



by michael sitanggang

# CREDIT SCORING ID/X - RAKAMIN

FINAL PROJECT





# BUSINESS UNDERSTANDING

A credit score is based on credit history: number of open accounts, total levels of debt, repayment history, and other factors. Lenders use credit scores to evaluate the probability that an individual will repay loans in a timely manner



# PROBLEM SOLUTION

- minimize the possibility of default
- determine potential target markets
- determine the attributes that influence the credit score



loan\_amnt  
funded\_amnt  
funded\_amnt\_inv  
term  
int\_rate  
installment  
grade  
sub\_grade  
emp\_title  
emp\_length  
home\_ownership  
annual\_inc  
verification\_status  
issue\_d  
loan\_status  
pymnt\_plan  
url  
desc  
purpose  
title  
zip\_code  
addr\_state  
dti  
delinq\_2yrs  
earliest\_cr\_line  
inq\_last\_6mths  
mths\_since\_last\_delinq  
mths\_since\_last\_record  
open\_acc  
pub\_rec  
revol\_bal  
revol\_util  
total\_acc  
initial\_list\_status  
out\_prncp  
out\_prncp\_inv

## DATA UNDERSTANDING

- 466285 records, 72 Columns
- Imbalanced Target Labels
- dtypes: float64(46), int64(4), object(22)

**OPPORTUNITY TO  
ACCEPT LOAN  
APPLICATION**

total\_pymnt  
total\_pymnt\_inv  
total\_rec\_prncp  
total\_rec\_int  
total\_rec\_late\_fee  
recoveries  
collection\_recovery\_fee  
last\_pymnt\_d  
last\_pymnt\_amnt  
next\_pymnt\_d  
last\_credit\_pull\_d  
collections\_12\_mths\_ex\_med  
mths\_since\_last\_major\_derog  
policy\_code  
application\_type  
annual\_inc\_joint  
dti\_joint  
verification\_status\_joint  
acc\_now\_delinq  
tot\_coll\_amt  
tot\_cur\_bal  
open\_acc\_6m  
open\_il\_6m  
open\_il\_12m  
open\_il\_24m  
mths\_since\_rcnt\_il  
total\_bal\_il  
il\_util  
open\_rv\_12m  
open\_rv\_24m  
max\_bal\_bc  
all\_util  
total\_rev\_hi\_lim  
inq\_fi  
total\_cu\_tl  
inq\_last\_12m



# DATA PREPARATION



Percentage of missing value in each features



>10%

Cleaning if  $\geq 50\%$  and imputation by median, mode



Convert data types



Diff data date-time

#Noted = difference  
calculated until July 2022

Delete Outlier



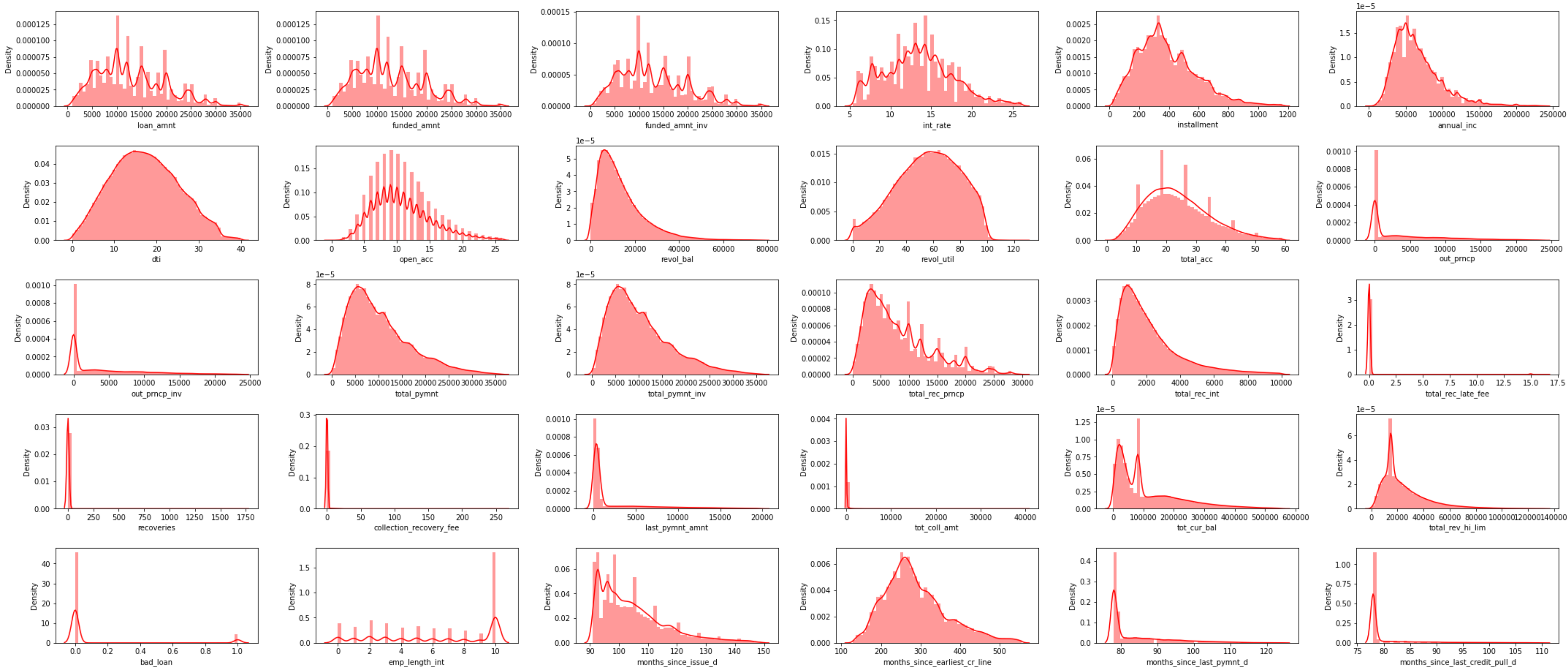
Label / Unique checking in each features

#Noted = found unreasonable categorical  
variables and single labels on several variables

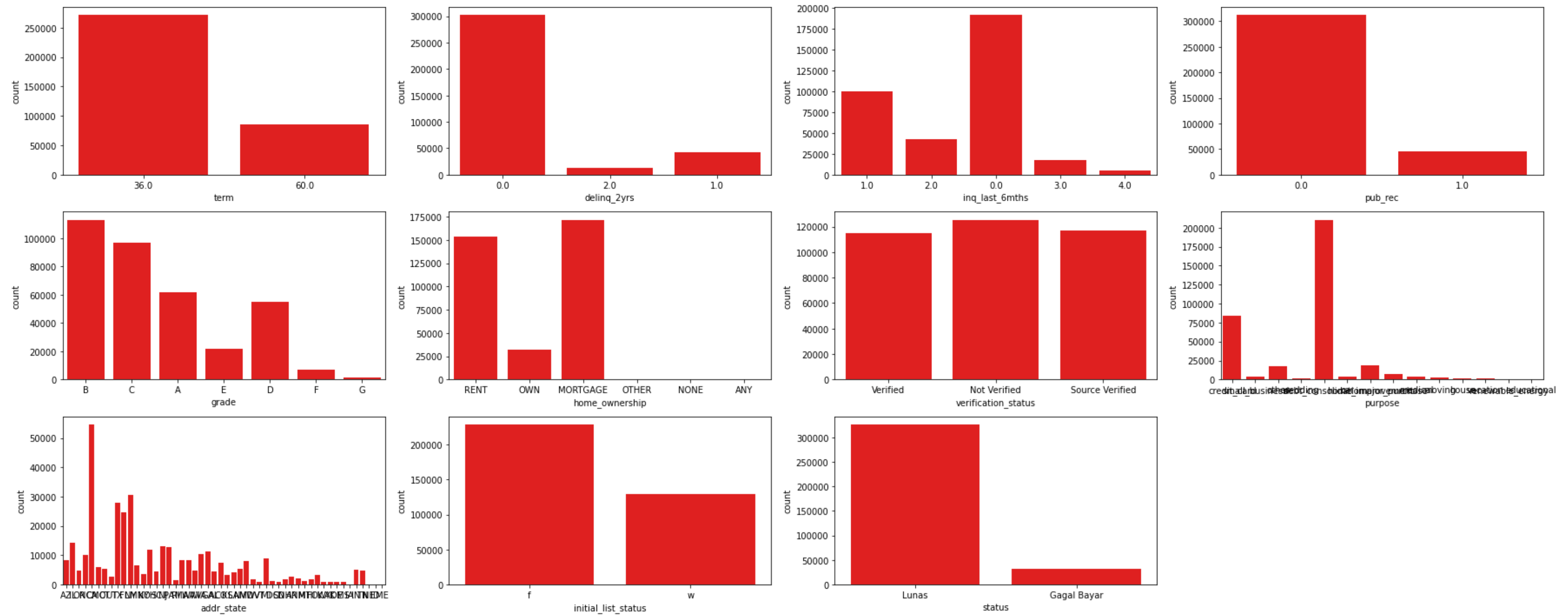
inq_last_12m	100.000000
total_bal_il	100.000000
dti_joint	100.000000
verification_status_joint	100.000000
annual_inc_joint	100.000000
open_acc_6m	100.000000
open_il_6m	100.000000
open_il_12m	100.000000
open_il_24m	100.000000
mths_since_rcnt_il	100.000000
il_util	100.000000
open_rv_24m	100.000000
total_cu_tl	100.000000
inq_fi	100.000000
max_bal_bc	100.000000
all_util	100.000000
open_rv_12m	100.000000
mths_since_last_record	86.566585
mths_since_last_major_derog	78.773926
desc	72.981546
mths_since_last_delinq	53.690554
next_pymnt_d	48.728567
tot_cur_bal	15.071469
tot_coll_amt	15.071469
total_rev_hi_lim	15.071469



# EXPLORATORY DATA ANALYSIS



- dist plot for all numeric variables and it can be seen that some variables are skew and multimodal
- for variables with high variance and high bias will be converted into categorical form |  $Q1=Q2=Q3=0$



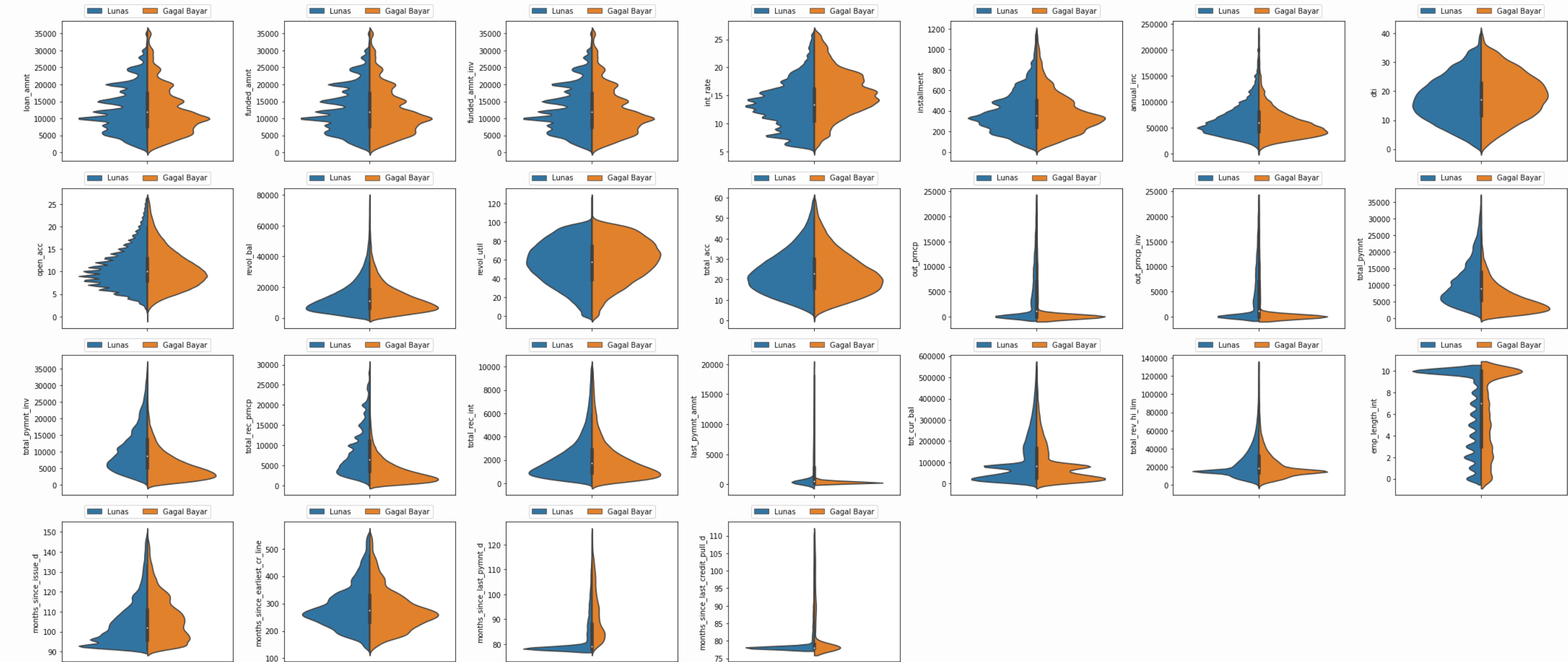
>>> It can be seen that the predictor variables with labels are not imbalanced and this will affect the accuracy of model, so a few labels are regrouped into several available labels or create other labels ;

```
data['home_ownership'].replace({'NONE':'RENT', 'ANY':'RENT', 'OTHER':'RENT'},inplace=True)

data['purpose'].replace({'educational':'major_purchase',
                        'house':'major_purchase',
                        'medical':'major_purchase',
                        'moving':'major_purchase',
                        'vacation':'other',
                        'wedding':'other',
                        'renewable_energy':'home_improvement'},inplace=True)

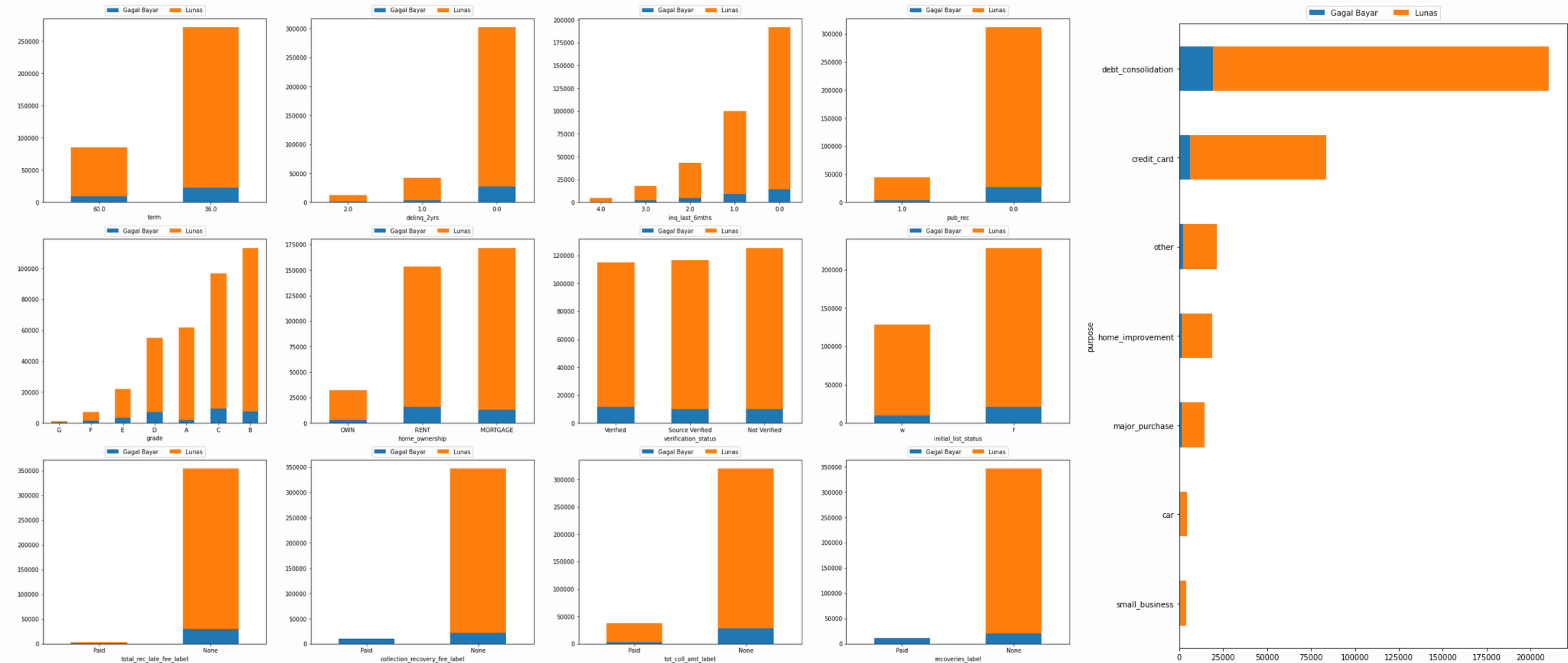
data['addr_state'].replace({'IA':'OTHER', 'ID':'OTHER', 'NE':'OTHER', 'ME':'OTHER'},inplace=True)
```

# Univariate Analysis - numeric variable



- dist plot for all numeric variables with paid and non-paid labels. The pattern of two labels so approaching
- there are several variables with high variance and high bias are maintained for log transformation

# Univariate Analysis - categorical variable



- count plot for all categorical variables. There are some labels of imbalance variables and will be maintained
- for 4 plots in bottom are the conversion results from numeric variable



# Bivariate Analysis - Multivariate Analysis

## ANOVA TEST

Find that for all numerical predictor variables of significance  $< 0.05$ , so  $H_0$  is rejected, means that there is an average difference between the tested target labels.

Conclusion = there is an effect of all numerical predictor variables on credit status.

## CHI SQUARED TEST

Find that almost all categorical predictor variables [except: pub\_rec] sign  $< 0.05$ , so  $H_0$  is rejected, means that there is an average difference between the tested target labels.

Conclusion = there is a categorical predictor variable has no effect on credit status and **the variable will be deleted.**

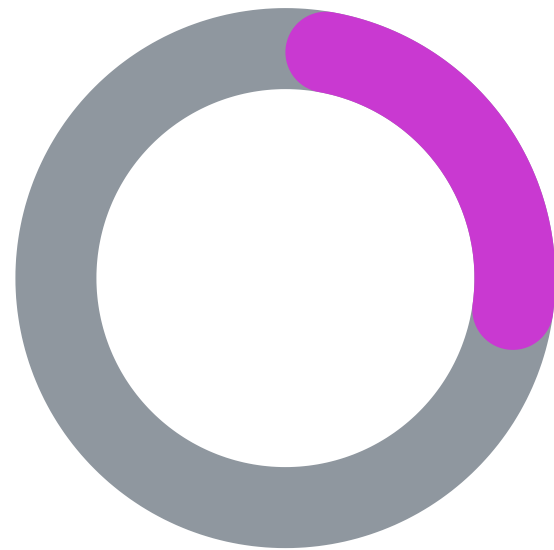
## PEARSON CORR

Find that for correlation all numerical predictor variables and there are several variables that have a correlation  $< 0.7$  and **the variables will be deleted**

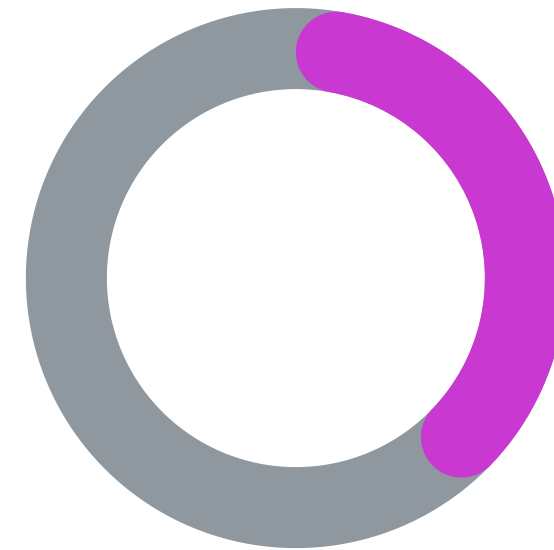




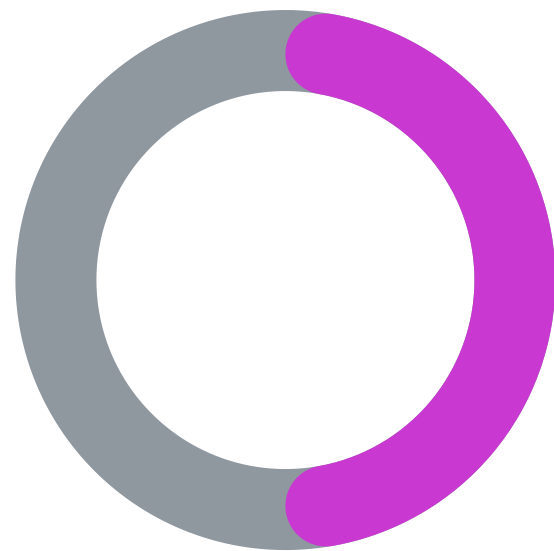
# DATA PREPROCESSING



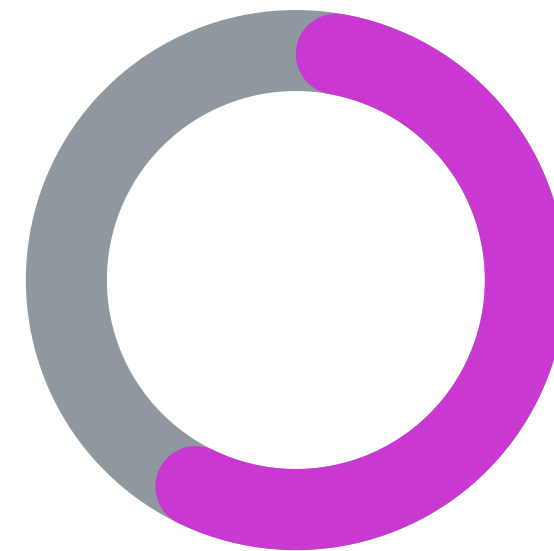
Label Encoding  
&  
One-Hot Encoding



Data Partitioning



Data Transformation



Undersampling for  
Imbalanced Target  
Labels



# Training Model ML - Compare Tree Algorithm

```
RFclassifier = RandomForestClassifier()  
RFclassifier.fit(X_train, y_train)  
y_pred_rf = RFclassifier.predict(X_train)  
y_pred_rf_test = RFclassifier.predict(X_test)
```

```
LRclassifier = LogisticRegression()  
LRclassifier.fit(X_train, y_train)  
y_pred_lr = LRclassifier.predict(X_train)  
y_pred_lr_test = LRclassifier.predict(X_test)
```

```
XTclassifier = ExtraTreesClassifier()  
XTclassifier.fit(X_train, y_train)  
y_pred_xt = XTclassifier.predict(X_train)  
y_pred_xt_test = XTclassifier.predict(X_test)
```

```
DTclassifier = DecisionTreeClassifier()  
DTclassifier.fit(X_train, y_train)  
y_pred_dt = DTclassifier.predict(X_train)  
y_pred_dt_test = DTclassifier.predict(X_test)
```

```
GBclassifier = GradientBoostingClassifier()  
GBclassifier.fit(X_train, y_train)  
y_pred_gb = GBclassifier.predict(X_train)  
y_pred_gb_test = GBclassifier.predict(X_test)
```

Model	Akurasi Train	Akurasi Test
RandomForest	1.000000	0.988540
Xtree	1.000000	0.976912
DecisionTree	1.000000	0.984462
GradientBoost	0.971550	0.975531
LogisticRegression	0.892719	0.936587

**then choose the important variable and underfit or overfit test**

Akurasi Train 1.0

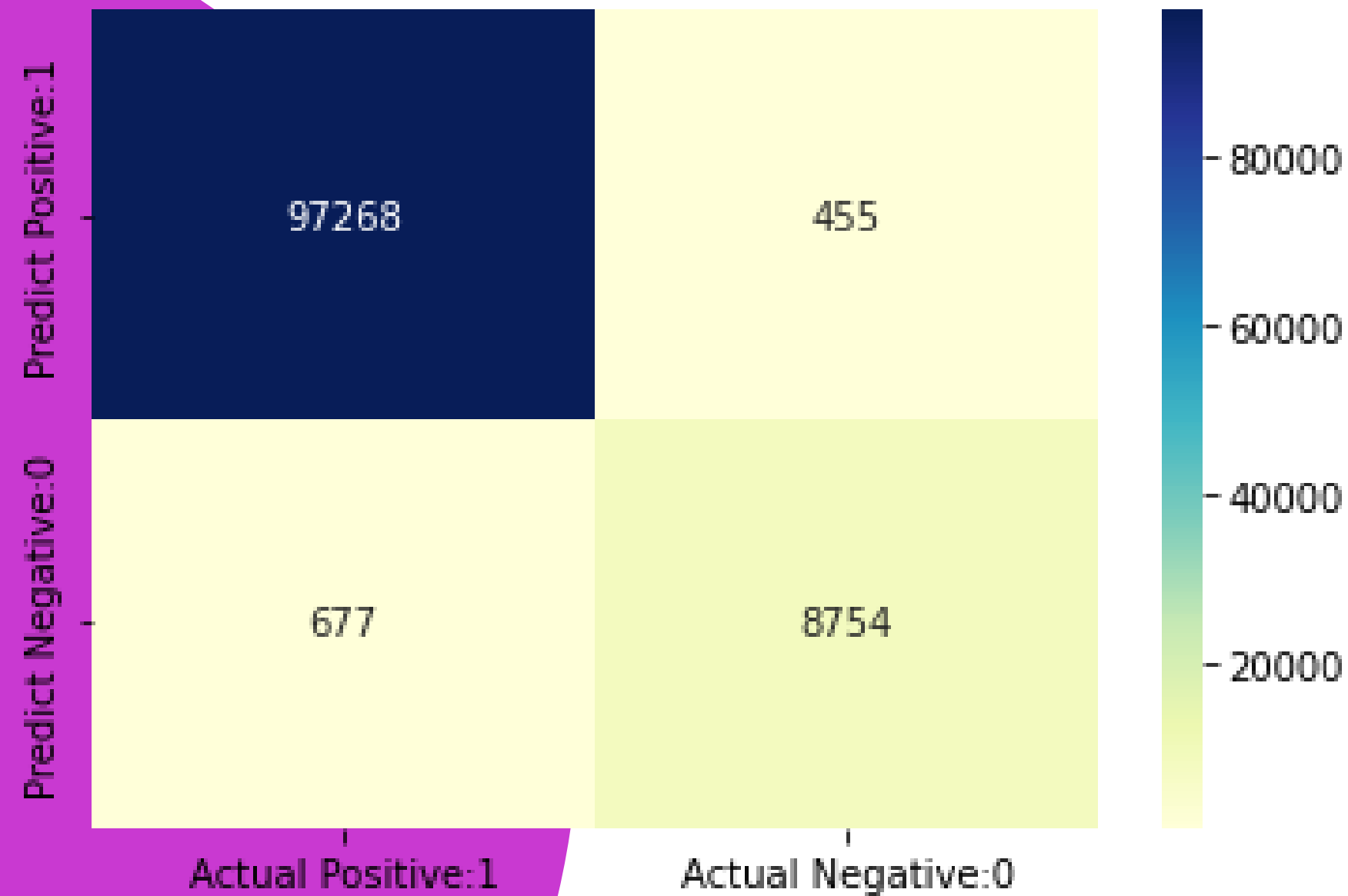
Akurasi Test 0.9894357653470706

**Model is not prone to underfit or overfit.**





# Model Evaluation

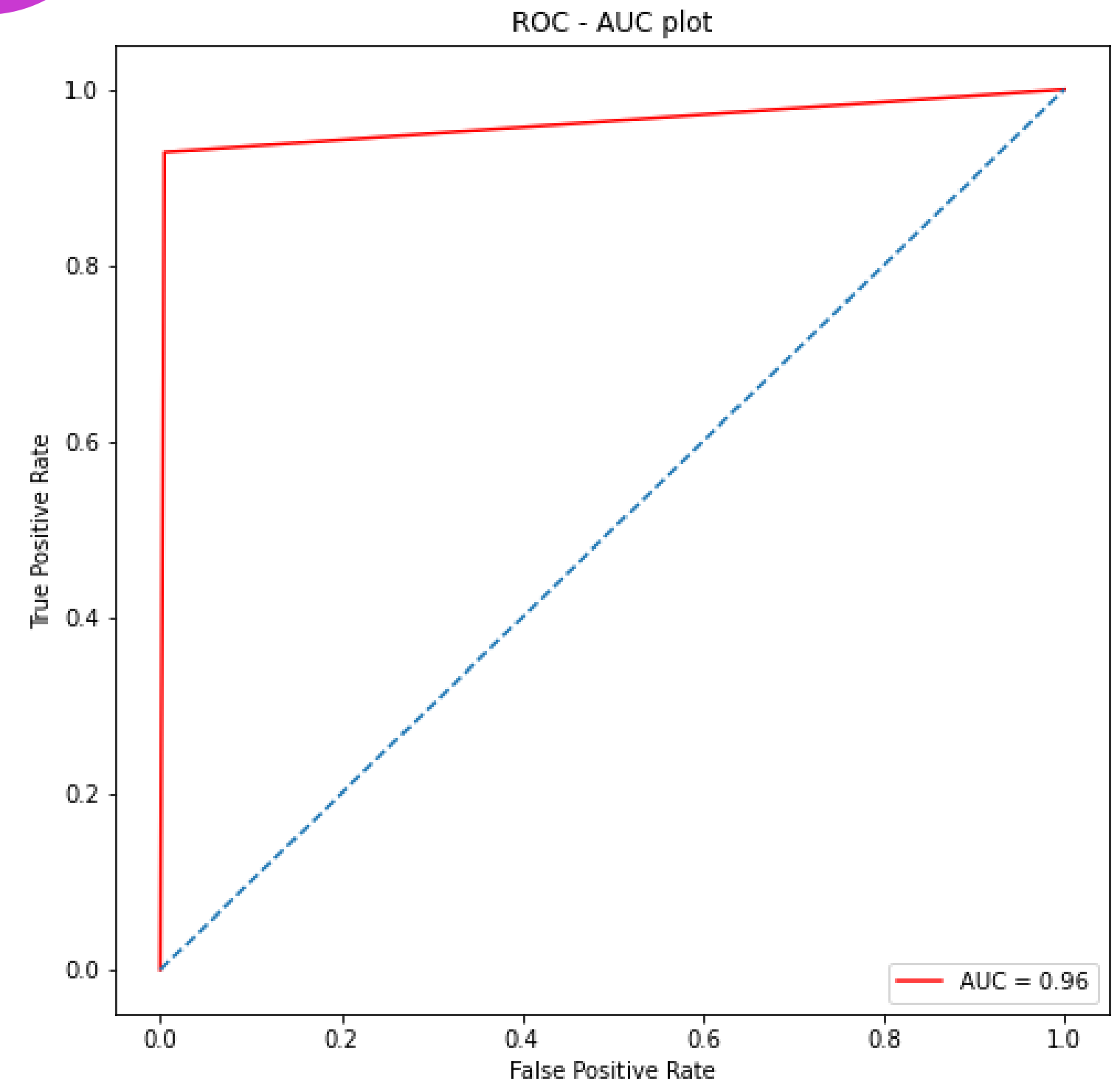


Akurasi Klasifikasi : 0.9894

Kesalahan Klasifikasi : 0.0106

Presisi : 0.9953

Sensitivitas : 0.9931



ROC- AUC more than 0.9 is considered outstanding classifier