
bureau_categorical['SK_ID_CURR'] = bureau['SK_ID_CURR']
grp = bureau_categorical.groupby(by = ['SK_ID_CURR']).mean().reset_index()
grp.columns = ['BUREAU_'+column if column != 'SK_ID_CURR' else column for column in grp.columns]
application_bureau = application_bureau.merge(grp, on='SK_ID_CURR', how='left')
application_bureau.update(application_bureau[grp.columns].fillna(0))

# Shape of application and bureau data combined
print('The shape application and bureau data combined:', application_bureau.shape)
```

The shape application and bureau data combined: (307511, 206)

Feature Engineering of Bureau Data:

In [58]:

```
# Number of past loans per customer
grp = bureau.groupby(by = ['SK_ID_CURR'])['SK_ID_BUREAU'].count().reset_index().rename(columns={'SK_ID_BUREAU': 'BUREAU_LOAN_COUNT'})
application_bureau = application_bureau.merge(grp, on='SK_ID_CURR', how='left')
application_bureau['BUREAU_LOAN_COUNT'] = application_bureau['BUREAU_LOAN_COUNT'].fillna(0)

# Number of types of past loans per customer
grp = bureau[['SK_ID_CURR', 'CREDIT_TYPE']].groupby(by = ['SK_ID_CURR'])['CREDIT_TYPE'].nunique().reset_index().rename(columns={'CREDIT_TYPE': 'BUREAU_LOAN_TYPES'})
application_bureau = application_bureau.merge(grp, on='SK_ID_CURR', how='left')
application_bureau['BUREAU_LOAN_TYPES'] = application_bureau['BUREAU_LOAN_TYPES'].fillna(0)

# Debt over credit ratio
bureau['AMT_CREDIT_SUM'] = bureau['AMT_CREDIT_SUM'].fillna(0)
bureau['AMT_CREDIT_SUM_DEBT'] = bureau['AMT_CREDIT_SUM_DEBT'].fillna(0)
grp1 = bureau[['SK_ID_CURR', 'AMT_CREDIT_SUM']].groupby(by=['SK_ID_CURR'])['AMT_CREDIT_SUM'].sum().reset_index().rename(columns={'AMT_CREDIT_SUM': 'TOTAL_CREDIT_SUM'})
grp2 = bureau[['SK_ID_CURR', 'AMT_CREDIT_SUM_DEBT']].groupby(by=['SK_ID_CURR'])['AMT_CREDIT_SUM_DEBT'].sum().reset_index().rename(columns={'AMT_CREDIT_SUM_DEBT': 'TOTAL_CREDIT_SUM_DEBT'})
grp1['DEBT_CREDIT_RATIO'] = grp2['TOTAL_CREDIT_SUM_DEBT']/grp1['TOTAL_CREDIT_SUM']
del grp1['TOTAL_CREDIT_SUM']
application_bureau = application_bureau.merge(grp1, on='SK_ID_CURR', how='left')
application_bureau['DEBT_CREDIT_RATIO'] = application_bureau['DEBT_CREDIT_RATIO'].fillna(0)
application_bureau['DEBT_CREDIT_RATIO'] = application_bureau.replace([np.inf, -np.inf], 0)
application_bureau['DEBT_CREDIT_RATIO'] = pd.to_numeric(application_bureau['DEBT_CREDIT_RATIO'], errors='coerce')

# Overdue over debt ratio
bureau['AMT_CREDIT_SUM_OVERDUE'] = bureau['AMT_CREDIT_SUM_OVERDUE'].fillna(0)
bureau['AMT_CREDIT_SUM_DEBT'] = bureau['AMT_CREDIT_SUM_DEBT'].fillna(0)
grp1 = bureau[['SK_ID_CURR', 'AMT_CREDIT_SUM_OVERDUE']].groupby(by=['SK_ID_CURR'])['AMT_CREDIT_SUM_OVERDUE'].sum().reset_index().rename(columns={'AMT_CREDIT_SUM_OVERDUE': 'TOTAL_CUSTOMER_OVERDUE'})
grp2 = bureau[['SK_ID_CURR', 'AMT_CREDIT_SUM_DEBT']].groupby(by=['SK_ID_CURR'])['AMT_CREDIT_SUM_DEBT'].sum().reset_index().rename(columns={'AMT_CREDIT_SUM_DEBT': 'TOTAL_CUSTOMER_DEBT'})
grp1['OVERDUE_DEBT_RATIO'] = grp1['TOTAL_CUSTOMER_OVERDUE']/grp2['TOTAL_CUSTOMER_DEBT']
del grp1['TOTAL_CUSTOMER_OVERDUE']
application_bureau = application_bureau.merge(grp1, on='SK_ID_CURR', how='left')
application_bureau['OVERDUE_DEBT_RATIO'] = application_bureau['OVERDUE_DEBT_RATIO'].fillna(0)
application_bureau['OVERDUE_DEBT_RATIO'] = application_bureau.replace([np.inf, -np.inf], 0)
application_bureau['OVERDUE_DEBT_RATIO'] = pd.to_numeric(application_bureau['OVERDUE_DEBT_RATIO'], errors='coerce')
```

In [59]:

```
application_bureau.head()
```

Out[59]:

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

5 rows × 210 columns

Using Previous Application Data:

In [60]:

```
print('Reading the data....', end='')
previous_applicaton = pd.read_csv('previous_application.csv')
print('done!!!')
print('The shape of data:', previous_applicaton.shape)
print('First 5 rows of data:')
previous_applicaton.head()
```

Reading the data....done!!!
 The shape of data: (1670214, 37)
 First 5 rows of data:

Out[60]:

	SK_ID_PREV	SK_ID_CURR	NAME_CONTRACT_TYPE	AMT_ANNUITY	AMT_APPLICATION	AMT_DOWN_PAYMENT
0	2030495	271877	Consumer loans	1730.430	17145.0	0.0
1	2802425	108129	Cash loans	25188.615	607500.0	0.0
2	2523466	122040	Cash loans	15060.735	112500.0	0.0
3	2819243	176158	Cash loans	47041.335	450000.0	0.0
4	1784265	202054	Cash loans	31924.395	337500.0	0.0

5 rows × 37 columns

Joining Previous Application data to Application Bureau data:

In [61]:

```
# Number of previous applications per customer
grp = previous_applicaton[['SK_ID_CURR', 'SK_ID_PREV']].groupby(by=['SK_ID_CURR'])['SK_ID_PREV'].count()
application_bureau_prev = application_bureau.merge(grp, on=['SK_ID_CURR'], how='left')
application_bureau_prev['PREV_APP_COUNT'] = application_bureau_prev['PREV_APP_COUNT'].fillna(0)

# Combining numerical features
grp = previous_applicaton.drop('SK_ID_PREV', axis=1).groupby(by=['SK_ID_CURR']).mean().reset_index()
prev_columns = ['PREV_'+column if column != 'SK_ID_CURR' else column for column in grp.columns]
grp.columns = prev_columns
application_bureau_prev = application_bureau_prev.merge(grp, on=['SK_ID_CURR'], how='left')
application_bureau_prev.update(application_bureau_prev[grp.columns].fillna(0))

# Combining categorical features
prev_categorical = pd.get_dummies(previous_applicaton.select_dtypes('object'))
prev_categorical['SK_ID_CURR'] = previous_applicaton['SK_ID_CURR']
prev_categorical.head()
grp = prev_categorical.groupby('SK_ID_CURR').mean().reset_index()
grp.columns = ['PREV_'+column if column != 'SK_ID_CURR' else column for column in grp.columns]
application_bureau_prev = application_bureau_prev.merge(grp, on=['SK_ID_CURR'], how='left')
application_bureau_prev.update(application_bureau_prev[grp.columns].fillna(0))
```

In [62]:

```
application_bureau_prev.head()
```

Out[62]:

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

5 rows × 373 columns

In [63]:

```
application_bureau_prev.head(0)
```

Out[63]:

	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY
--	------------	--------	--------------------	-------------	--------------	-----------------

0 rows × 373 columns

Using POS_CASH_balance data:

In [64]:

```
print('Reading the data....', end='')
pos_cash = pd.read_csv('POS_CASH_balance.csv')
print('done!!!')
print('The shape of data:', pos_cash.shape)
print('First 5 rows of data:')
pos_cash.head()
```

Reading the data....done!!!
 The shape of data: (10001358, 8)
 First 5 rows of data:

Out[64]:

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

Joining POS_CASH_balance data to application_bureau_prev_data:

In [65]:

```
# Combining numerical features
grp = pos_cash.drop('SK_ID_PREV', axis=1).groupby(by=['SK_ID_CURR']).mean().reset_index()
prev_columns = ['POS_'+column if column != 'SK_ID_CURR' else column for column in grp.columns]
grp.columns = prev_columns
application_bureau_prev = application_bureau_prev.merge(grp, on=['SK_ID_CURR'], how='left')
application_bureau_prev.update(application_bureau_prev[grp.columns].fillna(0))

# Combining categorical features
pos_cash_categorical = pd.get_dummies(pos_cash.select_dtypes('object'))
pos_cash_categorical['SK_ID_CURR'] = pos_cash['SK_ID_CURR']
grp = pos_cash_categorical.groupby('SK_ID_CURR').mean().reset_index()
grp.columns = ['POS_'+column if column != 'SK_ID_CURR' else column for column in grp.columns]
application_bureau_prev = application_bureau_prev.merge(grp, on=['SK_ID_CURR'], how='left')
application_bureau_prev.update(application_bureau_prev[grp.columns].fillna(0))
```

Using installments_payments data:

In [66]:

```
print('Reading the data....', end='')
insta_payments = pd.read_csv('installments_payments.csv')
print('done!!!')
print('The shape of data:',insta_payments.shape)
print('First 5 rows of data:')
insta_payments.head()
```

Reading the data....done!!!
 The shape of data: (13605401, 8)
 First 5 rows of data:

Out[66]:

	SK_ID_PREV	SK_ID_CURR	NUM_INSTALLMENT_VERSION	NUM_INSTALLMENT_NUMBER	DAYS
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	

Joining Installments Payments data to application_bureau_prev_data:

In [67]:

```
# Combining numerical features and there are no categorical features in this dataset
grp = insta_payments.drop('SK_ID_PREV', axis =1).groupby(by=['SK_ID_CURR']).mean().reset_index()
prev_columns = ['INSTA_'+column if column != 'SK_ID_CURR' else column for column in grp.columns]
grp.columns = prev_columns
application_bureau_prev = application_bureau_prev.merge(grp, on =['SK_ID_CURR'], how = 'left')
application_bureau_prev.update(application_bureau_prev[grp.columns].fillna(0))
```

Using Credit card balance data:

In [68]:

```
print('Reading the data....', end='')
credit_card = pd.read_csv('credit_card_balance.csv')
print('done!!!')
print('The shape of data:', credit_card.shape)
print('First 5 rows of data:')
credit_card.head()
```

Reading the data....done!!!
 The shape of data: (3840312, 23)
 First 5 rows of data:

Out[68]:

	SK_ID_PREV	SK_ID_CURR	MONTHS_BALANCE	AMT_BALANCE	AMT_CREDIT_LIMIT_ACTUA
0	2562384	378907	-6	56.970	13500
1	2582071	363914	-1	63975.555	4500
2	1740877	371185	-7	31815.225	4500
3	1389973	337855	-4	236572.110	22500
4	1891521	126868	-1	453919.455	4500

5 rows × 23 columns

Joining Credit card balance data to application_bureau_prev data:

In [69]:

```
# Combining numerical features
grp = credit_card.drop('SK_ID_PREV', axis =1).groupby(by=['SK_ID_CURR']).mean().reset_index()
prev_columns = ['CREDIT_'+column if column != 'SK_ID_CURR' else column for column in grp.columns]
grp.columns = prev_columns
application_bureau_prev = application_bureau_prev.merge(grp, on =['SK_ID_CURR'], how = 'left')
application_bureau_prev.update(application_bureau_prev[grp.columns].fillna(0))

# Combining categorical features
credit_categorical = pd.get_dummies(credit_card.select_dtypes('object'))
credit_categorical['SK_ID_CURR'] = credit_card['SK_ID_CURR']
grp = credit_categorical.groupby('SK_ID_CURR').mean().reset_index()
grp.columns = ['CREDIT_'+column if column != 'SK_ID_CURR' else column for column in grp.columns]
application_bureau_prev = application_bureau_prev.merge(grp, on=['SK_ID_CURR'], how='left')
application_bureau_prev.update(application_bureau_prev[grp.columns].fillna(0))
```

Dividing final data into train, valid and test datasets:

In [70]:

```
y = application_bureau_prev.pop('TARGET').values
X_train, X_temp, y_train, y_temp = train_test_split(application_bureau_prev.drop(['SK_ID_CURR'], axis=1), y_train, y_temp = train_test_split(X_temp, y_temp, stratify = y_temp, test_size=0.2)
print('Shape of X_train:',X_train.shape)
print('Shape of X_val:',X_val.shape)
print('Shape of X_test:',X_test.shape)
```

Shape of X_train: (215257, 418)

Shape of X_val: (46127, 418)

Shape of X_test: (46127, 418)

Featurizing the data:

In [71]:

```

from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import OneHotEncoder
from sklearn.impute import SimpleImputer
# Seperation of columns into numeric and categorical columns
types = np.array([dt for dt in X_train.dtypes])
all_columns = X_train.columns.values
is_num = types != 'object'
num_cols = all_columns[is_num]
cat_cols = all_columns[~is_num]
# Featurization of numeric data
imputer_num = SimpleImputer(strategy='median')
X_train_num = imputer_num.fit_transform(X_train[num_cols])
X_val_num = imputer_num.transform(X_val[num_cols])
X_test_num = imputer_num.transform(X_test[num_cols])
scaler_num = StandardScaler()
X_train_num1 = scaler_num.fit_transform(X_train_num)
X_val_num1 = scaler_num.transform(X_val_num)
X_test_num1 = scaler_num.transform(X_test_num)
X_train_num_final = pd.DataFrame(X_train_num1, columns=num_cols)
X_val_num_final = pd.DataFrame(X_val_num1, columns=num_cols)
X_test_num_final = pd.DataFrame(X_test_num1, columns=num_cols)
# Featurization of categorical data
imputer_cat = SimpleImputer(strategy='constant', fill_value='MISSING')
X_train_cat = imputer_cat.fit_transform(X_train[cat_cols])
X_val_cat = imputer_cat.transform(X_val[cat_cols])
X_test_cat = imputer_cat.transform(X_test[cat_cols])
X_train_cat1 = pd.DataFrame(X_train_cat, columns=cat_cols)
X_val_cat1 = pd.DataFrame(X_val_cat, columns=cat_cols)
X_test_cat1 = pd.DataFrame(X_test_cat, columns=cat_cols)
ohe = OneHotEncoder(sparse=False, handle_unknown='ignore')
X_train_cat2 = ohe.fit_transform(X_train_cat1)
X_val_cat2 = ohe.transform(X_val_cat1)
X_test_cat2 = ohe.transform(X_test_cat1)
cat_cols_ohe = list(ohe.get_feature_names(input_features=cat_cols))
X_train_cat_final = pd.DataFrame(X_train_cat2, columns = cat_cols_ohe)
X_val_cat_final = pd.DataFrame(X_val_cat2, columns = cat_cols_ohe)
X_test_cat_final = pd.DataFrame(X_test_cat2, columns = cat_cols_ohe)
# Final complete data
X_train_final = pd.concat([X_train_num_final, X_train_cat_final], axis = 1)
X_val_final = pd.concat([X_val_num_final, X_val_cat_final], axis = 1)
X_test_final = pd.concat([X_test_num_final, X_test_cat_final], axis = 1)
print(X_train_final.shape)
print(X_val_final.shape)
print(X_test_final.shape)

```

(215257, 548)

(46127, 548)

(46127, 548)

Saving the files for future use:

In [72]:

```
# Saving the Dataframes into CSV files for future use
X_train_final.to_csv('X_train_final.csv')
X_val_final.to_csv('X_val_final.csv')
X_test_final.to_csv('X_test_final.csv')
# Saving the numpy arrays into text files for future use
np.savetxt('y.txt', y)
np.savetxt('y_train.txt', y_train)
np.savetxt('y_val.txt', y_val)
np.savetxt('y_test.txt', y_test)
```

Selection of features:

In [73]:

```
import lightgbm as lgb
model_sk = lgb.LGBMClassifier(boosting_type='gbdt', max_depth=7, learning_rate=0.01, n_estimators=100,
                             class_weight='balanced', subsample=0.9, colsample_bytree=0.8, n_jobs=-1)
train_features, valid_features, train_y, valid_y = train_test_split(X_train_final, y_train,
                                                                    model_sk.fit(train_features, train_y, early_stopping_rounds=100, eval_set = [(valid_features, valid_y)]))
```

Training until validation scores don't improve for 100 rounds

```
[200] valid_0's auc: 0.755043 valid_0's binary_logloss: 0.59233
[400] valid_0's auc: 0.7688 valid_0's binary_logloss: 0.566255
[600] valid_0's auc: 0.774508 valid_0's binary_logloss: 0.551756
[800] valid_0's auc: 0.776809 valid_0's binary_logloss: 0.542226
[1000] valid_0's auc: 0.777925 valid_0's binary_logloss: 0.534798
[1200] valid_0's auc: 0.778786 valid_0's binary_logloss: 0.528271
[1400] valid_0's auc: 0.77892 valid_0's binary_logloss: 0.522623
Early stopping, best iteration is:
[1366] valid_0's auc: 0.778981 valid_0's binary_logloss: 0.523538
```

Out[73]:

```
LGBMClassifier(boosting_type='gbdt', class_weight='balanced',
               colsample_bytree=0.8, importance_type='split',
               learning_rate=0.01, max_depth=7, min_child_samples=20,
               min_child_weight=0.001, min_split_gain=0.0, n_estimators=200,
               n_jobs=-1, num_leaves=31, objective=None, random_state=None,
               reg_alpha=0.0, reg_lambda=0.0, silent=True, subsample=0.9,
               subsample_for_bin=200000, subsample_freq=0)
```

In [74]:

```
import pickle
feature_imp = pd.DataFrame(sorted(zip(model_sk.feature_importances_, X_train_final.columns))
                           features_df = feature_imp.sort_values(by="Value", ascending=False)
                           selected_features = list(features_df[features_df['Value'] >= 50]['Feature'])
# Saving the selected features into pickle file
with open('select_features.txt', 'wb') as fp:
    pickle.dump(selected_features, fp)
print('The no. of features selected:', len(selected_features))
```

The no. of features selected: 189

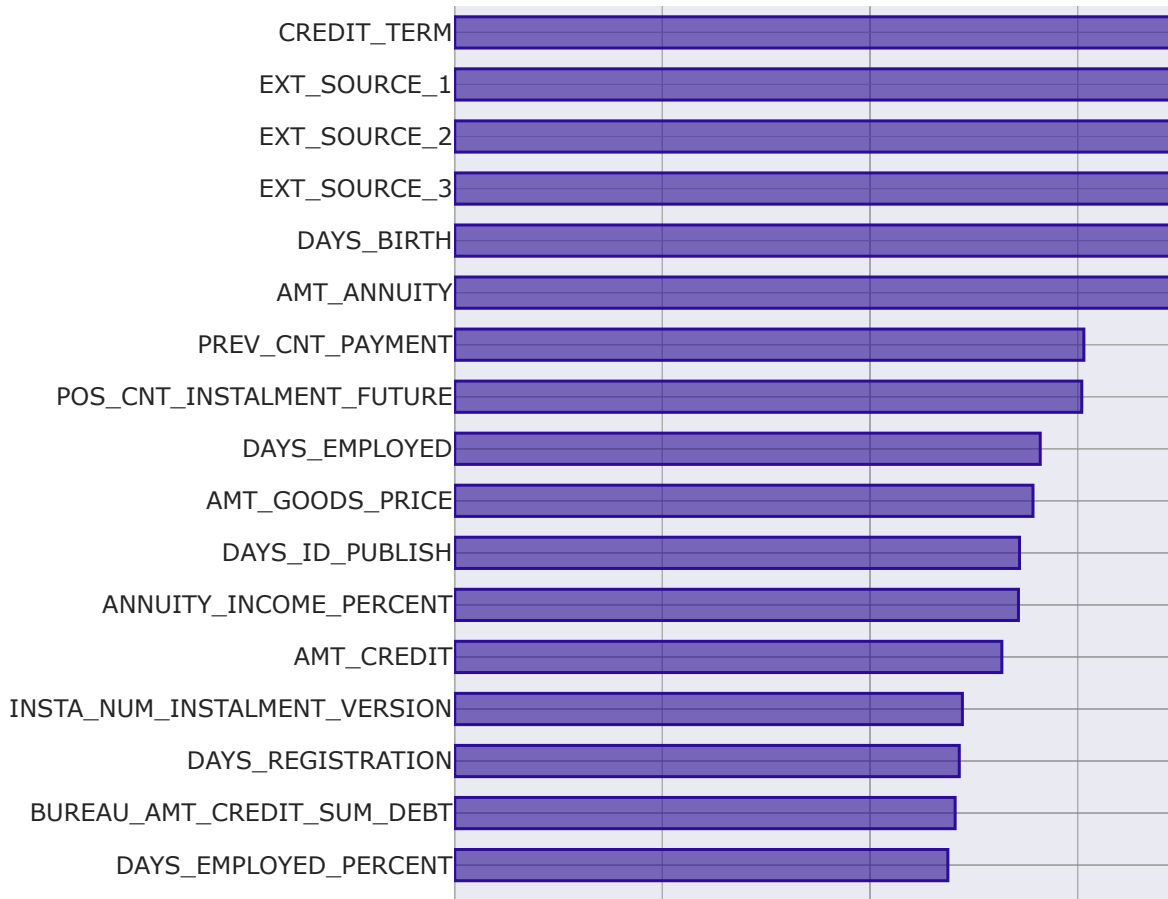
In [75]:

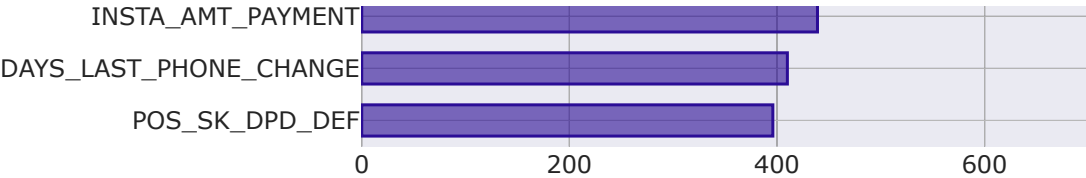
```

# Feature importance Plot
data1 = features_df.head(20)
data = [go.Bar(x = data1.sort_values(by='Value')['Value'] , y = data1.sort_values(by='Value')
              marker = dict(
                color = 'rgba(43, 13, 150, 0.6)',
                line = dict(
                  color = 'rgba(43, 13, 150, 1.0)',
                  width = 1.5)
            )) ]
layout = go.Layout(
    autosize=False,
    width=1300,
    height=700,
    title = "Top 20 important features",
    xaxis=dict(
        title='Importance value'
    ),
    yaxis=dict(
        automargin=True
    ),
    bargap=0.4
)
fig = go.Figure(data = data, layout=layout)
fig.layout.template = 'seaborn'
py.iplot(fig)

```

Top 20





Machine Learning Models:

In [76]:

```

def plot_confusion_matrix(test_y, predicted_y):
    # Confusion matrix
    C = confusion_matrix(test_y, predicted_y)

    # Recall matrix
    A = (((C.T)/(C.sum(axis=1))).T)

    # Precision matrix
    B = (C/C.sum(axis=0))

    plt.figure(figsize=(20,4))

    labels = ['Re-paid(0)', 'Not Re-paid(1)']
    cmap=sns.light_palette("purple")
    plt.subplot(1,3,1)
    sns.heatmap(C, annot=True, cmap=cmap, fmt="d", xticklabels = labels, yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title('Confusion matrix')

    plt.subplot(1,3,2)
    sns.heatmap(A, annot=True, cmap=cmap, xticklabels = labels, yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title('Recall matrix')

    plt.subplot(1,3,3)
    sns.heatmap(B, annot=True, cmap=cmap, xticklabels = labels, yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title('Precision matrix')

    plt.show()
def cv_plot(alpha, cv_auc):
    fig, ax = plt.subplots()
    ax.plot(np.log10(alpha), cv_auc, c='g')
    for i, txt in enumerate(np.round(cv_auc,3)):
        ax.annotate((alpha[i],str(txt)), (np.log10(alpha[i]),cv_auc[i]))
    plt.grid()
    plt.xticks(np.log10(alpha))
    plt.title("Cross Validation Error for each alpha")
    plt.xlabel("Alpha i's")
    plt.ylabel("Error measure")
    plt.show()

```

Logistic regression with selected features:

In [181]:

```

from sklearn.metrics import roc_auc_score
alpha = np.logspace(-4,4,9)
cv_auc_score = []
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l1', class_weight = 'balanced', loss='log', random
    clf.fit(X_train_final[selected_features], y_train)
    sig_clf = CalibratedClassifierCV(clf, method='sigmoid')
    sig_clf.fit(X_train_final[selected_features], y_train)
    y_pred_prob = sig_clf.predict_proba(X_val_final[selected_features])[:,1]
    cv_auc_score.append(roc_auc_score(y_val,y_pred_prob))
    print('For alpha {0}, cross validation AUC score {1}'.format(i,roc_auc_score(y_val,y_pr
cv_plot(alpha, cv_auc_score)
print('The Optimal C value is:', alpha[np.argmax(cv_auc_score)])

```

For alpha 0.0001, cross validation AUC score 0.7570320665909553

For alpha 0.001, cross validation AUC score 0.7567991219639147

For alpha 0.01, cross validation AUC score 0.7369149954506644

For alpha 0.1, cross validation AUC score 0.7021180843582268

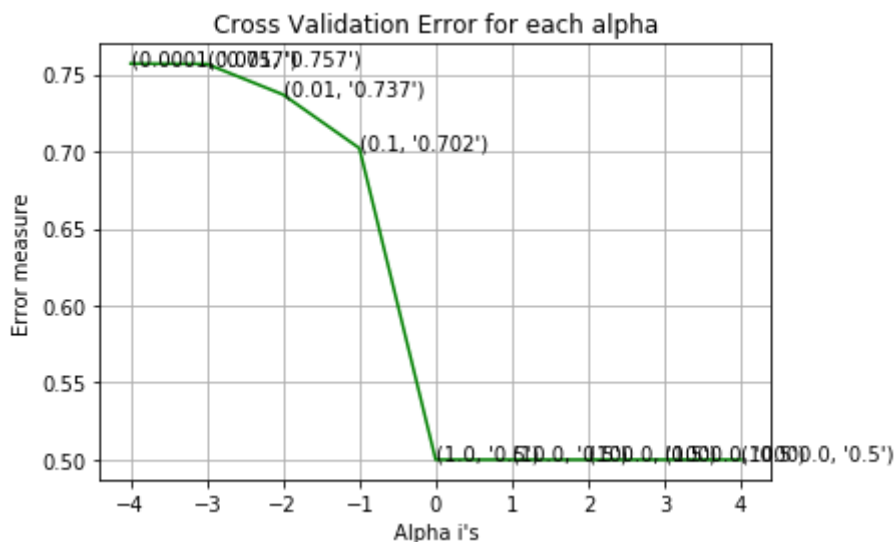
For alpha 1.0, cross validation AUC score 0.5

For alpha 10.0, cross validation AUC score 0.5

For alpha 100.0, cross validation AUC score 0.5

For alpha 1000.0, cross validation AUC score 0.5

For alpha 10000.0, cross validation AUC score 0.5



The Optimal C value is: 0.0001

In [182]:

```

best_alpha = alpha[np.argmax(cv_auc_score)]
logreg = SGDClassifier(alpha = best_alpha, class_weight = 'balanced', penalty = 'l1', loss=
logreg.fit(X_train_final[selected_features], y_train)
logreg_sig_clf = CalibratedClassifierCV(logreg, method='sigmoid')
logreg_sig_clf.fit(X_train_final[selected_features], y_train)
y_pred_prob = logreg_sig_clf.predict_proba(X_train_final[selected_features])[:,1]
print('For best alpha {0}, The Train AUC score is {1}'.format(best_alpha, roc_auc_score(y_t
y_pred_prob = logreg_sig_clf.predict_proba(X_val_final[selected_features])[:,1]
print('For best alpha {0}, The Cross validated AUC score is {1}'.format(best_alpha, roc_auc
y_pred_prob = logreg_sig_clf.predict_proba(X_test_final[selected_features])[:,1]
print('For best alpha {0}, The Test AUC score is {1}'.format(best_alpha, roc_auc_score(y_te
y_pred = logreg.predict(X_test_final[selected_features])
print('The test AUC score is :', roc_auc_score(y_test,y_pred_prob))
print('The percentage of misclassified points {:05.2f}% :'.format((1-accuracy_score(y_test,
plot_confusion_matrix(y_test, y_pred)

```

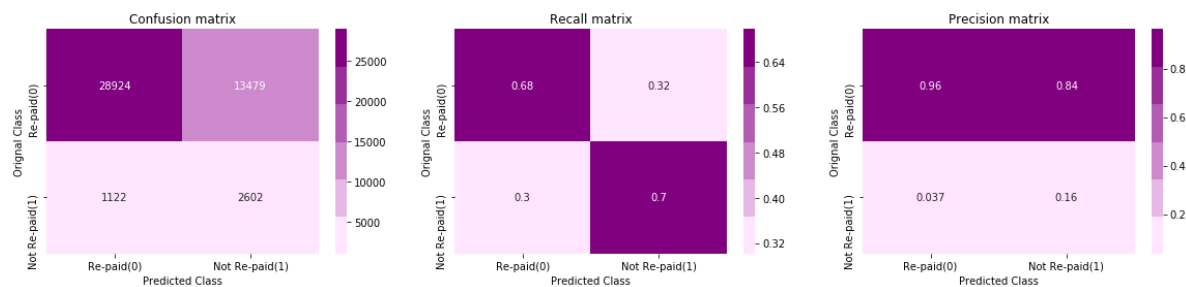
For best alpha 0.0001, The Train AUC score is 0.761944014332322

For best alpha 0.0001, The Cross validated AUC score is 0.7570320665909553

For best alpha 0.0001, The Test AUC score is 0.7596387045553111

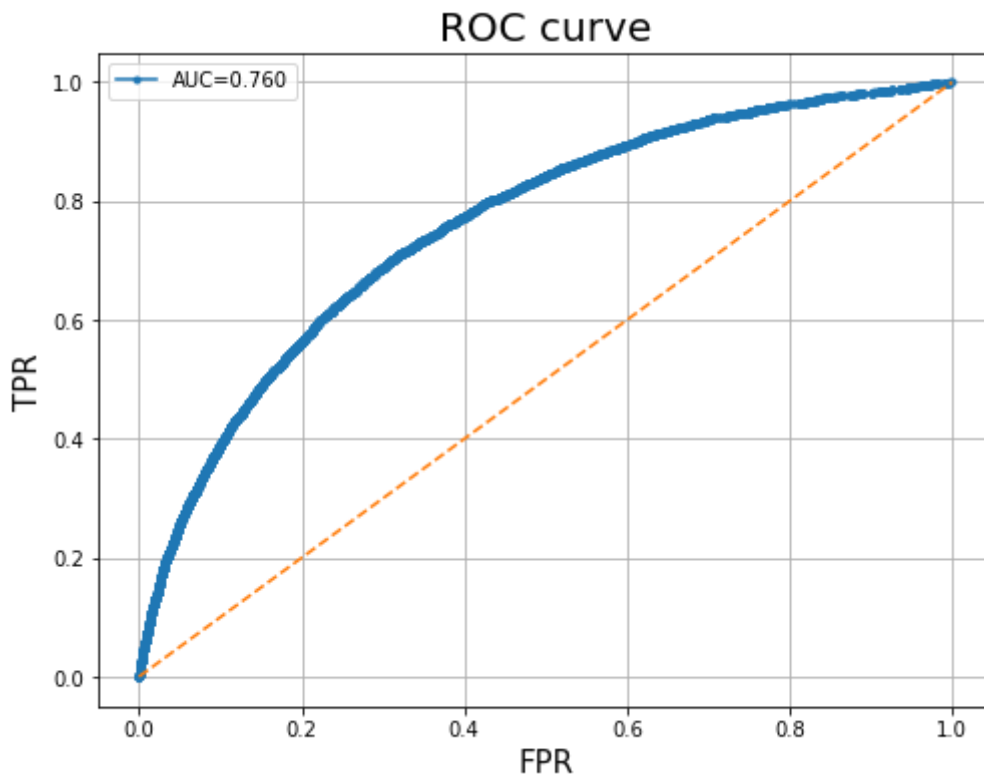
The test AUC score is : 0.7596387045553111

The percentage of misclassified points 31.65% :



In [183]:

```
from sklearn.metrics import roc_curve
fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob)
auc = roc_auc_score(y_test, y_pred_prob)
plt.figure(figsize=(8,6))
plt.plot(fpr, tpr, marker='.')
plt.plot([0, 1], [0, 1], linestyle='--')
plt.title('ROC curve', fontsize = 20)
plt.xlabel('FPR', fontsize=15)
plt.ylabel('TPR', fontsize=15)
plt.grid()
plt.legend(["AUC=%.3f"%auc])
plt.show()
```



Random Forest with selected features:

In [184]:

```
alpha = [200,500,1000,2000]
max_depth = [7, 10]
cv_auc_score = []
for i in alpha:
    for j in max_depth:
        clf = RandomForestClassifier(n_estimators=i, criterion='gini', max_depth=j, class_weight='balanced',
                                    random_state=42, n_jobs=-1)
        clf.fit(X_train_final[selected_features], y_train)
        sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
        sig_clf.fit(X_train_final[selected_features], y_train)
        y_pred_prob = sig_clf.predict_proba(X_val_final[selected_features])[:,1]
        cv_auc_score.append(roc_auc_score(y_val,y_pred_prob))
        print('For n_estimators {0}, max_depth {1} cross validation AUC score {2}'.
              format(i,j,roc_auc_score(y_val,y_pred_prob)))
```

```
For n_estimators 200, max_depth 7 cross validation AUC score 0.7452472874654
488
For n_estimators 200, max_depth 10 cross validation AUC score 0.749025247311
7833
For n_estimators 500, max_depth 7 cross validation AUC score 0.7455486703423
924
For n_estimators 500, max_depth 10 cross validation AUC score 0.749255937472
5554
For n_estimators 1000, max_depth 7 cross validation AUC score 0.745485025999
6956
For n_estimators 1000, max_depth 10 cross validation AUC score 0.74940844324
97516
For n_estimators 2000, max_depth 7 cross validation AUC score 0.745232038154
2832
For n_estimators 2000, max_depth 10 cross validation AUC score 0.74921271631
44553
```

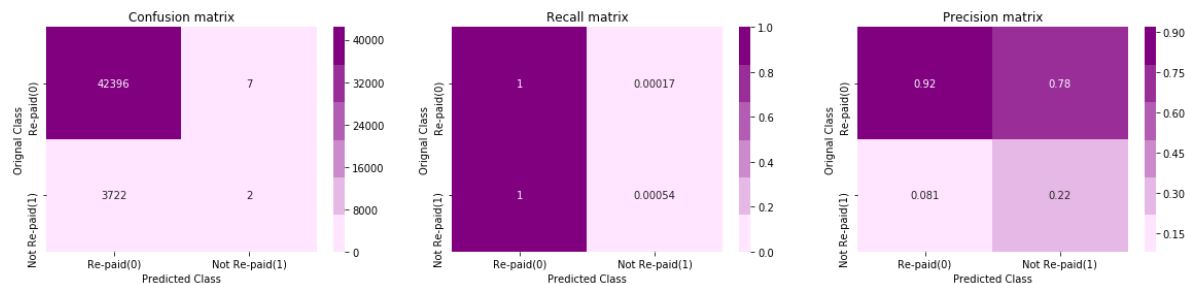
In [185]:

```

best_alpha = np.argmax(cv_auc_score)
print('The optimal values are: n_estimators {0}, max_depth {1} '.format(alpha[int(best_alpha/2)],
                                                                           max_depth[int(best_alpha/2)]))
rf = RandomForestClassifier(n_estimators=alpha[int(best_alpha/2)], criterion='gini', max_depth=max_depth[int(best_alpha/2)],
                           class_weight='balanced', random_state=42, n_jobs=-1)
rf.fit(X_train_final[selected_features], y_train)
rf_sig_clf = CalibratedClassifierCV(rf, method="sigmoid")
rf_sig_clf.fit(X_train_final[selected_features], y_train)
y_pred_prob = rf_sig_clf.predict_proba(X_train_final[selected_features])[:,1]
print('For best n_estimators {0} best max_depth {1}, The Train AUC score is {2}'.format(alpha[int(best_alpha/2)],
                                                                           max_depth[int(best_alpha/2)],roc_auc_score(y_train,y_pred_prob)))
y_pred_prob = rf_sig_clf.predict_proba(X_val_final[selected_features])[:,1]
print('For best n_estimators {0} best max_depth {1}, The Validation AUC score is {2}'.format(alpha[int(best_alpha/2)],
                                                                           max_depth[int(best_alpha/2)],roc_auc_score(y_val,y_pred_prob)))
y_pred_prob = rf_sig_clf.predict_proba(X_test_final[selected_features])[:,1]
print('For best n_estimators {0} best max_depth {1}, The Test AUC score is {2}'.format(alpha[int(best_alpha/2)],
                                                                           max_depth[int(best_alpha/2)],roc_auc_score(y_test,y_pred_prob)))
y_pred = rf_sig_clf.predict(X_test_final[selected_features])
print('The test AUC score is :', roc_auc_score(y_test,y_pred_prob))
print('The percentage of misclassified points {:05.2f}% :'.format((1-accuracy_score(y_test,y_pred))))
plot_confusion_matrix(y_test, y_pred)

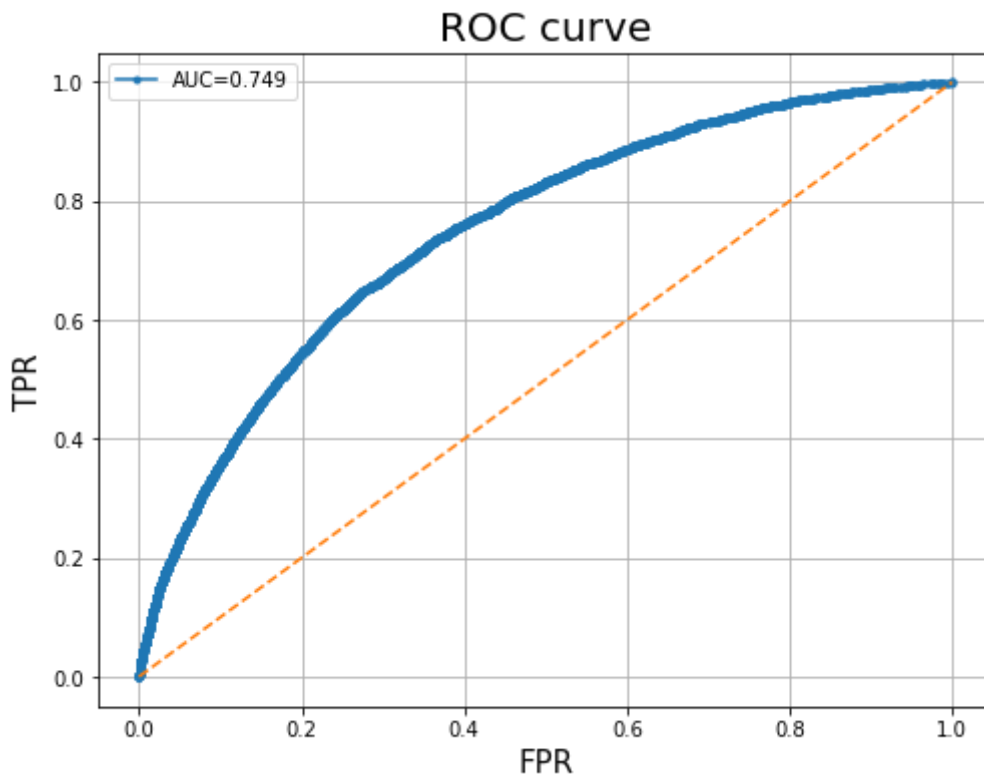
```

The optimal values are: n_estimators 1000, max_depth 10
 For best n_estimators 1000 best max_depth 10, The Train AUC score is 0.8411687775440093
 For best n_estimators 1000 best max_depth 10, The Validation AUC score is 0.7494084432497516
 For best n_estimators 1000 best max_depth 10, The Test AUC score is 0.7491606356105411
 The test AUC score is : 0.7491606356105411
 The percentage of misclassified points 08.08% :



In [186]:

```
from sklearn.metrics import roc_curve
fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob)
auc = roc_auc_score(y_test, y_pred_prob)
plt.figure(figsize=(8,6))
plt.plot(fpr, tpr, marker='.')
plt.plot([0, 1], [0, 1], linestyle='--')
plt.title('ROC curve', fontsize = 20)
plt.xlabel('FPR', fontsize=15)
plt.ylabel('TPR', fontsize=15)
plt.grid()
plt.legend(["AUC=%.3f"%auc])
plt.show()
```



LightGBM with selected features:

In [187]:

```

weight = np.ones((len(X_train_final),), dtype=int)
for i in range(len(X_train_final)):
    if y_train[i] == 0:
        weight[i] = 1
    else:
        weight[i] = 11

train_data = lgb.Dataset(X_train_final[selected_features], label = y_train, weight = weight)
valid_data = lgb.Dataset(X_val_final[selected_features], label = y_val)
cv_auc_score = []
max_depth = [3, 5, 7, 10]
for i in max_depth:

    params = {'boosting_type': 'gbdt',
              'max_depth': i,
              'objective': 'binary',
              'nthread': 5,
              'num_leaves': 32,
              'learning_rate': 0.05,
              'max_bin': 512,
              'subsample_for_bin': 200,
              'subsample': 0.7,
              'subsample_freq': 1,
              'colsample_bytree': 0.8,
              'reg_alpha': 20,
              'reg_lambda': 20,
              'min_split_gain': 0.5,
              'min_child_weight': 1,
              'min_child_samples': 10,
              'scale_pos_weight': 1,
              'num_class': 1,
              'metric': 'auc'
             }

    lgbm = lgb.train(params,
                     train_data,
                     2500,
                     valid_sets=valid_data,
                     early_stopping_rounds=100,
                     verbose_eval=10
                    )

    y_pred_prob = lgbm.predict(X_val_final[selected_features])
    cv_auc_score.append(roc_auc_score(y_val, y_pred_prob))
    print('For max_depth {0} and some other parameters, cross validation AUC score {1}'.format(i, cv_auc_score[-1]))
print('The optimal max_depth: ', max_depth[np.argmax(cv_auc_score)])
params = {'boosting_type': 'gbdt',
          'max_depth': max_depth[np.argmax(cv_auc_score)],
          'objective': 'binary',
          'nthread': 5,
          'num_leaves': 32,
          'learning_rate': 0.05,
          'max_bin': 512,
          'subsample_for_bin': 200,
          'subsample': 0.7,
          'subsample_freq': 1,
          'colsample_bytree': 0.8,
          'reg_alpha': 20,
          'reg_lambda': 20,
          'min_split_gain': 0.5,

```



```

        'min_child_weight': 1,
        'min_child_samples': 10,
        'scale_pos_weight': 1,
        'num_class' : 1,
        'metric' : 'auc'
    }
lgbm = lgb.train(params,
                 train_data,
                 2500,
                 valid_sets=valid_data,
                 early_stopping_rounds= 100,
                 verbose_eval= 10
                )
y_pred_prob = lgbm.predict(X_train_final[selected_features])
print('For best max_depth {0}, The Train AUC score is {1}'.format(max_depth[np.argmax(cv_auc)],
                                                                    roc_auc_score(y_train,y_pred_prob)))
y_pred_prob = lgbm.predict(X_val_final[selected_features])
print('For best max_depth {0}, The Cross validated AUC score is {1}'.format(max_depth[np.argmax(cv_auc)],
                                                                    roc_auc_score(y_val,y_pred_prob)))
y_pred_prob = lgbm.predict(X_test_final[selected_features])
print('For best max_depth {0}, The Test AUC score is {1}'.format(max_depth[np.argmax(cv_auc)],
                                                                    roc_auc_score(y_test,y_pred_prob)))
y_pred = np.ones((len(X_test_final),), dtype=int)
for i in range(len(y_pred_prob)):
    if y_pred_prob[i]<=0.5:
        y_pred[i]=0
    else:
        y_pred[i]=1
print('The test AUC score is :', roc_auc_score(y_test,y_pred_prob))
print('The percentage of misclassified points {:05.2f}% :'.format((1-accuracy_score(y_test,y_pred))))
plot_confusion_matrix(y_test, y_pred)

```

Training until validation scores don't improve for 100 rounds

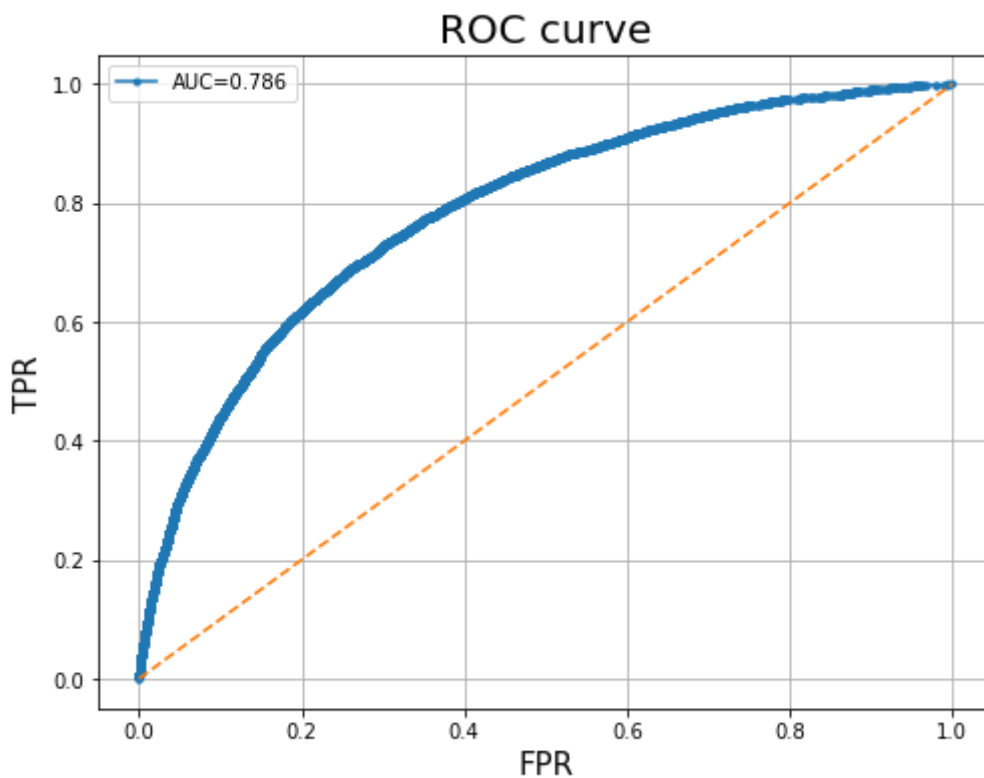
```

[10]  valid_0's auc: 0.713549
[20]  valid_0's auc: 0.722387
[30]  valid_0's auc: 0.725972
[40]  valid_0's auc: 0.731676
[50]  valid_0's auc: 0.737617
[60]  valid_0's auc: 0.741776
[70]  valid_0's auc: 0.745808
[80]  valid_0's auc: 0.748856
[90]  valid_0's auc: 0.75141
[100] valid_0's auc: 0.753801
[110] valid_0's auc: 0.755493
[120] valid_0's auc: 0.757107
[130] valid_0's auc: 0.758353
[140] valid_0's auc: 0.759626
[150] valid_0's auc: 0.760839
[160] valid_0's auc: 0.762015
[170] valid_0's auc: 0.763146
[180] valid_0's auc: 0.76396
[190] valid_0's auc: 0.764671
[200] valid_0's auc: 0.765271
[210] valid_0's auc: 0.765771
[220] valid_0's auc: 0.766171
[230] valid_0's auc: 0.766571
[240] valid_0's auc: 0.766971
[250] valid_0's auc: 0.767371
[260] valid_0's auc: 0.767771
[270] valid_0's auc: 0.768171
[280] valid_0's auc: 0.768571
[290] valid_0's auc: 0.768971
[300] valid_0's auc: 0.769371
[310] valid_0's auc: 0.769771
[320] valid_0's auc: 0.770171
[330] valid_0's auc: 0.770571
[340] valid_0's auc: 0.770971
[350] valid_0's auc: 0.771371
[360] valid_0's auc: 0.771771
[370] valid_0's auc: 0.772171
[380] valid_0's auc: 0.772571
[390] valid_0's auc: 0.772971
[400] valid_0's auc: 0.773371
[410] valid_0's auc: 0.773771
[420] valid_0's auc: 0.774171
[430] valid_0's auc: 0.774571
[440] valid_0's auc: 0.774971
[450] valid_0's auc: 0.775371
[460] valid_0's auc: 0.775771
[470] valid_0's auc: 0.776171
[480] valid_0's auc: 0.776571
[490] valid_0's auc: 0.776971
[500] valid_0's auc: 0.777371
[510] valid_0's auc: 0.777771
[520] valid_0's auc: 0.778171
[530] valid_0's auc: 0.778571
[540] valid_0's auc: 0.778971
[550] valid_0's auc: 0.779371
[560] valid_0's auc: 0.779771
[570] valid_0's auc: 0.780171
[580] valid_0's auc: 0.780571
[590] valid_0's auc: 0.780971
[600] valid_0's auc: 0.781371
[610] valid_0's auc: 0.781771
[620] valid_0's auc: 0.782171
[630] valid_0's auc: 0.782571
[640] valid_0's auc: 0.782971
[650] valid_0's auc: 0.783371
[660] valid_0's auc: 0.783771
[670] valid_0's auc: 0.784171
[680] valid_0's auc: 0.784571
[690] valid_0's auc: 0.784971
[700] valid_0's auc: 0.785371
[710] valid_0's auc: 0.785771
[720] valid_0's auc: 0.786171
[730] valid_0's auc: 0.786571
[740] valid_0's auc: 0.786971
[750] valid_0's auc: 0.787371
[760] valid_0's auc: 0.787771
[770] valid_0's auc: 0.788171
[780] valid_0's auc: 0.788571
[790] valid_0's auc: 0.788971
[800] valid_0's auc: 0.789371
[810] valid_0's auc: 0.789771
[820] valid_0's auc: 0.790171
[830] valid_0's auc: 0.790571
[840] valid_0's auc: 0.790971
[850] valid_0's auc: 0.791371
[860] valid_0's auc: 0.791771
[870] valid_0's auc: 0.792171
[880] valid_0's auc: 0.792571
[890] valid_0's auc: 0.792971
[900] valid_0's auc: 0.793371
[910] valid_0's auc: 0.793771
[920] valid_0's auc: 0.794171
[930] valid_0's auc: 0.794571
[940] valid_0's auc: 0.794971
[950] valid_0's auc: 0.795371
[960] valid_0's auc: 0.795771
[970] valid_0's auc: 0.796171
[980] valid_0's auc: 0.796571
[990] valid_0's auc: 0.796971

```

In [188]:

```
from sklearn.metrics import roc_curve
fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob)
auc = roc_auc_score(y_test, y_pred_prob)
plt.figure(figsize=(8,6))
plt.plot(fpr, tpr, marker='.')
plt.plot([0, 1], [0, 1], linestyle='--')
plt.title('ROC curve', fontsize = 20)
plt.xlabel('FPR', fontsize=15)
plt.ylabel('TPR', fontsize=15)
plt.grid()
plt.legend(["AUC=%.3f"%auc])
plt.show()
```



In []: