CREDIT RISK PREDICTION

ID/X PARTNERS DATA SCIENTIST PROJECT BASED INTERNSHIP PROGRAM

BY:

Saskia Dwi Ulfah





PRESENTATION CONTENT

#1

Data Understanding

#2

Data Preparation

#3

Data Modeling

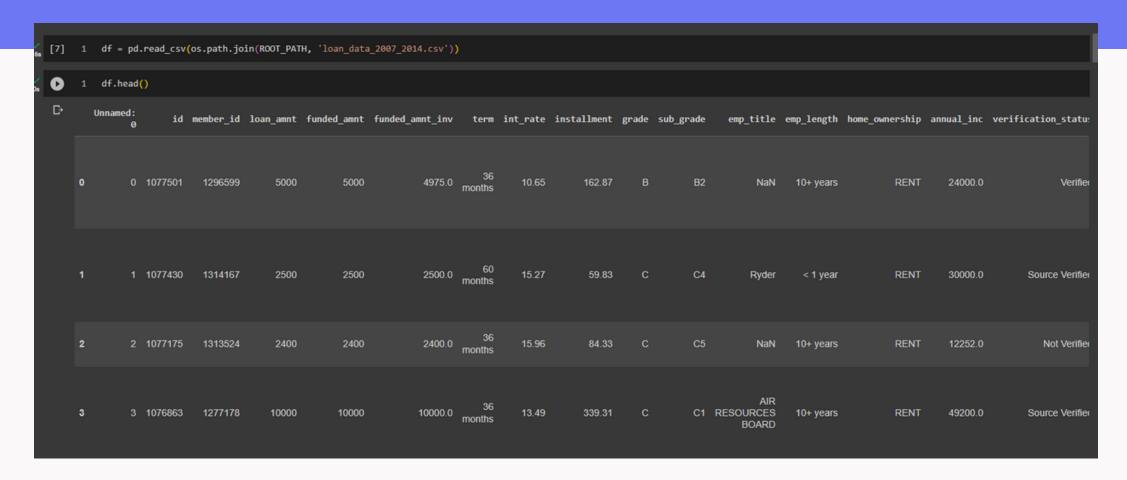
#4

Model Evaluation and Saving Final Model

DATA UNDERSTANDING



DATA UNDERSTANDING



In this task, we were given data. The data itself contained **75 columns** and **466285 entries**. The features related to the lendee's information and their credit activities, such as employment, home ownership, loan amount, purpose of credit activity, payment plan, etc. Based on those features, we were asked to **build a model for credit risk prediction**.

```
df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 466285 entries, 0 to 466284
Data columns (total 75 columns):
                                 Non-Null Count
    Column
                                                 Dtype
    Unnamed: 0
                                 466285 non-null int64
     id
                                 466285 non-null int64
     member id
                                 466285 non-null int64
     loan amnt
                                 466285 non-null int64
     funded amnt
                                 466285 non-null int64
     funded amnt inv
                                 466285 non-null float64
    term
                                 466285 non-null object
     int rate
                                 466285 non-null float64
     installment
                                 466285 non-null float64
     grade
                                 466285 non-null object
    sub grade
                                 466285 non-null object
     emp title
                                 438697 non-null
                                                 object
    emp length
                                 445277 non-null object
    home ownership
                                 466285 non-null object
    annual inc
                                 466281 non-null float64
    verification status
                                 466285 non-null object
    issue d
                                 466285 non-null
                                                 object
    loan status
                                 466285 non-null object
    pymnt plan
                                 466285 non-null object
    url
                                 466285 non-null object
    desc
                                 125983 non-null
                                                 object
                                 466285 non-null object
    purpose
    title
                                 466265 non-null
                                                 obiect
 23 zip code
                                 466285 non-null object
                                 466285 non-null
    addr state
                                                 object
 25 dti
                                 466285 non-null float64
    deling 2yrs
                                 466256 non-null float64
    earliest cr line
                                 466256 non-null
                                                 obiect
    ing last 6mths
                                 466256 non-null float64
    mths since last deling
                                 215934 non-null float64
    mths since last record
                                 62638 non-null
                                                 float64
                                 466256 non-null float64
    open acc
 32 pub_rec
                                 466256 non-null float64
```

DATA PREPARATION



CHECK DUPLICATED ROWS

```
Check for duplicated records.

[10] 1 df[df.duplicated()]

Unnamed: id member_id loan_amnt funded_amnt_inv term int_rate installment grade sub_grade emp_title

No duplicated rows.
```

This data **contained no duplicated rows**. Hence, every record in this data corresponded to one credit activity.

CHECK AND DROP NULL VALUES

In this phase, we checked columns whose percentage of null values was 100% or less. We **dropped columns whose percentage of null values was > 40%**.

FILL-IN NULL VALUES

```
cols_with_null = df.columns[df.isnull().any()]
null_num_cols = df[cols_with_null].select_dtypes(exclude='object').columns
null_non_num_cols = df[cols_with_null].select_dtypes(include='object').columns
```

```
for column in null_num_cols:
    median = df[column].median()
    df[column].fillna(median, inplace=True)
```

```
for column in null_non_num_cols:
    mode_value = df[column].mode()[0]
    df[column].fillna(mode_value, inplace=True)
```

For the column whose percentage of null values was <40%:

- For **numerical data**, we filled in with **median value**.
- For **non-numerical data**, we filled in with **mode value**.

DROP UNNECESSARY COLUMNS

```
unncss_cols = ['Unnamed: 0', 'id', 'member_id', 'url', 'policy_code', 'application_type']

# drop

df = df.drop(columns=unncss_cols)
```

```
df = df.drop(columns=['emp_title','title'])
```

```
df = df.drop(columns=['issue_d','last_credit_pull_d'])
```

We also dropped columns that belong to one of these categories:

- Columns with high cardinality.
- Columns with one variation of value.
- Free text column.
- Columns that require value from the future.

ENCODE NON-NUMERICAL COLUMNS

```
# Term column
df["term"] = df["term"].str.replace(" 36 months", "36")
df["term"] = df["term"].str.replace(" 60 months", "60")
df["term"] = df["term"].astype(int)
```

The computer worked with numbers. So, we needed to transform text data into the form of numbers. To do this, we utilized LabelEncoder() from sklearn.

FILTERING

```
1 df = df[(df['loan_status']=='Fully Paid') | (df['loan_status']=='Charged Off')]
 1 df.shape
(227214, 30)
 1 df['loan_status'].value_counts()
Fully Paid 184739
Charged Off 42475
Name: loan_status, dtype: int64
 1 # Map the risk
 3 risk_mapping = {
        'Fully Paid' : 'Low Risk',
        'Charged Off': 'High Risk'
 1 df['loan_status'] = df['loan_status'].map(risk_mapping)
 1 df['loan_status'].value_counts()
Low Risk
           184739
High Risk 42475
Name: loan_status, dtype: int64
 1 # Map the risk to numerical value
    num_risk_mapping = {
        'Low Risk': 0,
         'High Risk': 1
 1 df['loan_status'] = df['loan_status'].map(num_risk_mapping)
 1 df['loan_status'].value_counts()
    42475
Name: loan_status, dtype: int64
```

In this task, we only considered "Fully Paid" and "Charged Off" categories. Those two categories, we turned into its number representation:

- Fully Paid \rightarrow Low Risk \rightarrow 0
- Charged Off → High Risk → 1

ASSIGN LABEL, FEATURES, AND SMOTE

```
1 X = df.drop(columns=['loan_status'])
      2 y = df['loan_status']
[69] 1 scaler= StandardScaler()
      2 X= scaler.fit transform(X)
[70]   1    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
[71] 1 y_train.value_counts()
         147757
     Name: loan_status, dtype: int64
         print("Before oversampling: ", Counter(y_train))
         smote = SMOTE()
          X_train,y_train= smote.fit_resample(X_train,y_train)
          print("After oversampling: ",Counter(y_train))
     Before oversampling: Counter({0: 147757, 1: 34014})
     After oversampling: Counter({0: 147757, 1: 147757})
```

- Last step, we assigned features as X and loan_status as target (y).
- Next, we standardized X values with StandardScaler() from sklearn.
- Next, we split data into 2 sets: train set and test set.
- Because we have imbalanced classes, we used SMOTE from train set. SMOTE oversampled the minority class to have the same amount as majority class.

DATA MODELING



TRAINING

s = setup(X_train, target=y_train, session_id = 123, use_gpu=True)

best = compare_models()

	Model	Accuracy	AUC	Recall	Prec.	F1	Карра	MCC	TT (Sec)
et	Extra Trees Classifier	0.9063	0.9649	0.9027	0.9093	0.9060	0.8127	0.8127	52.9780
rf Random Forest Classifier		0.8840	0.9460	0.8296	0.9308	0.8773	0.7679	0.7725	88.3580
lightgbm Light Gradient Boosting Mad		0.8829	0.9328	0.7845	0.9768	0.8701	0.7658	0.7811	6.4750
gbc	Gradient Boosting Classifier	0.8467	0.9161	0.7893	0.8916	0.8373	0.6933	0.6979	144.7530
dt	Decision Tree Classifier	0.8049	0.8049	0.8175	0.7974	0.8073	0.6098	0.6100	7.2290
ada	Ada Boost Classifier	0.7813	0.8654	0.7513	0.7993	0.7745	0.5625	0.5636	29.7240
knn	K Neighbors Classifier	0.7547	0.8713	0.9488	0.6835	0.7946	0.5094	0.5527	52.9510
Ir	Logistic Regression	0.6544	0.7113	0.6589	0.6530	0.6559	0.3087	0.3088	2.7510
ridge	Ridge Classifier	0.6522	0.0000	0.6510	0.6526	0.6518	0.3044	0.3044	0.3940
lda	Linear Discriminant Analysis	0.6522	0.7086	0.6510	0.6526	0.6518	0.3044	0.3044	1.6900
svm	SVM - Linear Kernel	0.6495	0.0000	0.6403	0.6528	0.6458	0.2990	0.2996	2.1490
nb	Naive Bayes	0.5218	0.6979	0.9657	0.5161	0.6690	0.0437	0.0858	0.7180
qda	Quadratic Discriminant Analysis	0.5207	0.7107	0.9788	0.5123	0.6717	0.0413	0.0909	1.2470
dummy	Dummy Classifier	0.5000	0.5000	0.9000	0.4500	0.6000	0.0000	0.0000	0.2380

- For the sake of simplicity, we used PyCarer which automatically trained and compared models.
- Based on the experiment, the Extra Trees
 Classifier reached the best
 performance amongst models.

```
ExtraTreesClassifier

ExtraTreesClassifier(bootstrap=False, ccp_alpha=0.0, class_weight=None, criterion='gini', max_depth=None, max_features='sqrt', max_leaf_nodes=None, max_samples=None, min_impurity_decrease=0.0, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=-1, oob_score=False,
```

random state=123, verbose=0, warm start=False)

TRAINING

```
xt = ExtraTreesClassifier(bootstrap=False, ccp_alpha=0.0, class_weight=None,
criterion='gini', max_depth=None, max_features='sqrt',
max_leaf_nodes=None, max_samples=None,
min_impurity_decrease=0.0, min_samples_leaf=1,
min_samples_split=2, min_weight_fraction_leaf=0.0,
n_estimators=100, n_jobs=-1, oob_score=False,
random_state=123, verbose=0, warm_start=False)
```

1 xt.fit(X_train, y_train)

```
ExtraTreesClassifier
```

```
ExtraTreesClassifier(bootstrap=False, ccp_alpha=0.0, class_weight=None, criterion='gini', max_depth=None, max_features='sqrt', max_leaf_nodes=None, max_samples=None, min_impurity_decrease=0.0, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=-1, oob_score=False, random_state=123, verbose=0, warm_start=False)
```

We trained Extra Tree Classifier by passing X_train and y_train as parameter.

MODEL EVALUATION AND SAVING FINAL MODEL



MODEL EVALUATION AND SAVING THE MODEL

0	<pre>print(classification_report(y_test, xt.predict(X_test)))</pre>									
₽		precision	recall	f1-score	support					
	0 1	0.83 0.42	0.94 0.18	0.89 0.25	36982 8461					
	accuracy macro avg weighted avg	0.63 0.76	0.56 0.80	0.80 0.57 0.77	45443 45443 45443					
[]	[] dump(xt, "/content/drive/MyDrive/Data Science/Rakamin/IDX Partners Virtual Internship/extra_tree.joblib") ['/content/drive/MyDrive/Data Science/Rakamin/IDX Partners Virtual Internship/extra tree.joblib']									
	[/concent/ul	1vc/Hybi 1ve/	Data Stie	nee/ NakaiiiI	II, IDA FAI CI	ici 3 vii cuui internamip/extra_cree.jooiio j				

We **gained an accuracy of 0.80 in the test set**. It seemed like the model outperformed in synthetical training data and underperformed in real testing data. For future improvement, we could experiment with various sampling strategies and tune the model.

THANK YOU