

# Credit Risk Prediction Model

ID/X Partners - Data Scientist

Presented by Muhammad Nazaruddin





### **Muhammad Nazaruddin**

### **Data Scientist**

Informatika kampus Gunadarma, dengan Lulusan ketertarikan mendalam pada Machine Learning dan Data Scientist. Saya memiliki sertifikasi resmi TensorFlow Developer dalam bidang Data Scientist. Saya terbiasa menggunakan tools seperti Python, TensorFlow, dan Darknet. Memiliki pengalaman dalam bekerja secara dalam tim maupun mandiri. Saya berpengalaman dalam ikut membuat projek Object Detection, System Recommendation, Sales Forecasting, Price Optimization dan selain itu memiliki pengalaman dalam AWS CLoud.

**Jakarta Barat** 





## **About Company**

id/x partners

ID/X Partners didirikan pada tahun 2002 oleh mantan bankir dan konsultan manajemen yang memiliki pengalaman luas dalam siklus dan proses kredit, pengembangan scoring, serta manajemen kinerja. Pengalaman gabungan mereka telah melayani berbagai perusahaan di kawasan Asia dan Australia, serta di berbagai industri seperti jasa keuangan, telekomunikasi, manufaktur, dan ritel.

ID/X Partners menyediakan layanan konsultasi yang mengkhususkan diri dalam pemanfaatan analitik data dan solusi pengambilan keputusan (DAD) yang dikombinasikan dengan disiplin manajemen risiko dan pemasaran terintegrasi untuk membantu klien mengoptimalkan profitabilitas portofolio dan proses bisnis mereka.

Layanan konsultasi yang komprehensif dan solusi teknologi yang ditawarkan menjadikan id/x partners sebagai penyedia layanan yang lengkap dan terpercaya



## **Project Portfolio**

Proyek ini bertujuan untuk memahami dan menganalisis data peminjaman (loan), dalam hal ini dataset yang digunakan adalah data pinjaman yang diberikan kepada anggota. Tujuan utama proyek ini adalah mengidentifikasi faktor-faktor yang berkontribusi terhadap pinjaman bermasalah (bad loans), yaitu pinjaman yang gagal bayar atau tidak dilunasi sesuai ketentuan. Oleh karena itu, tujuan analisis ini adalah:

- 1. Mengklasifikasikan mana pinjaman yang termasuk bad loan (gagal bayar) berdasarkan status pinjaman (loan\_status) menggunakan label target bad\_loan.
- 2. Menganalisis pola dan karakteristik peminjam yang cenderung memiliki bad loan.
- 3. Memberikan wawasan untuk membantu pengambilan keputusan lebih baik dalam pemberian pinjaman dengan memprediksi kemungkinan risiko kegagalan pembayaran.

## 1. Data Understanding

### **Membaca Dataset**

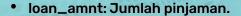
Dataset terdiri dari 466.285 baris dan 75 kolom dengan tipe data campuran (numerik dan kategorikal). Beberapa kolom memiliki banyak nilai kosong, bahkan ada yang sepenuhnya kosong. Informasi ini menunjukkan perlunya pembersihan data sebelum analisis lebih lanjut.

	Unnamed: 0	i	d member_id	loan_amnt	funded_amnt	funded_amnt_inv	term	int_rate	installment	grade	 total_bal_il
0	0	107750	1 1296599	5000	5000	4975.0	36 months	10.65	162.87	В	 NaN
1	1	107743	0 1314167	2500	2500	2500.0	60 months	15.27	59.83	С	 NaN
2	2	107717	5 1313524	2400	2400	2400.0	36 months	15.96	84.33	С	 NaN
3	3	107686	3 1277178	10000	10000	10000.0	36 months	13,49	339.31	С	 NaN
4	4	107535	8 1311748	3000	3000	3000.0	60 months	12.69	67.79	В	 NaN
5 ro	we v 75 col	lumne									



<class 'pandas.core.frame.DataFrame'> RangeIndex: 466285 entries, 0 to 466284 Data columns (total 75 columns):

#	Column	Non-Null Count	Dtype
0	Unnamed: 0	466285 non-null	int64
1	id	466285 non-null	int64
2	member_id	466285 non-null	int64
3	loan_amnt	466285 non-null	int64
4	funded_amnt	466285 non-null	int64
5	funded_amnt_inv	466285 non-null	float64
6	term	466285 non-null	object
7	int_rate	466285 non-null	float64
8	installment	466285 non-null	float64
9	grade	466285 non-null	object
10	sub_grade	466285 non-null	object
11	emp_title	438697 non-null	object
12	emp_length	445277 non-null	object
13	home_ownership	466285 non-null	object
14	annual_inc	466281 non-null	float64
15	verification_status	466285 non-null	object
16	issue_d	466285 non-null	9
17	loan_status	466285 non-null	_
18	pymnt_plan	466285 non-null	9
19	url	466285 non-null	
20	desc	125983 non-null	_
21	purpose	466285 non-null	object
22	title	466265 non-null	object
23	zip_code	466285 non-null	object
24	addr_state	466285 non-null	object
25	dti	466285 non-null	float64
26	delinq_2yrs	466256 non-null	
27	earliest_cr_line	466256 non-null	
28	inq_last_6mths	466256 non-null	
29	mths_since_last_delinq	215934 non-null	
30	mths_since_last_record	62638 non-null	
31	open_acc	466256 non-null	
32	pub rec	466256 non-null	float64



- funded\_amnt: Jumlah pendanaan yang disetujui.
- term: Durasi pinjaman (misalnya 36 atau 60 bulan).
- int\_rate: Suku bunga pinjaman.
- installment: Pembayaran angsuran bulanan.
- grade & sub\_grade: Klasifikasi risiko kredit dari peminjam.
- emp\_title, emp\_length: Informasi pekerjaan peminjam.
- home\_ownership: Jenis kepemilikan tempat tinggal.
- annual\_inc: Penghasilan tahunan peminjam.
- verification\_status: Status verifikasi penghasilan.
- purpose: Tujuan pinjaman.
- addr\_state: Negara bagian tempat peminjam tinggal.
- dti: Rasio utang terhadap penghasilan (Debt-to-Income Ratio).
- loan\_status: Status akhir dari pinjaman.





### Menentukan Target Variabel

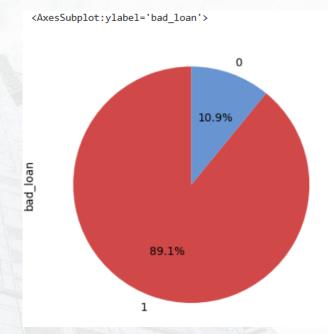
Kolom target untuk klasifikasi adalah loan\_status. Loan Status sendiri adalah status akhir dari pinjaman (seperti "Sudah membayar", "Belum/gagal membayar")

Perbandingan dari data yang sudah membayar dan belum membayar adalah:

10,9% sudah

89,1% belum

Dari Kolom target data bad\_loan ini maka hal yang menjadi tujuan utama adalah mengidentifikasi faktor-faktor yang memengaruhi pinjaman bermasalah (bad loans).





## 2. Feature Engineering

### **Melakukan Transformasi Variabel**

Fitur	Transformasi
bad_loan	Membuat kolom target berdasarkan loan_status .
term	Konversi dari string ke numerik.
emp_length	Konversi ke numerik dan penyimpanan dalam kolom baru <code>emp_length_int</code> .
<pre>issue_d, earliest_cr_line, last_pymnt_d, last_credit_pull_d</pre>	Konversi ke datetime dan pembuatan kolom baru dalam satuan bulan (mnth_since_*).
Missing Values	Diisi dengan mean (numerik) atau mode (kategorikal), serta penghapusan kolom dengan missing values > 40%.
Kolom tidak relevan	Dihapus dari dataset.

ini bertujuan untuk membersihkan data dan membuatnya lebih mudah digunakan dalam model prediksi atau analisis lanjutan.



## **Mengatasi Missing Variabel**

Jenis Data	Metode Imputasi	Alasan
Numerik	Mean	Menjaga stabilitas distribusi data
Kategorikal	Mode	Nilai yang paling sering muncul

Cara untuk mengatasi missing value pada dataframe bisa dengan melakukan imputasi yaitu mengisi nilai yang hilang dengan nilai tertentu.

#### Imputasi Data Hilang pada Kolom Numerik

Metode: Menggunakan rata-rata (mean).

Kenapa mean? Untuk menjaga distribusi nilai tetap stabil dan tidak mengganggu analisis statistik.

#### Imputasi Data Hilang pada Kolom Kategorikal

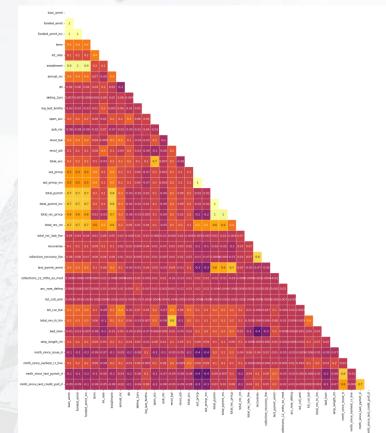
**Metode:** Menggunakan modus (*mode*), yaitu nilai yang paling sering muncul.

Kenapa mode? Cocok untuk data diskrit atau kategorikal karena merupakan nilai yang paling umum.

## 3. Exploratory Data Analysis

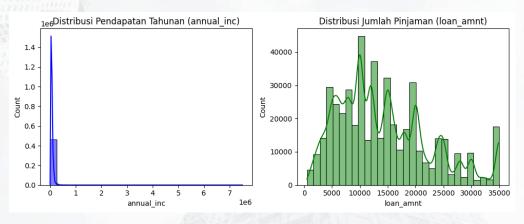
Rakamin

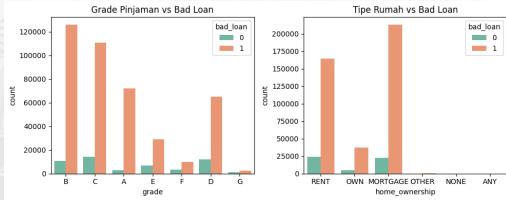
Memeriksa kolerasi antar variable, jika tidak terdapat kolerasi maka akan dihapus



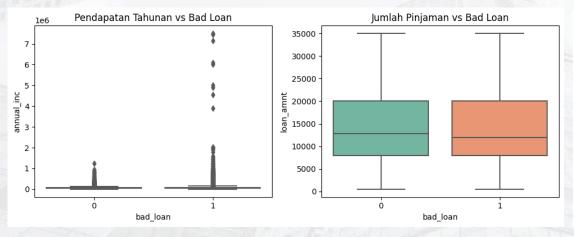


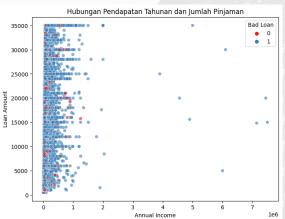




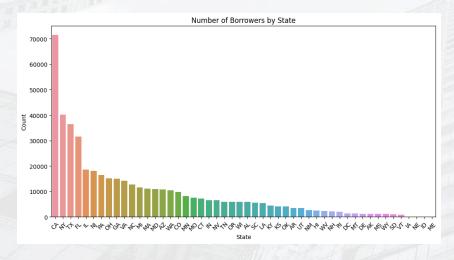


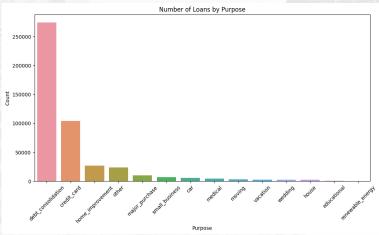














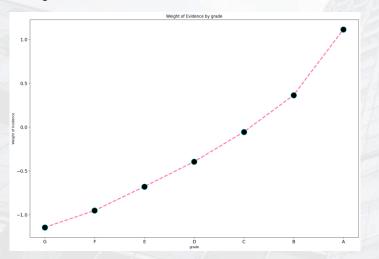


Weight of Evidence (WoE) digunakan untuk mengevaluasi hubungan antara variabel independen dengan target mengukur kekuatan suatu bin (kelompok data) dalam membedakan antara customer yang baik dan buruk. Nilai WoE < 0 menunjukkan bahwa bin variabel tersebut menangkap proporsi akun buruk yang lebih tinggi.

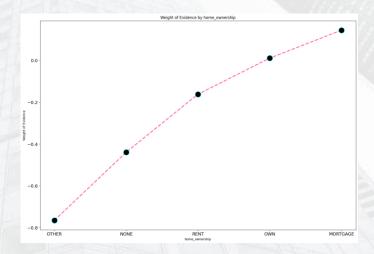
$$WoE = [ln(\frac{\text{Relative frequecy of Goods}}{\text{Relative frequecy of Bads}})] * 100$$

Information Value (IV) digunakan untuk mengevaluasi prediktifitas variabel independen terhadap target. IV (Information Value) adalah ukuran kekuatan prediktif keseluruhan dari suatu variabel dan sangat berguna untuk proses seleksi fitur.  $IV = \sum (DistributionGood_i - DistributionBad_i) * WoE_i$ 

· grade variable

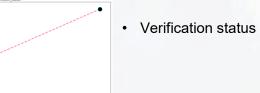


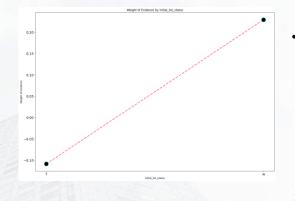
· home ownership variable

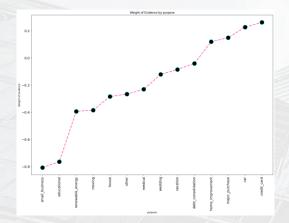




initial list status

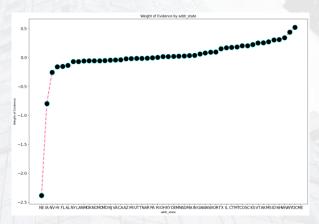




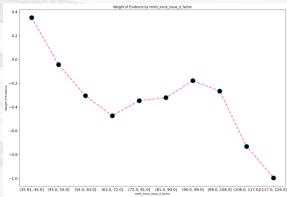


Source Verified verification status

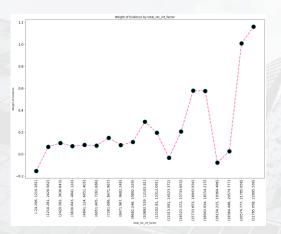
Purpose



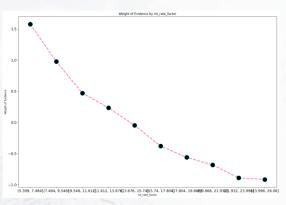
Address state



#### · months since issued

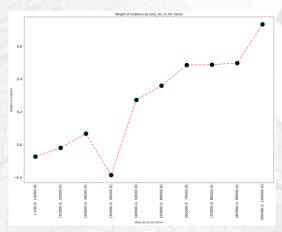


total\_rec\_int

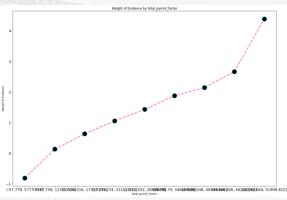


Rakamin Academy

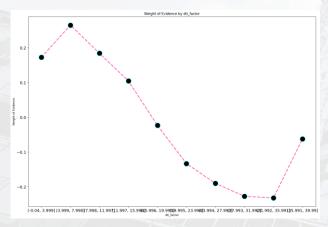
Interest rate



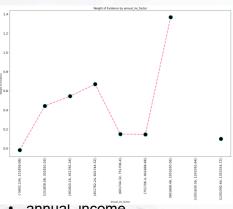
revolving\_high\_limit





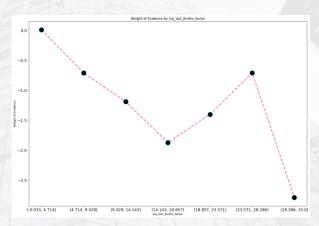


• dti

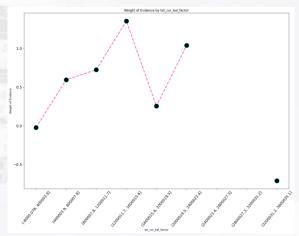




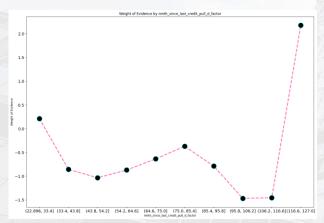
• annual\_income



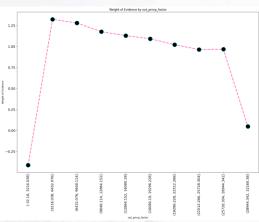
Income last 6 months



Total current balance

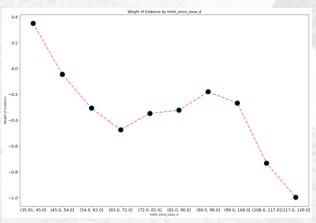


Months since last credit pulled





Outstanding principal



· months since last issued

## 4. Data Preparation

### **Drop Data yang memiliki Low IV**

Range: 0-1	Predictive powers
IV < 0.02	No power
0.02 < IV < 0.1	Weak power
0.1 < IV < 0.3	Medium power
0.3 < IV < 0.5	Strong power
0.5 < IV	Suspiciously high, too good to be true

Sebagai aturan, semua variabel dengan nilai IV < 0,02 dianggap tidak berguna untuk prediksi, sedangkan variabel dengan IV > 0,5 memiliki kekuatan prediktif yang mencurigakan. Oleh karena itu, variabel-variabel berikut tidak akan disertakan: out\_prncp, last\_pymnt\_amnt, delinq\_2yrs, mths\_since\_last\_delinq, open\_acc, pub\_rec, total\_acc, collections\_12\_mths\_ex\_med, acc\_now\_delinq, tot\_coll\_amt, dan mths\_since\_last\_pymnt\_d.



```
Information value of loan amnt is 0.004021
Information value of funded_amnt is 0.005335
Information value of funded amnt inv is 0.008519
Information value of term is 0.03886
Information value of int rate is 0.330892
Information value of installment is 0.00716
Information value of grade is 0.290782
Information value of home ownership is 0.021672
Information value of annual inc is 0.055095
Information value of verification status is 0.020831
Information value of purpose is 0.03698
Information value of addr state is 0.012518
Information value of dti is 0.026295
Information value of deling 2vrs is 4.5e-05
Information value of ing last 6mths is 0.036969
Information value of open acc is 0.000743
Information value of pub rec is 0.00058
Information value of revol bal is 0.006575
Information value of revol util is 0.027617
Information value of total acc is 0.007425
Information value of initial list status is 0.024801
Information value of out prncp is 0.76565
Information value of out_prncp_inv is 0.765591
Information value of total pymnt is 0.641949
Information value of total pymnt inv is 0.647435
Information value of total_rec_prncp is 1.490755
Information value of total rec int is 0.025457
Information value of total rec late fee is 0.0
Information value of recoveries is 6.156862
Information value of collection recovery fee is 0.0
Information value of last pymnt amnt is 1.549867
Information value of collections 12 mths ex med is 0.000593
Information value of acc now deling is 5.2e-05
Information value of tot coll amt is 0.031008
Information value of tot cur bal is 0.060884
Information value of total rev hi lim is 0.055222
Information value of emp length int is 0.006028
Information value of mnth since issue d is 0.11029
Information value of mnth since earliest cr line is 0.014745
Information value of mnth since last pymnt d is 2.036035
Information value of mnth since last credit pull d is 0.227195
```



### **Drop Not Use Variabel & Missing Value**

hapus data yang tidak relevan

```
col = ['id','member_id','url','sub_grade','zip_code']
data.drop(col, axis=1, inplace= True)
```

Menghapus variable yang tidak membantu dalam pelatihan model nantinya.

```
missing_values = data.isna().mean()*100
col_missingvalues = missing_values[missing_values > 40].index
col_missingvalues
```

hapus data dengan missingvalue mendekati 50%

Menghapus variable yang memiliki terlalu banyak missing value agar tidak menghasilkan prediksi yang ambigu



### Implementasi Fitur Baru Menggunakan Hasil WoE

Sebelumnya sudah dilakukan Eksplorasi Data Analisis mengunakan metode WoE, kemudian pada tahap ini Data frame dapat dikembangkan lagi dengan cara mengelompokan Variabel dan melakukan Encoding.

• **Grade**; mengubah nilai grade yang sebelumnya object menjadi data binnery



• **Home ownership**; Kategori OTHER, NONE, dan ANY memiliki sangat sedikit observasi dan sebaiknya digabungkan dengan kategori yang memiliki risiko gagal bayar tinggi, yaitu RENT.

home_ownership:OWN	home_ownership:OTHER_NONE_RENT_ANY	home_ownership:MORTGAGE
0	1	0
0	1	0
0	1	0



- verification\_status; membagi semua unique value menjadi kolom tersendiri NOT VERIFIED, SOURCE VERIFIED, dan VERIFIED
- purpose; mengelompokan kolom menjadi:
- SMALL BUSINESS EDUCATIONAL RENEWABLE ENERGY MOVING
- OTHER HOUSE MEDICAL
- WEDDING VACATION
- HOME IMPROVEMENT MAJOR PURCHASE
- CAR CREDIT CARD
- addr\_state; Dapat digabungkan menjadi satu kategori baru. Proses ini akan digunakan untuk analisis selanjutnya. Adapun kategori-kategori yang akan digabungkan adalah sebagai berikut:
- NE, IA, NV, HI, FL, AL
- NY
- LA, NM, OK, NC, MO, MD, NJ, VA
- CA
- AZ, MI, UT, TN, AR, PA
- RI, OH, KY, DE, MN, SD, MA, IN
- GA, WA
- WI, OR
- TX
- IL, CT,MT
- CO, SC
- KS, VT, AK, MS
- NH, WV, WY, DC



- initial\_list\_status; Memiliki nilai WoE yang sangat berbeda antar kategori, sehingga setiap kategori sebaiknya diperlakukan sebagai variabel terpisah F dan W.
- term; memisahkan variabel term yang hanya memiliki nilai 36 dan 60
- total\_rec\_int; Dibuat berdasarkan nilai WoE dan jumlah observasi
- <1000
- 1000-2000
- 2000-9000
- >9000
- total\_rev\_hi\_lim; Dibuat berdasarkan nilai WoE dan jumlah observasi
- <10000
- 10000-20000
- 20000-40000
- 40000-60000
- 60000-80000
- 80000-100000
- <100000

- out\_prncp; Dibuat berdasarkan nilai WoE dan jumlah observasi
- <3000
- 3000-6000
- 6000-10000
- 10000-12000
- >12000



- total\_pymnt; Dibuat berdasarkan nilai WoE dan jumlah observasi
- <5000
- 5000-11000
- 11000-16000
- 16000-22000
- >22000
- int\_rate; Dibuat berdasarkan nilai WoE dan jumlah observasi
- <7.484
- 7.484 9.548
- 9.548 11.612
- 11.612 13.676
- 13.676 15.74
- 15.74 17.804
- 17.804 19.868
- 7.19.868 21.932
- 21.932 26.06

- dti; Dibuat berdasarkan nilai WoE dan jumlah observasi
- (<4)
- (4-8)
- (8-12)
- -(12-16)
- -(16-20)
- -(20-23)
- -(23-27)
- -(27-40)
- annual\_inc; Dibuat berdasarkan nilai WoE dan jumlah observasi
- <32000
- 32000-50000
- 50000-60000
- 60000-75000
- 75000-90000
- 90000-120000
- 120000-135000
- 135000-150000
- >150000



## inq\_last\_6mths; Dibuat berdasarkan nilai WoE dan jumlah observasi

- <1
- 1-2
- 2-3
- 4-7
- tot\_cur\_bal; Dibuat berdasarkan nilai WoE dan jumlah observasi
- <40000
- 40000-80000
- 80000-120000
- 120000-160000
- 160000-200000
- 200000-240000
- 240000-320000
- 320000-400000
- >400000

- mnth\_since\_last\_credit\_pull\_d; Dibuat berdasarkan nilai WoE dan jumlah observasi
- <65
- 65 76
- >76
- mnth\_since\_issue\_d; Dibuat berdasarkan nilai WoE dan jumlah observasi
- (67.97, 70.8)
- (70.8, 73.6)
- (73.6-76.4)
- (76.4.- 79.2)
- (79.2 82)
- (82 84)
- (84 90.4)
- (90.4 96)



## 5. Data Modeling dan Evaluasi

#### Memisahkan data training dan data testing

```
X = df_train1.drop(columns='bad_loan', axis=1)
y = df_train1['bad_loan']

# Membagi data menjadi training dan testing set
x_train, x_test, y_train, y_test = train_test_split(X, y, random_state=42)

# Ubah semua -1.0 menjadi 0 agar kompatibel dengan semua model (terutama XGBoost)
y_train = y_train.replace(-1.0, 0)
y_test = y_test.replace(-1.0, 0)
```



## 5. Data Modeling dan Evaluasi

#### => Random Forest Model

```
# Random Forest Classifier
randof = RandomForestClassifier(random_state = 42)
randof.fit(x_train, y_train)
y_pred = randof.predict(x_test)
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0 1	0.74 0.92	0.34 0.99	0.46 0.95	12750 103822
accuracy macro avg weighted avg	0.83 0.90	0.66 0.91	0.91 0.71 0.90	116572 116572 116572

#### => Decision Tree Model

```
# Decision Tree Classifier
dtree = DecisionTreeClassifier(random_state = 42)
dtree = dtree.fit(x_train, y_train)
y_pred = dtree.predict(x_test)
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0 1	0.41 0.93	0.44 0.92	0.42 0.93	12750 103822
accuracy macro avg weighted avg	0.67 0.87	0.68 0.87	0.87 0.67 0.87	116572 116572 116572



### => K-Nearest Neighbors (KNN)

# K-Nearest Neighbors
knn = KNeighborsClassifier()
knn.fit(x\_train, y\_train)
y\_pred = knn.predict(x\_test)
print(classification\_report(y\_test, y\_pred))

	precision	recall	f1-score	support
0 1	0.56 0.91	0.24 0.98	0.33 0.94	12750 103822
accuracy macro avg weighted avg	0.73 0.87	0.61 0.90	0.90 0.64 0.88	116572 116572 116572

#### => XGBoost Model

# Now train the XGBoost model
xgb = XGBClassifier(random\_state=42)
xgb.fit(x\_train, y\_train)
y\_pred = xgb.predict(x\_test)
print(classification\_report(y\_test, y\_pred))

	precision	recall	f1-score	support
0 1	0.77 0.93	0.35 0.99	0.49 0.96	12750 103822
accuracy macro avg weighted avg	0.85 0.91	0.67 0.92	0.92 0.72 0.90	116572 116572 116572



#### => Logistic Regression Model

```
# Logistic Regression
from sklearn.linear_model import LogisticRegression
logreg = LogisticRegression(max_iter=1000, random_stalogreg.fit(x_train, y_train)
y_pred = logreg.predict(x_test)
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0 1	0.74 0.92	0.31 0.99	0.43 0.95	12750 103822
accuracy macro avg weighted avg	0.83 0.90	0.65 0.91	0.91 0.69 0.90	116572 116572 116572

#### Mengecek overfit dan underfit dari model Logistic Regression

```
: X_train2 = x_train
X_test2 = x_test

logreg = RandomForestClassifier()
logreg.fit(X_train2, y_train)
y_pred_rf = logreg.predict(X_train2)
y_pred_rf_test = logreg.predict(X_test2)

print('Akurasi Train',accuracy_score(y_train, y_pred_rf))
print('Akurasi Test',accuracy_score(y_test, y_pred_rf_test))

Akurasi Train 0.9981584899617687
Akurasi Test 0.914130322890574

model tidak mengalami overfitting maupun underfitting
```



ModelAUC ScoreXGBoost0.98Random Forest0.97Logistic Regression0.95K-Nearest Neighbors0.93Decision Tree0.87



# 7. Conclusion

Model	Akurasi	F1-score (Kelas 0)	F1-score (Kelas 1)	Catatan
XGBoost	0.92	0.49	0.96	Paling seimbang dan akurat, <b>model</b> <b>terbaik</b>
Random Forest	0.91	0.46	0.95	Akurat untuk kelas 1, lemah di kelas 0
Logistic Regression	0.91	0.43	0.95	Mirip Random Forest, tapi lebih ringan
KNN	0.90	0.33 (terendah)	0.94	Buruk untuk kelas 0, tidak cocok untuk data besar
Decision Tree	0.87	0.42	0.93	Performa keseluruhan paling rendah

# **Thank You**





Logo Company