# CREDIT RISK MODELLING WITH MACHINE LEARNING

Fidya Almira Suheri

# OBJECTIVE

create a model to determine or predict
credit risk based on the classification
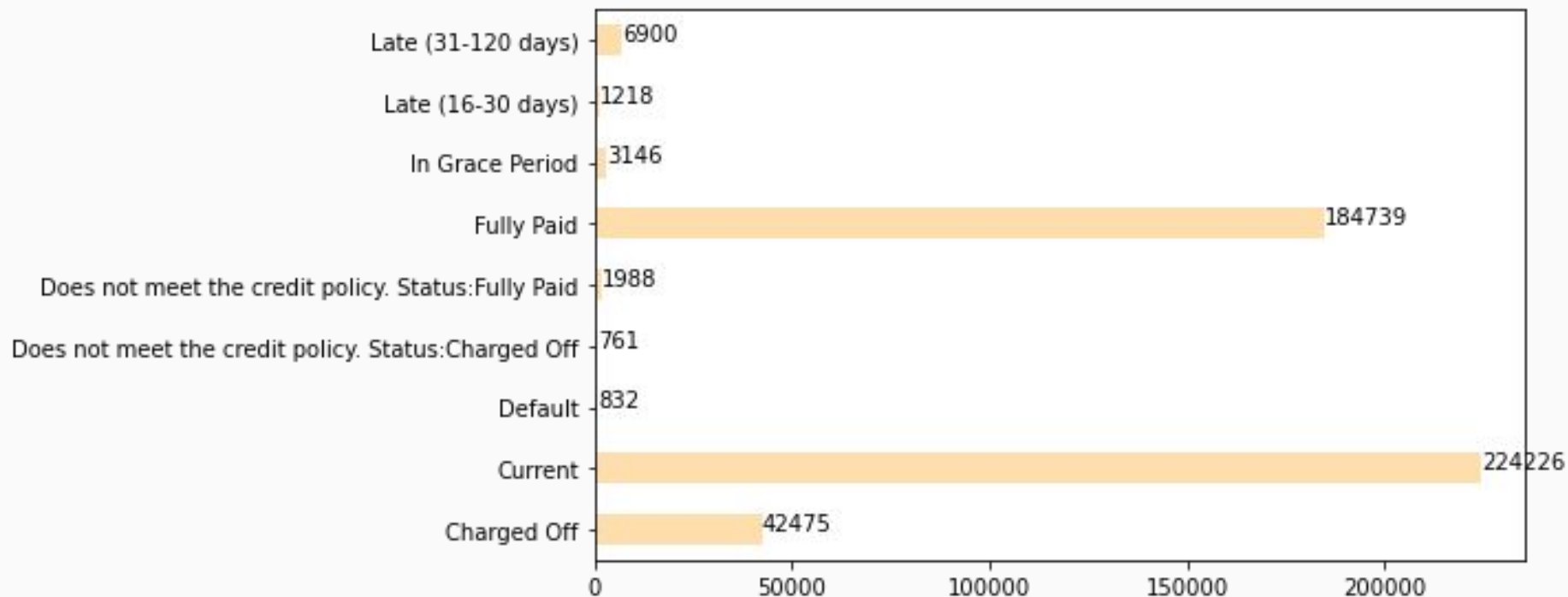of good and bad loans

01
TARGET

03
FEATURES

02
MODELLING

04
EVALUATION

Late (31-120 days) — 6900
Late (16-30 days) — 1218
In Grace Period — 3146
Fully Paid — 184739
Does not meet the credit policy. Status:Fully Paid — 1988
Does not meet the credit policy. Status:Charged Off — 761
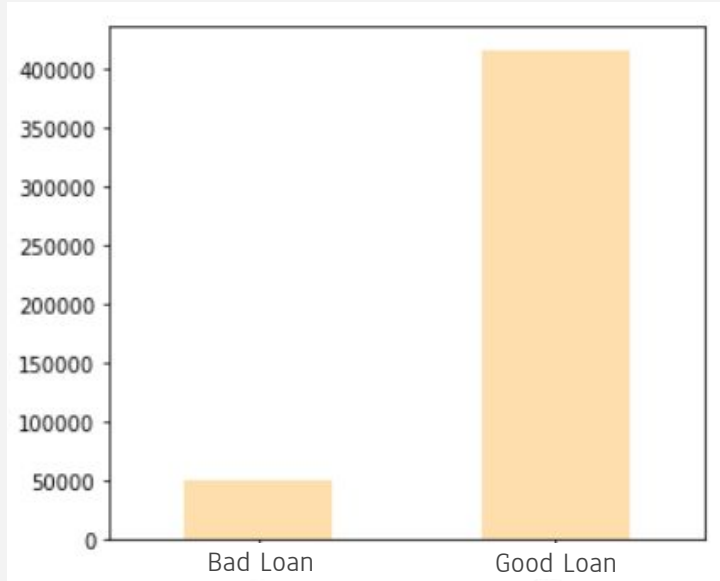Default — 832
Current — 224226
Charged Off — 42475

Status loan dibagi menjadi 2 tipe:

**Good Loan**
- Current
- Fully Paid
- In Grace Period
- Late (16-30) days
- Does not meet the credit policy. Status: Fully Paid

**Bad Loan**
- Charged Off
- Default'
- Late (31-120 days)
- Does not meet the credit policy. Status:Charged Off'

# PERSENTASE FEATURES DENGAN MISSING VALUE > 70%

| Features | Persentase | Features | Persentase |
|---|---|---|---|
| desc | 72.9815% | open_il_6m | 100% |
| mths_since_last_record | 86.5666% | open_il_12m | 100% |
| mths_since_last_major_derog | 78.7739% | open_il_24m | 100% |
| annual_inc_joint | 100% | mths_since_rcnt_il | 100% |
| dti_joint | 100% | total_bal_il | 100% |
| verification_status_joint | 100% | il_util | 100% |
| open_acc_6m | 100% | open_rv_12m | 100% |
| open_rv_24m | 100% | max_bal_bc | 100% |
| all_util | 100% | total_cu_tl | 100% |
| inq_fi | 100% | inq_last_12m | 100% |

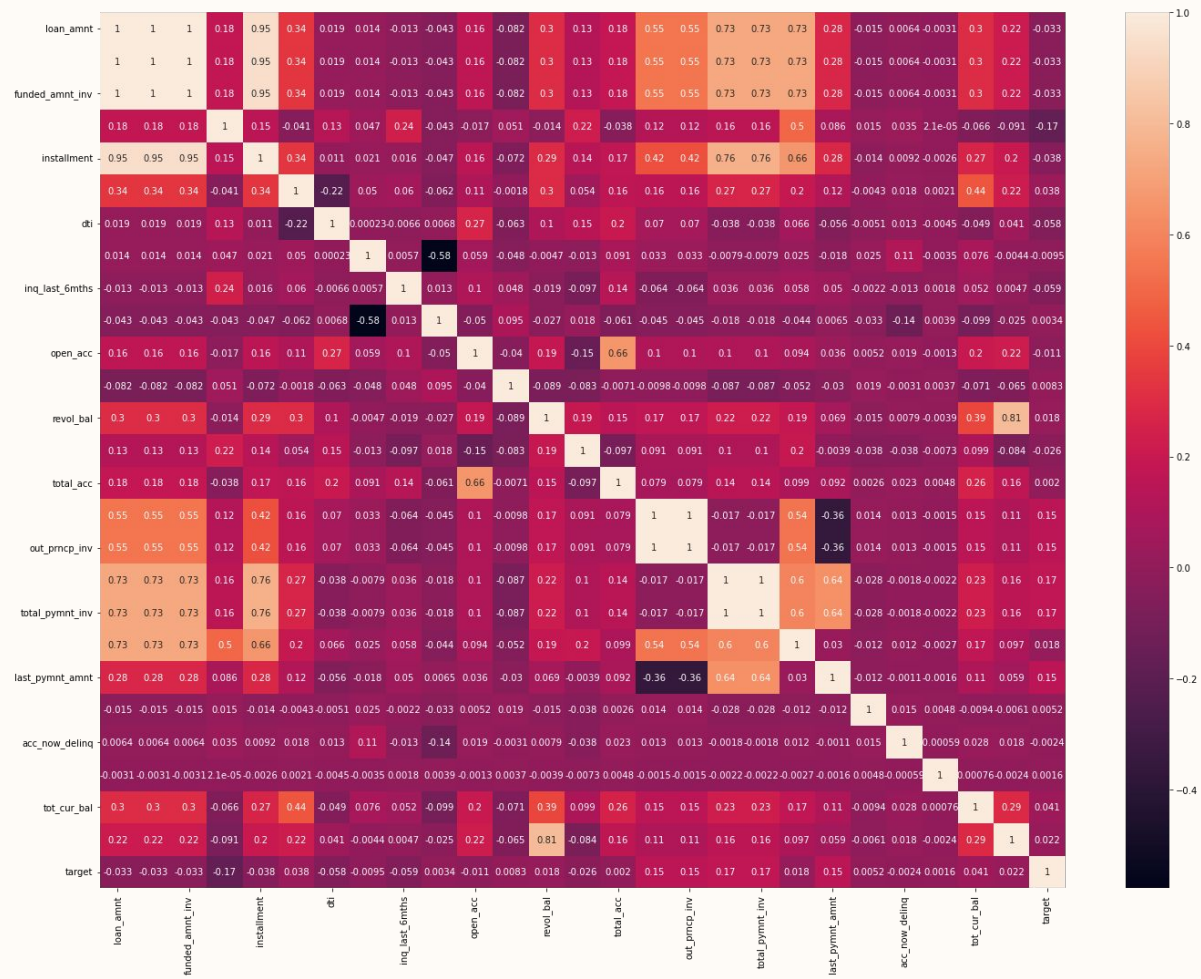| Features | Description |
|---|---|
| id | A unique LC assigned ID for the loan listing. |
| member_id | A unique LC assigned Id for the borrower member. |
| url | URL for the LC page with listing data. |
| zip_code | The first 3 numbers of the zip code provided by the borrower in the loan application. |
| policy_code | publicly available policy_code=1<br>new products not publicly available policy_code=2 |
| emp_title | The job title supplied by the Borrower when applying for the loan. |
| title | The loan title provided by the borrower |

*Features yang dapat diabaikan

| Features | Description |
|---|---|
| collection_recovery_fee | post charge off collection fee |
| next_pymnt_d | Next scheduled payment date |
| recoveries | post charge off gross recovery |
| total_rec_prncp | Principal received to date |
| total_rec_late_fee | Late fees received to date |
| sub_grade | LC assigned loan subgrade |

*Features yang dapat diabaikan

KORELASI ANTAR FEATURES

$$IV = \sum (\text{Event}\% - \text{Non Event}\%) * \ln\left(\frac{\text{Event}\%}{\text{Non Event}\%}\right)$$

$$IV = \sum (\text{Event}\% - \text{Non Event}\%) * (\text{WOE})$$

## Information Value>0.5

| Variable | IV |
|---|---|
| out_prncp | 0.703375 |
| total_pymnt | 0.515794 |
| last_pymnt_amnt | 1.491828 |
| mths_since_last_pymnt_d | 2.331187 |

## Information Value<0.02

| Variable | IV |
|---|---|
| emp_length | 0.007174 |
| home_ownership | 0.017952 |
| pymnt_plan | 0.000309 |
| addr_state | 0.010291 |
| delinq_2yrs | 0.001039 |
| mths_since_last_delinq | 0.002487 |
| open_acc | 0.004499 |
| pub_rec | 0.000504 |

| Variable | IV |
|---|---|
| revol_util | 0.008858 |
| initial_list_status | 0.011513 |
| total_rec_int | 0.011108 |
| collections_12_mths_ex_med | 0.000733 |
| application_type | 0.000000 |
| acc_now_delinq | 0.000200 |
| tot_coll_amt | 0.000738 |
| total_rev_hi_lim | 0.018835 |

**Kesimpulan**

- Terdapat 75 features pada data asli dengan dengan 40 feature yang memiliki *Missing Value.*
- Feature yang memiliki 70% *missing value* akan dihapus karena akan menimbulkan ketidakakuratan ketika pemodelan.
- Feature yang merupakan data kejadian yang terjadi dimasa depan akan dihapus karena tidak memiliki pengaruh dalam pemodelan
- Feature yang memiliki korelasi yang sangat kuat dengan feature lain akan dihapus karena akan menghasilkan model yang kurang maksimal pada *logistic regression*
- Feature yang bertipe kategorikal akan dikonversi menjadi numerik disesuaikan dengan data yang terkait.
- Feature yang memiliki *Information Value* jauh dibawah 0.2 dan diatas 0.5 akan dihapus karena tidak dapat digunakan untuk pemodelan.

**TRAIN TEST SPLIT** :
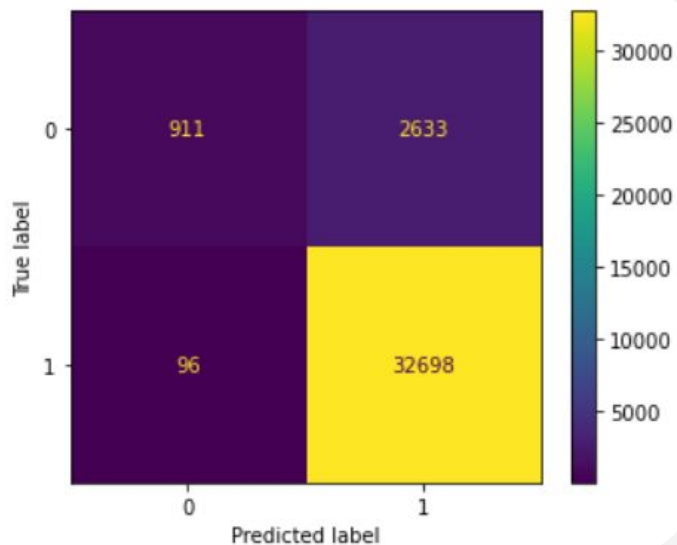
Data dibagi menjadi 2 bagian, yaitu:
- Training Dataset = 80%
- Test Dataset = 20%

**MODEL MACHINE LEARNING:**

Classification : Logistic Regression

## Confusion Matrix



## Classification Report

|  | Precision | Recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.90 | 0.26 | 0.40 | 3544 |
| 1 | 0.93 | 1 | 0.96 | 32794 |
| accuracy |  |  | 0.92 | 36338 |
| Macro avg | 0.92 | 0.63 | 0.68 | 36338 |
| Weigh avg | 0.92 | 0.92 | 0.91 | 36338 |

Training Accuracy : 0.9262390951371879
Testing Accuracy : 0.9248995541857009