FINAL PROJECT PRESENTATION

TAI PHAM DA17-2023

Credit Card Approval Prediction

DATA: https://www.kaggle.com/datasets/rikdifos/credit-card-approval-prediction

NAME_INCOME_TYPE

	. >	1 ~	1 * 4	
MA	t a	div	القال	۰
1 4 10	ιa	чu	liệu	۰

application_record.c			credit_re		
Feature name	Explanation	Remarks			
ID	Client number		Feature name	Explan ation	Remarks
CODE_GENDER	Gender			Client	
FLAG_OWN_CAR	Is there a car		ID	numbe r	
FLAG_OWN_REALTY	Is there a property				
()= ()=()=()=()			MONTHS_ BALANCE	Record	The month of the extra month, and so on
CNT_CHILDREN	Number of children		DALANOL	month	month, and 30 on
AMT_INCOME_TOTAL	Annual income				0: 1-29 days past due

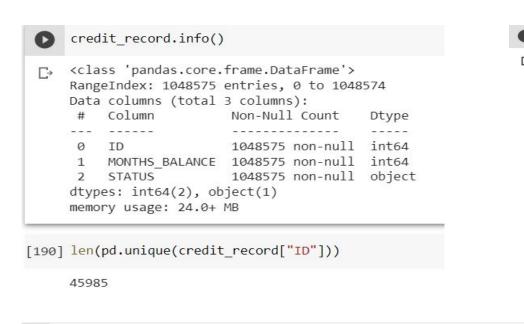
Load data and check data

```
from google.colab import drive drive.mount('/content/drive')

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).
```

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, confusion_matrix
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn import svm
from sklearn.svm import SVC
from sklearn.metrics import classification_report
```

```
[128] application_record =pd.read_csv("/content/drive/MyDrive/archive/application_record.csv") credit_record = pd.read_csv("/content/drive/MyDrive/archive/credit_record.csv")
```



new["approval"].value counts()

Name: approval, dtype: int64

45318 667



45985 rows × 1 columns

EDA

-Kiểm tra dữ liệu bị thiếu

```
new_data.isna().sum()
    CODE GENDER
    FLAG_OWN_CAR
    FLAG_OWN_REALTY
    CNT_CHILDREN
    AMT INCOME TOTAL
    NAME_INCOME_TYPE
    NAME_EDUCATION_TYPE
    NAME_FAMILY_STATUS
    NAME_HOUSING_TYPE
    DAYS_BIRTH
    DAYS EMPLOYED
    FLAG_MOBIL
    FLAG WORK PHONE
    FLAG_PHONE
    FLAG EMAIL
    OCCUPATION TYPE
                          11323
    CNT FAM MEMBERS
                              0
    approval
                              0
    dtype: int64
```

[143] new_data=pd.merge(application_record,new,how="inner",on="ID")

0

new_data

₽		ID	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_INCOME_TOTAL
	0	5008804	М	Υ	Υ	0	427500.0
	1	5008805	М	Y	Υ	0	427500.0
	2	5008806	М	Υ	Υ	0	112500.0
	3	5008808	F	N	Υ	0	270000.0
	4	4 5008809	F	N	Υ	0	270000.0
	•••	1966	***	***	1996	1865	(aux)
	36452	5149828	М	Υ	Y	0	315000.0
	36453	5149834	F	N	Υ	0	157500.0
	36454	5149838	F	N	Υ	0	157500.0
	36455	5150049	F	N	Υ	0	283500.0
	36456	5150337	М	N	Y	0	112500.0

36457 rows × 19 columns

Pre Processing Data

```
[26] from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
```

							Correl	ation He	atmap						
CODE_GENDER -	1	0.36							0.17						
FLAG_OWN_CAR -		1				-0.1			-0.16						
FLAG_OWN_REALTY		-0.015	1	0.033	-0.047		-0.18	-0.13		-0.21	0.067				
AMT_INCOME_TOTAL -			0.033	1	-0.073	0.23			0.17						
NAME_INCOME_TYPE -			-0.047	-0.073	1	0.057			-0.36				-0.12		
NAME_EDUCATION_TYPE -				-0.23	0.057	1	-0.036	-0.17			-0.045	-0.098		-0.041	
NAME_HOUSING_TYPE -			-0.18			-0.036	1		0.11						
DAYS_BIRTH -			-0.13			-0.17	0.21	1	-0.62	0.18			-0.18		
DAYS_EMPLOYED -	-0.17	-0 16		-0.17	-0.36	0.12	-0.11	-0.62	1	-0.24				-0.22	
FLAG_WORK_PHONE -			-0.21			-0.021		0.18	0.24	1	0.31				
FLAG_PHONE -						0.045				0.31	1	0.01		-0.0042	
FLAG_EMAIL -											0.01	1	-0.038		
OCCUPATION_TYPE -		-0.094			-0.12			-0.18				-0.038	1	0.071	
CNT_FAM_MEMBERS						0.041			-0.22		-0.0042		-0.071	1	-0.0057
approval -														-0.0057	1
	CODE_GENDER -	FLAG_OWN_CAR -	FLAG_OWN_REALTY -	AMT_INCOME_TOTAL_	NAME_INCOME_TYPE -	EDUCATION_TYPE -	1E_HOUSING_TYPE -	DAYS_BIRTH -	DAYS_EMPLOYED -	FLAG_WORK_PHONE -	FLAG_PHONE -	FLAG_EMAIL -	OCCUPATION_TYPE -	CNT_FAM_MEMBERS -	approval -
	8	FLAG	FLAG_OV	AMT_INCO	NAME_INC	NAME_EDUCA	NAME_HOU	۵	DAYS	FLAG_WO	2	. ≖	OCCUPA	CNT_FAM	

-Drop columns, Scale dữ liệu

	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	AMT_INCOME_TOTAL	NAME_INCOME_TYPE	NAME_EDUCATION_TYPE	NAME_HOUSING_TYPE	DAYS_BIRTH	DAYS_EMPLOYED	FLAG_WORK_PHONE	FLAG_PHONE
0	1.0	1.0	1.0	0.258721	1.00	0.25	0.8	0.744324	0.029324	1.0	0.0
1	1.0	1.0	1.0	0.258721	1.00	0.25	0.8	0.744324	0.029324	1.0	0.0
2	1.0	1.0	1.0	0.055233	1.00	1.00	0.2	0.208232	0.038270	0.0	0.0
3	0.0	0.0	1.0	0.156977	0.00	1.00	0.2	0.342071	0.033237	0.0	1.0
4	0.0	0.0	1.0	0.156977	0.00	1.00	0.2	0.342071	0.033237	0.0	1.0
	***										•••
36452	1.0	1.0	1.0	0.186047	1.00	1.00	0.2	0.441828	0.034894	0.0	0.0
36453	0.0	0.0	1.0	0.084302	0.00	0.25	0.2	0.722697	0.037768	0.0	1.0
36454	0.0	0.0	1.0	0.084302	0.25	0.25	0.2	0.722697	0.037768	0.0	1.0
36455	0.0	0.0	1.0	0.165698	1.00	1.00	0.2	0.407292	0.039527	0.0	0.0
36456	1.0	0.0	1.0	0.055233	1.00	1.00	0.8	0.903810	0.038115	0.0	0.0

36457 rows × 14 columns

Building Model

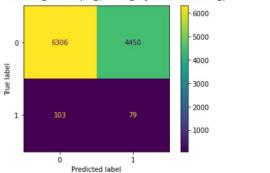
- Logistic Regression
- SVM
- Decision Tree
- RandomForest

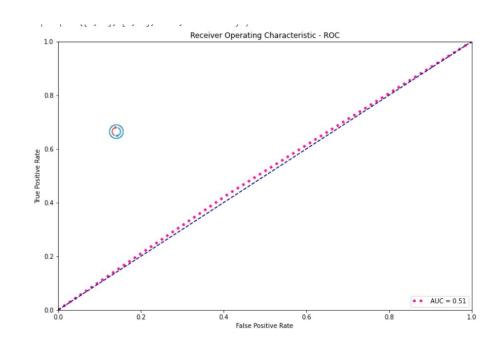
**Upsampling bằng SMOTE

• Logistic Regression

	precision	recall	f1-score	support
0	0.98	0.59	0.73	10756
1	0.02	0.43	0.03	182
accuracy			0.58	10938
macro avg	0.50	0.51	0.38	10938
weighted avg	0.97	0.58	0.72	10938

/usr/local/lib/python3.8/dist-packages/sklearn/utils/depreca warnings.warn(msg, category=FutureWarning)



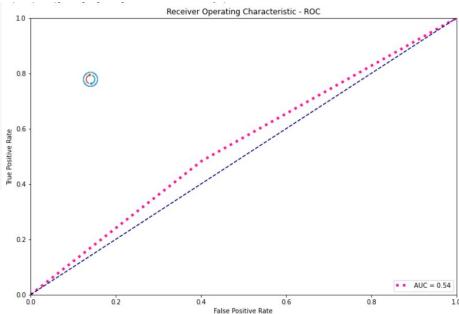


• SVM

Accuracy Score is 0.596

0 1

0 6431 4325
1 04 88



Decision Tree

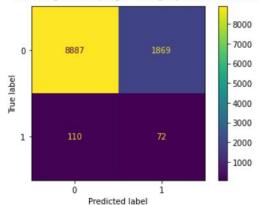
```
model = DecisionTreeClassifier(max depth=12,
                               min samples split=8,
                               random state=1024)
model.fit(X train resample, Y train resample)
y predict = model.predict(X test)
print(classification_report( Y_test, y_predict))
plot confusion matrix(model, X test, Y test)
print('Accuracy Score is {:.5}'.format(accuracy score(Y test, y predict)))
print(pd.DataFrame(confusion matrix(Y test,y predict)))
```

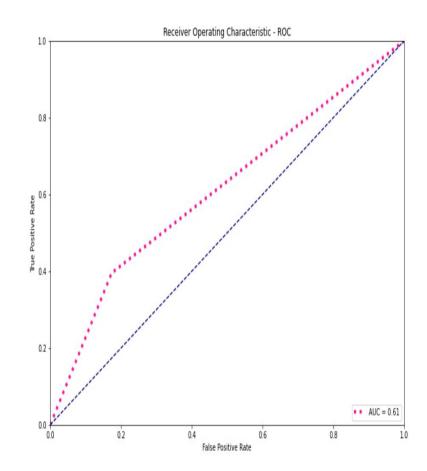
		precision	recall	f1-score	support
	0	0.99	0.83	0.90	10756
	1	0.04	0.40	0.07	182
accur	racy			0.82	10938
macro	avg	0.51	0.61	0.48	10938
weighted	avg	0.97	0.82	0.89	10938

Accuracy Score is 0.81907

0 8887 1869 1 110 72

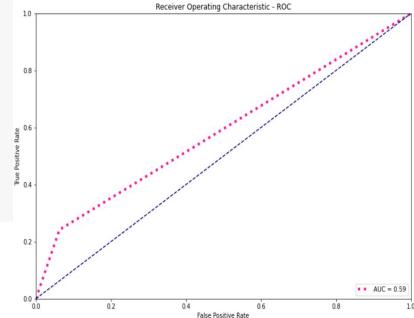
/usr/local/lib/python3.8/dist-packages/sklearn/utils/deprecation.py:87: warnings.warn(msg, category=FutureWarning)





• Random Forest

```
model = RandomForestClassifier(n estimators=250,
                             max depth=12,
                             min samples leaf=16
model.fit(X train resample, Y train resample)
y predict 2 = model.predict(X test)
print('Accuracy Score is {:.5}'.format(accuracy_score(Y_test, y_predict_2)))
print(pd.DataFrame(confusion_matrix(Y_test,y_predict_2)))
Accuracy Score is 0.92558
```



THANK YOU FOR YOUR ATTENTION!!