

Virtual Internship Experience

Data Scientist

Presented by
Yoga Aria Sena



Yoga Aria Sena

About Me

I am an enthusiastic individual deeply intrigued by the world of data science. My journey is fueled by a profound interest in unraveling insights from data and transforming them into meaningful solutions. Aspiring to become a skilled data scientist, I am dedicated to expanding my expertise in programming, statistical analysis, and machine learning algorithms. My aspiration is to harness the power of data to create innovative solutions that make a tangible impact on the world.

My Experience

- Project Based Virtual Internship:
Data Scientist Kalbe Nutritionals x
Rakamin Academy
(August 2023 – September 2023)
- Project Based Virtual Internship:
Data Scientist ID/X Partners x Rakamin
Academy
(September 2023 – Present)

About Rakamin Academy



Rakamin Academy adalah Platform akselerasi karir, tempat bagi siapapun untuk belajar dan membangun karir di bidang digital dan teknologi.

Rakamin Academy menyediakan layanan Pendidikan yang paling terjangkau (efisien secara model bisnis), berdampak (dihubungkan langsung dengan ahli, praktik standar industri), dan bermakna (mengajarkan soft skills, membantu pengembangan karir, membangun komunitas).

About ID/X Partners



id/x partners didirikan pada tahun 2002 oleh mantan banker dan konsultan manajemen yang memiliki pengalaman luas dalam siklus kredit dan manajemen proses, pengembangan skor, dan manajemen kinerja. Pengalaman gabungan kami telah melayani perusahaan di berbagai wilayah Asia dan Australia, serta dalam berbagai industri, khususnya layanan keuangan, telekomunikasi, manufaktur, dan ritel.

id/x partners menyediakan layanan konsultasi yang mengkhususkan diri dalam penggunaan solusi analitik data dan pengambilan keputusan (DAD) yang dikombinasikan dengan disiplin manajemen risiko terintegrasi dan pemasaran untuk membantu klien mengoptimalkan profitabilitas portofolio dan proses bisnis.

Layanan konsultasi komprehensif dan solusi teknologi yang ditawarkan oleh id/x partners menjadikannya sebagai penyedia layanan satu atap.

Case Study

- **Latar Belakang**

Sebagai bagian dari tugas akhir kontrak sebagai intern Data Scientist di ID/X Partners, proyek ini melibatkan kerjasama dengan Lending Company.

- **Alat**

Python, Jupyter Notebook

- **Tujuan:**

Membangun model prediksi risiko kredit menggunakan dataset berisi data peminjaman yang diterima dan ditolak.

Menyajikan solusi secara visual kepada klien.

Menyediakan solusi teknologi yang dapat membantu perusahaan mengoptimalkan portofolio dan proses bisnis mereka.

Dataset dan Tujuan



Dataset dan Tujuan

| | id | member_id | loan_amnt | funded_amnt | funded_amnt_inv | term | int_rate |
|---|---------|-----------|-----------|-------------|-----------------|-----------|----------|
| 0 | 1077501 | 1296599 | 5000 | 5000 | 4975.0 | 36 months | 10.65 |
| 1 | 1077430 | 1314167 | 2500 | 2500 | 2500.0 | 60 months | 15.27 |
| 2 | 1077175 | 1313524 | 2400 | 2400 | 2400.0 | 36 months | 15.96 |
| 3 | 1076863 | 1277178 | 10000 | 10000 | 10000.0 | 36 months | 13.49 |
| 4 | 1075358 | 1311748 | 3000 | 3000 | 3000.0 | 60 months | 12.69 |

5 rows x 74 columns

- **Dataset**

Terdiri dari data status peminjaman yang mencakup berbagai status, seperti Charged Off, Current, Default, Does not meet the credit policy. Status:Charged Off, Does not meet the credit policy. Status:Fully Paid, Fully Paid, In Grace Period, Late (16-30 hari), dan Late (31-120 hari)

- **Tujuan**

Membangun peta risiko kredit yang akurat dan informatif.

Membantu dalam pengambilan keputusan terkait peminjaman dari setiap risiko kredit.

Data Cleaning



Data Cleaning

```
1 # Mengecek kolom-kolom duplikat
2
3 df.duplicated().sum()
```

```
1 # Menghapus kolom-kolom yang tidak diperlukan
2
3 unnecessary_columns = ['id', 'member_id', 'url', 'pymnt_plan', 'policy_code', 'application_type', 'acc_now_delinq']
4
5 df.drop(columns=unnecessary_columns, axis=1, inplace=True)
```

```
1 # Melakukan perubahan tipe data untuk kolom-kolom tanggal
2
3 df['term_months'] = df['term'].str.split(' ').str[1].astype(int)
4 df.drop('term', axis=1, inplace=True)
5
6 df['earliest_cr_line'] = pd.to_datetime(df['earliest_cr_line'], format='%b-%y', errors='coerce').dt.year
7
8 date_columns = ['issue_d', 'last_pymnt_d', 'last_credit_pull_d']
9 for column in date_columns:
10     df[column] = pd.to_datetime(df[column], format='%b-%y', errors='coerce').dt.month
```

Data Cleaning

```
1 # Menghapus kolom-kolom dengan missing value lebih dari 20%
2
3 columns_to_drop = []
4
5 # Loop melalui setiap kolom dalam DataFrame
6 for column in df.columns:
7     total_missing = df[column].isnull().sum()
8     percentage_missing = (total_missing / len(df)) * 100
9     if percentage_missing > threshold:
10         columns_to_drop.append(column)
11
12 # Drop kolom-kolom yang memiliki persentase missing value lebih dari 20%
13 df = df.drop(columns=columns_to_drop)
```

```
1 # Imputasi missing values kolom-kolom numerical
2
3 numeric_columns = df.select_dtypes(exclude='object').columns
4
5 for column in numeric_columns:
6     df[column].fillna(df[column].mean(), inplace=True)
```

Data Cleaning

```

1 # Imputasi missing values kolom-kolom kategorikal
2
3 categorical_columns = df.select_dtypes(include='object').columns
4
5 for column in categorical_columns:
6     df[column].fillna(df[column].mode()[0], inplace=True)

```

Data: First Five rows

| | loan_amnt | funded_amnt | funded_amnt_inv | int_rate | installment | grade | sub_g |
|---|-----------|-------------|-----------------|-----------|-------------|-------|-------|
| 0 | 5000 | 5000 | 4975.000000 | 10.650000 | 162.870000 | B | |
| 1 | 2500 | 2500 | 2500.000000 | 15.270000 | 59.830000 | C | |
| 2 | 2400 | 2400 | 2400.000000 | 15.960000 | 84.330000 | C | |
| 3 | 10000 | 10000 | 10000.000000 | 13.490000 | 339.310000 | C | |
| 4 | 3000 | 3000 | 3000.000000 | 12.690000 | 67.790000 | B | |

Data: Unique Value Counts In Each Column

| | loan_amnt | funded_amnt | funded_amnt_inv | int_rate | installment | grade | sub_g |
|--------------------|-----------|-------------|-----------------|----------|-------------|-------|-------|
| Unique Value Count | 1352 | 1354 | 9854 | 506 | 55622 | 7 | |

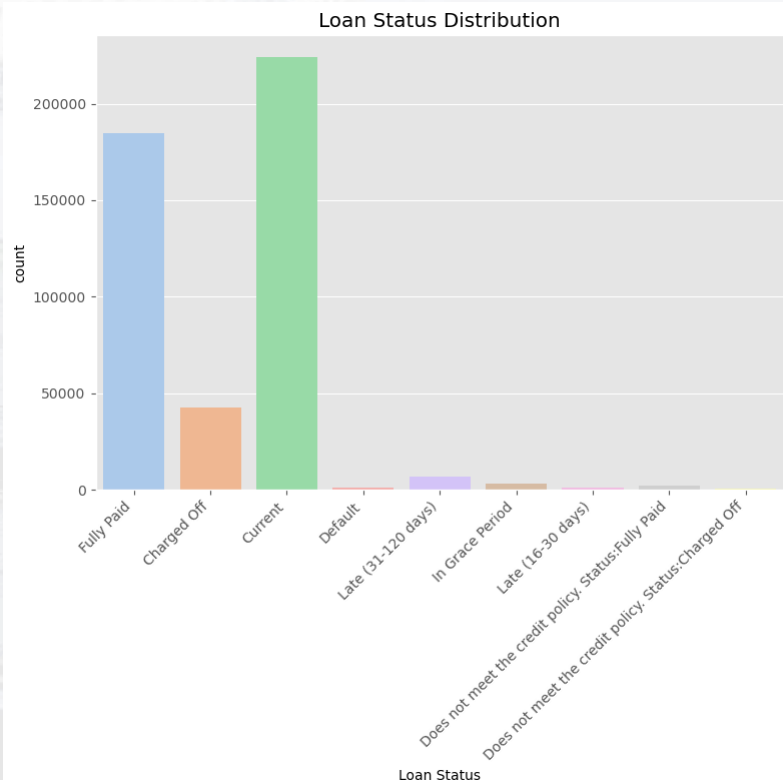
Data: Columns With Nan

| | loan_amnt | funded_amnt | funded_amnt_inv | int_rate | installment | grade | sub_g |
|---|-----------|-------------|-----------------|----------|-------------|-------|-------|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Exploratory Data Analysis



Loan Status Distribution



Distribusi Status Pinjaman

- **Current** sejumlah 224,226
- **Fully Paid** sejumlah 184,739
- **Charged Off** sejumlah 42,475
- **Late (31-120 days)** sejumlah 6,900
- **In Grace Period** sejumlah 3,146
- **Does not meet the credit policy. Status: Fully Paid** sejumlah 1,988
- **Late (16-30 days)** sejumlah 1,218
- **Default** sejumlah 832
- **Does not meet the credit policy. Status: Charged Off** sejumlah 761

Low Risk

Medium Risk

High Risk

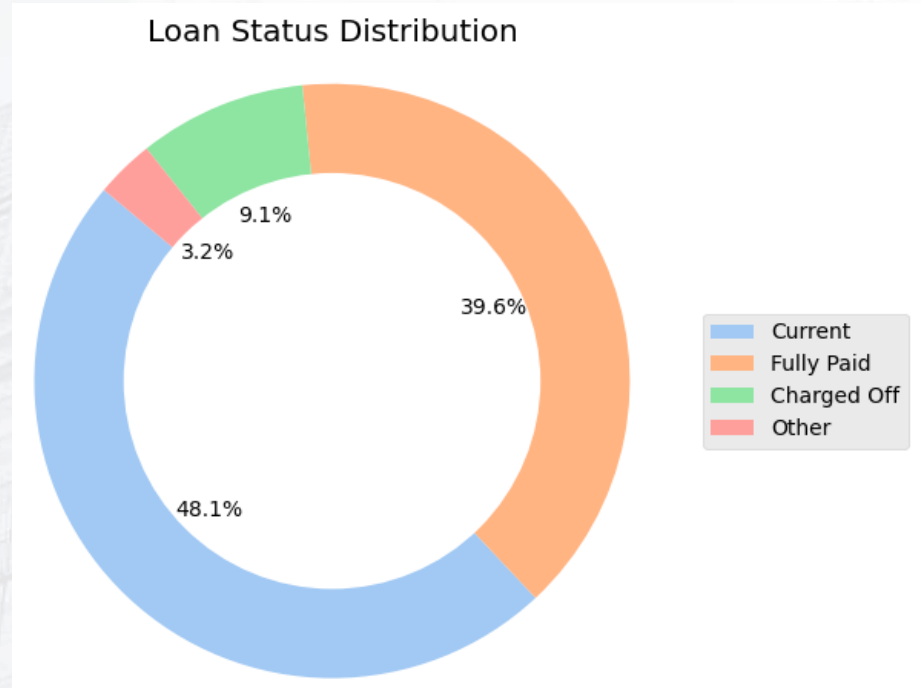
Loan Status Distribution

- **Current** sebesar 48.1%
- **Fully Paid** sebesar 39.6%
- **Charged Off** sebesar 9.1%
- Dan sisanya sebesar 3.2%

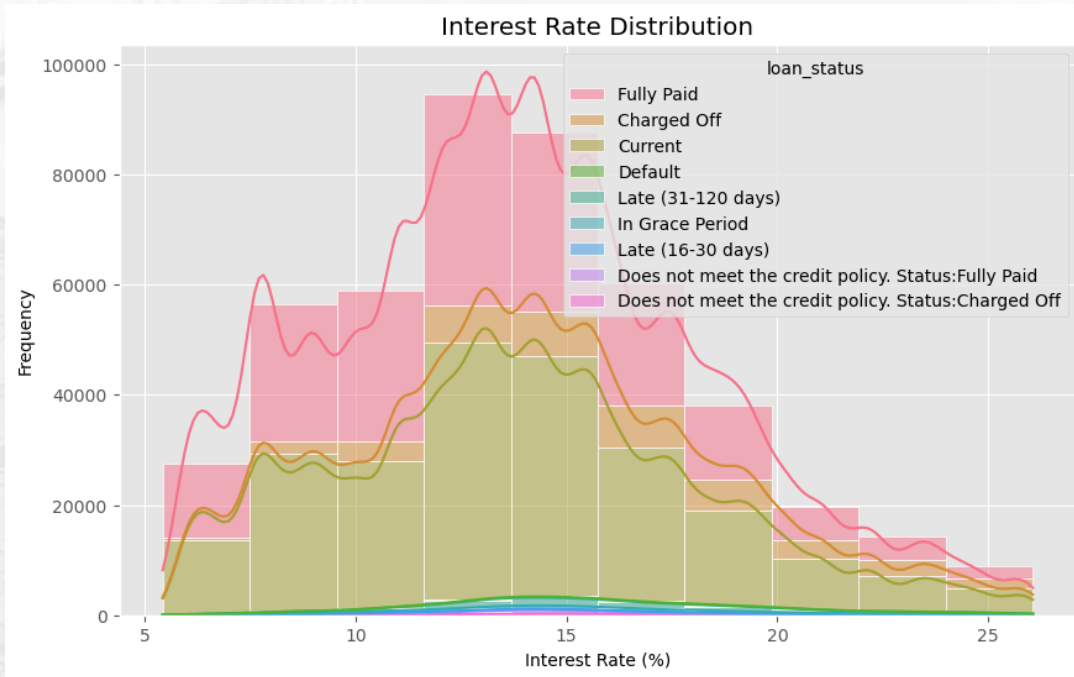
Low Risk

Medium Risk

High Risk

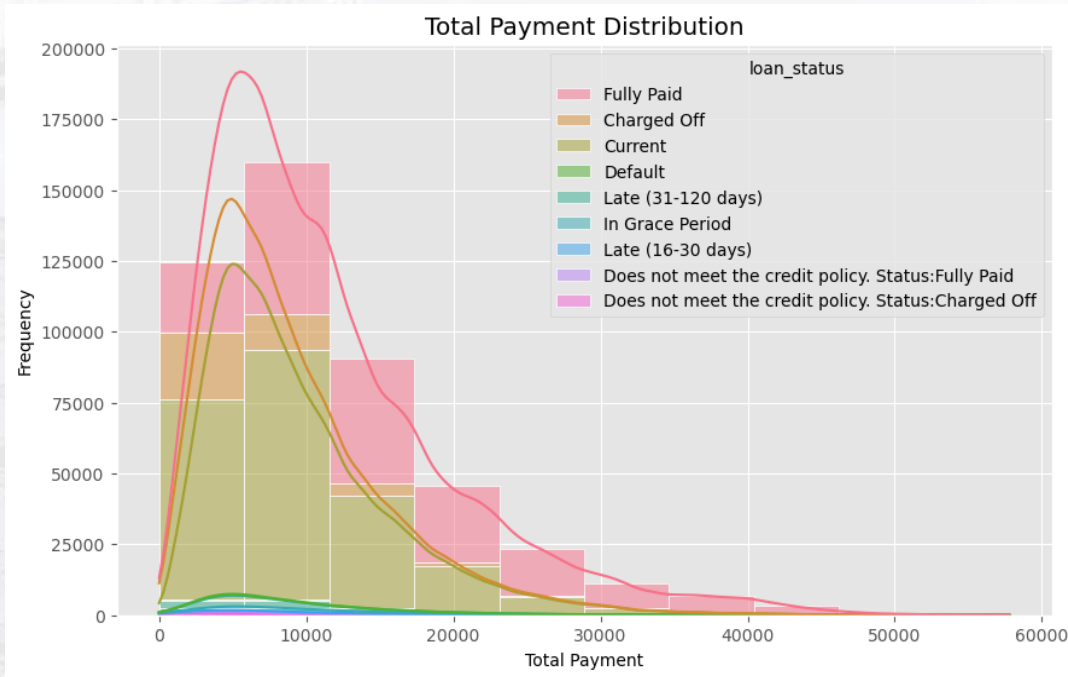


Interest Rate Distribution



- Analisis distribusi tingkat bunga pada dataset ini mengindikasikan bahwa sebagian besar tingkat bunga berkisar antara **12%** hingga **18%**.

Total Payment Distribution



- Analisis distribusi total pembayaran pada dataset ini menggambarkan variasi jumlah pembayaran. Terlihat bahwa puncak distribusi terjadi di sekitar **5,000 USD**.

Data Preprocessing



Data Preprocessing

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler, LabelEncoder
3 from sklearn.decomposition import PCA
```

```
1 X = df.drop('loan_status', axis=1)
2 y = df['loan_status']
```

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_
```

```
1 mapping = {label: index for index, label in enumerate(label_enc
2
3 print("Mapping hasil Label Encoding:")
4 for label, value in mapping.items():
5     print(f"{label} : {value}")
```

Mengimport preprocessing module yang diperlukan

Melakukan train-test split dengan variable target

```
Charged Off : 0
Current : 1
Default : 2
Does not meet the credit policy. Status:Charged Off : 3
Does not meet the credit policy. Status:Fully Paid : 4
Fully Paid : 5
In Grace Period : 6
Late (16-30 days) : 7
Late (31-120 days) : 8
```

Data Preprocessing

```
1 from sklearn.base import BaseEstimator, TransformerMixin
2
3 class MultiColumnLabelEncoder(BaseEstimator, TransformerMixin):
4     def __init__(self, columns=None):
5         self.columns = columns
6
7     def fit(self, X, y=None):
8         return self
9
10    def transform(self, X):
11        output = X.copy()
12        if self.columns is not None:
13            for col in self.columns:
14                output[col] = LabelEncoder().fit_transform(output[col])
15        else:
16            for colname, col in output.iteritems():
17                output[colname] = LabelEncoder().fit_transform(col)
18        return output
```

```
1 encoder = MultiColumnLabelEncoder(columns=X.select_dtypes(include=[object]))
2 X_train = encoder.fit_transform(X_train)
3 X_test = encoder.transform(X_test)
```

Melakukan encode label pada semua variabel kategorikal untuk memastikan konsistensi dalam pemrosesan data.

Data Preprocessing

```
1 X_num_cols = df.select_dtypes(exclude='object').columns
```

```
1 scaler = StandardScaler()  
2  
3 X_train[X_num_cols] = scaler.fit_transform(X_train[X_num_cols])  
4 X_test[X_num_cols] = scaler.transform(X_test[X_num_cols])
```

```
1 from imblearn.over_sampling import RandomOverSampler  
2  
3 ros = RandomOverSampler(sampling_strategy={8: 10000, 6: 10000},  
4  
5 X_train_oversampled, y_train_oversampled = ros.fit_resample(X_t  
6  
7 print("Distribusi kelas setelah Oversampling:", Counter(y_train
```

Melakukan standarisasi pada variabel-variabel numerik untuk memastikan data memiliki distribusi yang normal.

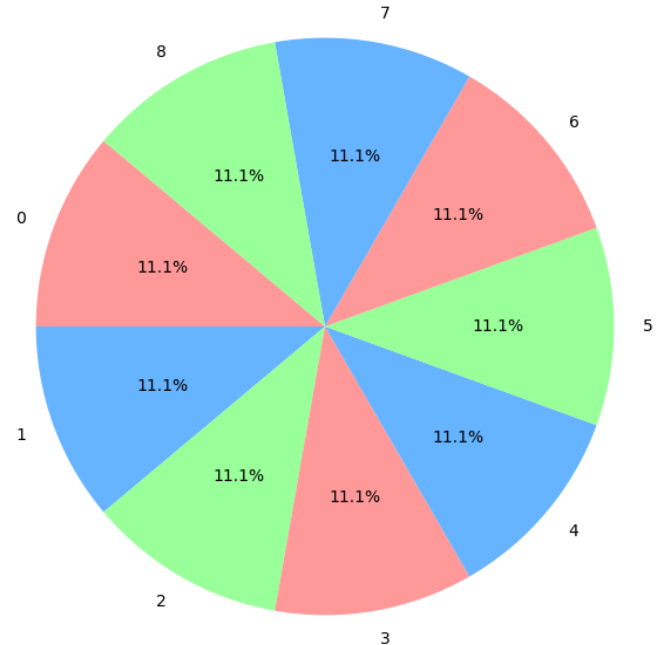
Melakukan oversampling untuk kelas-kelas minority pada variable target menjadi **10,000 sampel**

Data Preprocessing

```
1 from imblearn.under_sampling import RandomUnderSampler
2
3 rus = RandomUnderSampler(sampling_strategy={1: 10000, 5: 10000},
4
5 X_train_resampled, y_train_resampled = rus.fit_resample(X_train
6
7 print("Distribusi kelas setelah undersampling:", Counter(y_train
```

Melakukan undersampling untuk kelas-kelas majority pada variable target menjadi 10,000 sampel untuk mengatasi imbalance data.

Distribusi Kelas Setelah Resampling



Model Building and Evaluation



1. XGB Classifier

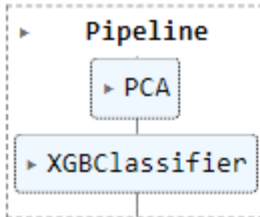
```
1 from sklearn.pipeline import make_pipeline
2 from sklearn.decomposition import PCA
3 from xgboost import XGBClassifier
```

Import modul yang diperlukan

```
1 pipeline = make_pipeline(
2     PCA(),
3     XGBClassifier(random_state=42, num_class=9)
4 )
```

Membangun model menggunakan pipeline

```
1 pipeline.fit(X_train_resampled, y_train_resampled)
```



```
1 y_pred = pipeline.predict(X_test)
```

XGB Classifier Model Evaluation

```
1 from sklearn.metrics import classification_report
2 print(classification_report(y_test, y_pred))
```

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.96 | 0.97 | 0.96 | 8495 |
| 1 | 0.97 | 0.82 | 0.89 | 44845 |
| 2 | 0.25 | 0.25 | 0.25 | 166 |
| 3 | 0.66 | 0.68 | 0.67 | 152 |
| 4 | 0.36 | 0.92 | 0.51 | 398 |
| 5 | 0.99 | 0.95 | 0.97 | 36948 |
| 6 | 0.04 | 0.34 | 0.08 | 629 |
| 7 | 0.02 | 0.10 | 0.03 | 244 |
| 8 | 0.26 | 0.73 | 0.38 | 1380 |
| accuracy | | | 0.88 | 93257 |
| macro avg | 0.50 | 0.64 | 0.53 | 93257 |
| weighted avg | 0.95 | 0.88 | 0.91 | 93257 |

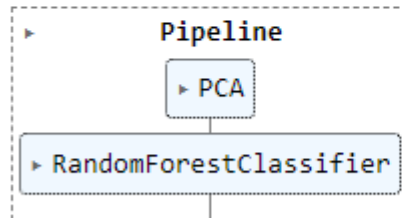
Model memiliki kinerja yang baik dalam mengklasifikasikan kategori risiko rendah (low risk) seperti **"Fully Paid"** dan **"Current"**, dengan presisi dan recall yang tinggi. Kinerja sedang terlihat pada kategori risiko menengah (medium risk) seperti **"Does not meet the credit policy. Status: Fully Paid"** dan **"In Grace Period"**. Model memiliki kinerja yang baik dalam mengklasifikasikan kategori risiko tinggi (high risk) seperti **"Charged Off"**, namun kinerja rendah untuk **"Default"**, dan **"Does not meet the credit policy. Status: Charged Off"**. Dan model ini memiliki performa yang buruk Dalam memprediksi status pinjaman yang mengalami keterlambatan.

2. Random Forest Classifier

```
1 from sklearn.ensemble import RandomForestClassifier
```

```
1 pipeline2 = make_pipeline(  
2     PCA(),  
3     RandomForestClassifier(random_state=42)  
4 )
```

```
1 pipeline2.fit(X_train_resampled, y_train_resampled)
```



```
1 y_pred_2 = pipeline2.predict(X_test)
```

Import modul yang diperlukan

Membangun model menggunakan pipeline

Random Forest Classifier Model Evaluation

```
1 from sklearn.metrics import classification_report
2 print(classification_report(y_test, y_pred_2))
```

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.91 | 0.92 | 0.91 | 8495 |
| 1 | 0.96 | 0.81 | 0.87 | 44845 |
| 2 | 0.00 | 0.00 | 0.00 | 166 |
| 3 | 0.75 | 0.37 | 0.49 | 152 |
| 4 | 0.32 | 0.82 | 0.46 | 398 |
| 5 | 0.94 | 0.93 | 0.94 | 36948 |
| 6 | 0.05 | 0.11 | 0.07 | 629 |
| 7 | 0.00 | 0.00 | 0.00 | 244 |
| 8 | 0.13 | 0.75 | 0.22 | 1380 |
| accuracy | | | 0.86 | 93257 |
| macro avg | 0.45 | 0.52 | 0.44 | 93257 |
| weighted avg | 0.92 | 0.86 | 0.88 | 93257 |

Model ini berhasil dengan baik dalam mengklasifikasikan risiko rendah (low risk) dengan tingkat presisi dan recall yang tinggi pada kategori **'Fully Paid'** dan **'Current'**. Kategori risiko menengah (medium risk) juga berhasil diidentifikasi, terutama pada **'Does not meet the credit policy. Status: Fully Paid'** dan **'In Grace Period'**. Kinerja model terlihat baik dalam mengklasifikasikan risiko tinggi (high risk) seperti **'Charged Off'**. Namun, terdapat penurunan kinerja pada **'Default'** dan **'Does not meet the credit policy. Status: Charged Off'**. Paling menonjol, model mengalami kesulitan memprediksi status pinjaman yang mengalami keterlambatan (**late status**) dengan akurasi yang rendah.



Connect With Me!



[yogaariasena \(Yoga Aria Sena\)](https://github.com/yogaariasena)
[\(github.com\)](https://github.com/yogaariasena)



[https://www.linkedin.com/in/yoga](https://www.linkedin.com/in/yogaariasena)
[aariasena](https://www.linkedin.com/in/yogaariasena)



Thank You



Rakamin
Academy



id/x partners