

# Prediction Model

## ID/X Partners - Data Scientist

Presented by  
Hamzah Mulyana

## Hamzah Mulyana

### Data Scientist



Data Scientist with under one year of experience and certified in Data Science & Machine Learning from Purwadhika Digital Technology School. Skilled in data collection, cleaning, processing, and visualization using Python (Pandas, Numpy, Matplotlib, Seaborn) and Microsoft Excel (Basic Formula, IF Condition, Lookup Value, Pivot Table), along with tools like Tableau, Google Looker Studio. Proficient in MySQL, Spreadsheets, and Microsoft Office, with statistical analytics capabilities to support data-driven decision making. Detail-oriented and committed to applying technical skills to solve business problems through data-driven approaches.



Cirebon, Jawa Barat



[hamzahmulyana88@gmail.com](mailto:hamzahmulyana88@gmail.com)



[Hamzah Mulyana](#)

# Courses and Certification

**Linear Models in Machine Learning | [Link](#)**

**May, 2025**

**Purwadhika Digital Technology School | [Link](#)**

**Oct, 2024**

**Machine Learning | [Link](#)**

**Aug, 2024**

**Data Analysis | [Link](#)**

**Jul, 2024**

# About Company

**ID/X Partners** didirikan pada tahun 2002 oleh *ex-banker* dan *management consultants* yang memiliki pengalaman dalam *credit cycle and process management, scoring development, and performance management*. Pengalaman gabungan kami telah melayani perusahaan-perusahaan di wilayah Asia dan Australia di berbagai industri, khususnya jasa keuangan, telekomunikasi, manufaktur, dan ritel.

**ID/X Partners** menyediakan layanan konsultasi yang mengkhususkan diri dalam pemanfaatan solusi *Data Analytic and Decisioning (DAD)*, dikombinasikan dengan pendekatan terpadu dalam manajemen risiko dan pemasaran, untuk membantu klien mengoptimalkan profitabilitas portofolio dan proses bisnis mereka.

Dengan layanan konsultasi yang komprehensif dan solusi teknologi yang ditawarkan, id/x partners menjadi penyedia layanan *one-stop solution* bagi kebutuhan bisnis klien.

# Project Portfolio

**Proyek ini bertujuan membangun model machine learning untuk memprediksi risiko kredit pada perusahaan multifinance guna meningkatkan akurasi penilaian kelayakan peminjam dan mengurangi potensi kredit macet. Dataset yang digunakan mencakup riwayat pinjaman beserta berbagai atribut terkait seperti profil nasabah, riwayat pembayaran, dan karakteristik pinjaman. Solusi ini diharapkan dapat membantu perusahaan dalam mengoptimalkan proses pengambilan keputusan pemberian pinjaman sekaligus meminimalkan kerugian akibat gagal bayar, dengan menghasilkan model prediktif yang akurat beserta insight bisnis yang relevan untuk mendukung strategi penyaluran kredit yang lebih efektif.**



# Data Understanding



# Data Understanding

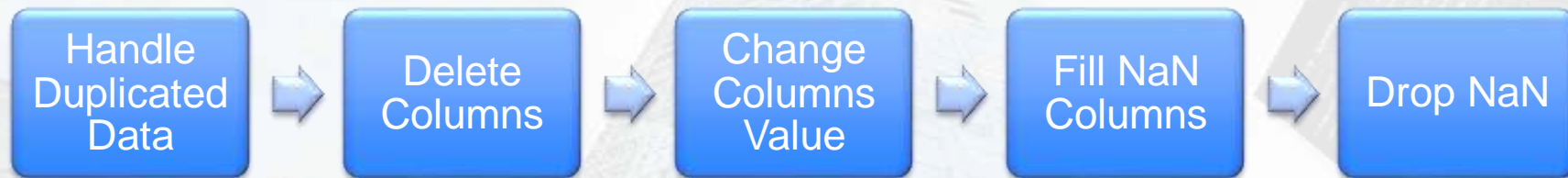
## Data Size

- Data terdiri dari 466285 baris dan 74 kolom

## Data Condition

- Data types: Object, Float, Int
- Missing Value: Lebih dari 10 kolom tidak memiliki Nilai
- Duplicated Data: Tidak ada duplikat data

# Data Preparation





# Data Preparation

## Handle Duplicated Data

- Tidak terdapat duplikat dalam data

## Delete Columns

- Menghapus kolom yang tidak digunakan untuk analisis dan prediksi, seperti:

*'id', 'member\_id', 'funded\_amnt', 'funded\_amnt\_inv', 'url', 'desc', 'zip\_code', 'title',  
'mths\_since\_last\_delinq', 'mths\_since\_last\_record', 'mths\_since\_last\_major\_derog',  
'policy\_code', 'application\_type', 'annual\_inc\_joint', 'dti\_joint', 'verification\_status\_joint',  
'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', 'pymnt\_plan', 'out\_prncp\_inv',  
'total\_pymnt\_inv', 'sub\_grade', 'earliest\_cr\_line', 'emp\_title', 'last\_credit\_pull\_d'*

# Data Preparation

## Change Columns Value

- **loan\_status:** Fully Paid, Current, In Grace Periode -> Diterima (1)
- **verification\_status:** Source Verified -> Verified
- **purpose:** house -> home\_improvement
- **home\_ownership:** ANY -> OTHER

## Fill NaN

- Mengisi nilai dengan menggunakan metode modus(mode) untuk kolom:  
'delinq\_2yrs', 'inq\_last\_6mths', 'open\_acc', 'pub\_rec', 'revol\_util', 'total\_acc',  
'collections\_12\_mths\_ex\_med', 'acc\_now\_delinq'

## Drop NaN

- Menghapus baris yg memiliki nilai NaN karena tidak ada acuan untuk mengisinya

# Exploratory Data Analysis

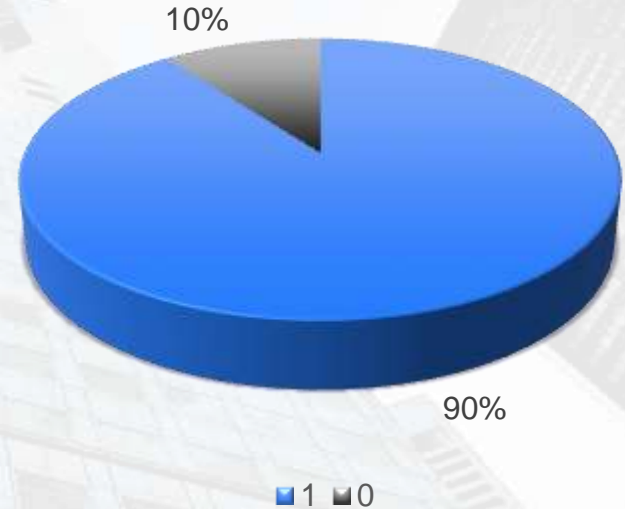


Target  
Distribution

Numerical  
Distribution

Categorical  
Distribution

# Exploratory Data Analysis



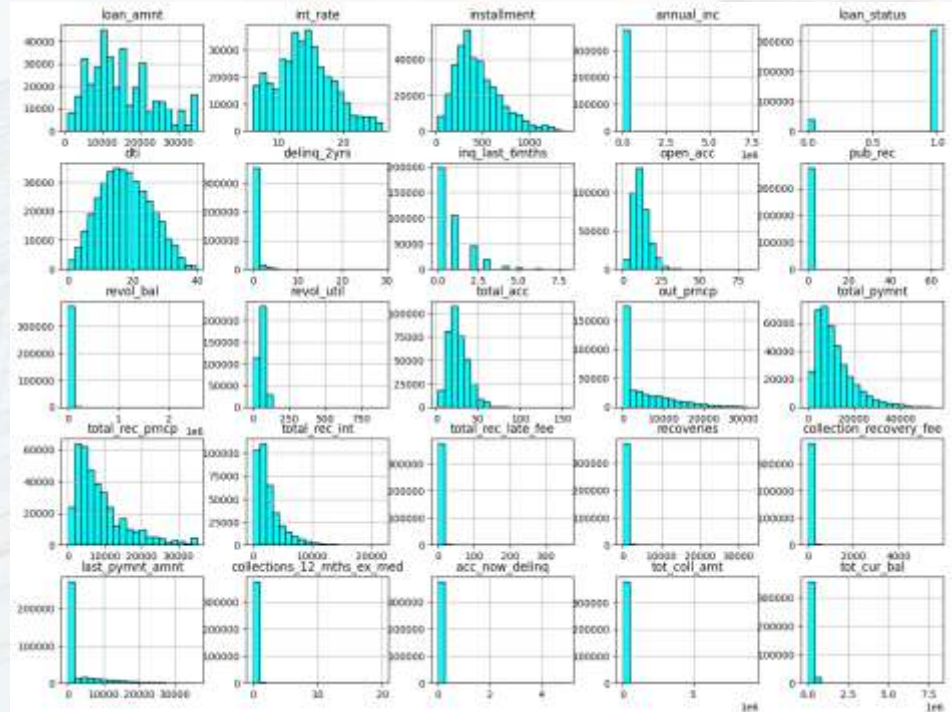


# Exploratory Data Analysis

id/x partners

Dapat dilihat pada visualisasi disamping, untuk data numerical setiap fitur/kolom memiliki distribusi data yang tidak normal, maka dari itu analisis statistiknya menggunakan Non Parametrik.

Setelah dilakukan uji Non Parametrik ada kolom yang tidak akan dimasukkan dalam Feature Engineering.





# Exploratory Data Analysis

## Categorical Feature

- Untuk categorical feature, hasilnya sebagai berikut:

**term -> 36 months**

**grade -> B**

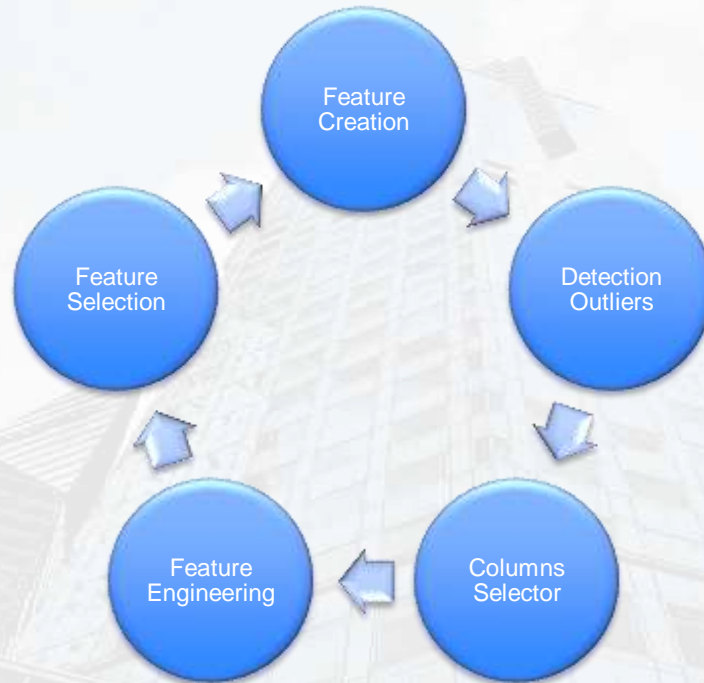
**emp\_length -> 10+ years**

**home\_ownership -> mortgage**

**verification\_status -> verified**

**purpose -> debt\_consolidation**

# Feature Engineering



# Feature Engineering

## Feature Creation

- Membuat fitur tambahan untuk analisis, kasus kali ini yaitu memisahkan bulan dan tanggal kedalam kolom berbeda

## Outliers Detection

- Untuk kolom yang memiliki nilai outliers tinggi, maka akan dilakukan penanganan dalam Outliers Detection ini, yaitu akan diubah ke rentang percentil 10 sampai percentil 90

## Columns Selector

- Untuk kolom yang tidak akan dilakukan apa-apa maka Columns Selector ini dapat digunakan, alasannya karena ada kolom yang tidak perlu dilakukan feature engineering maka dari itu dibuatlah Columns Selector ini.

# Feature Engineering

## Feature Engineering

- Onehot Encoding, Ordinal Encoding, Scaling (Robust), Simple Imputer

## Feature Selection

- Memilih feature yang berpengaruh untuk dilakukan modeling, pemilihan feature ini menggunakan Feature Importance pada algoritma Random Forest. Didapatkan hasil feature-feature yang akan diambil sebagai berikut:

*'mnth\_last\_pymnt\_d', 'mnth\_next\_pymnt\_d', 'loan\_amnt', 'total\_rec\_prncp',  
'collection\_recovery\_fee', 'last\_pymnt\_amnt', 'recoveries', 'total\_pymnt', 'out\_prncp',  
'installment', 'total\_rec\_int', 'total\_rec\_late\_fee', 'int\_rate', 'day\_next\_pymnt\_d',  
'day\_last\_pymnt\_d'*

# Data Modeling

Linear

- Logistic Regression

Ensemble

- Decision Tree
- Random Forest
- AdaBoost Classifier

Cross Validation



Hyperparameter  
Tuning



# Data Modeling

## Before Tunning (Cross Validation)

Nama Model	Precision	
	Train	Test
Logistic Regression	98%	98%
Decision Tree	99%	99%
Random Forest	99%	99%
AdaBoost Classifier	99%	99%

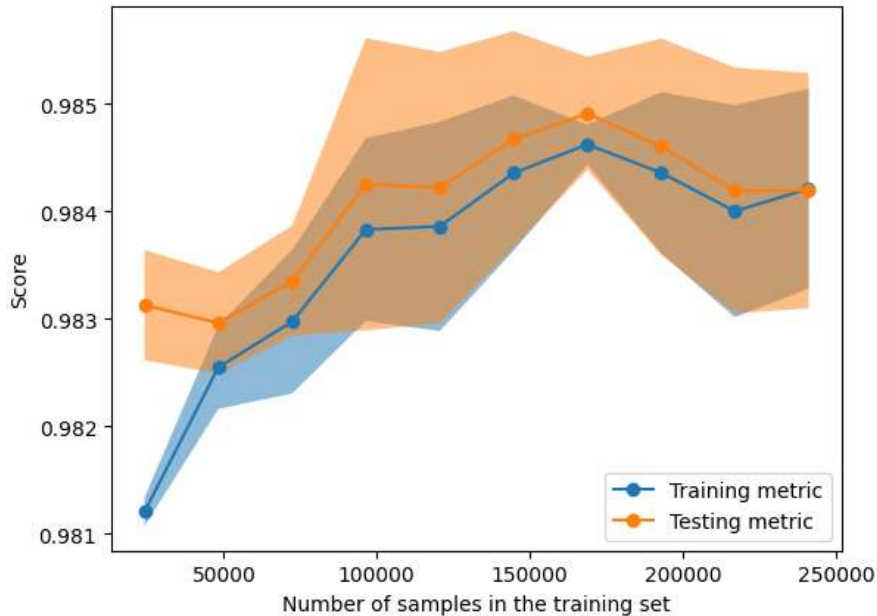
# Data Modeling

## After Tunning (Hyperparameter Tunning)

Nama Model	Precision	
	Train	Test
Logistic Regression	95%	95%
Decision Tree	99%	99%
Random Forest	99%	99%
AdaBoost Classifier	99%	99%

# Evaluation

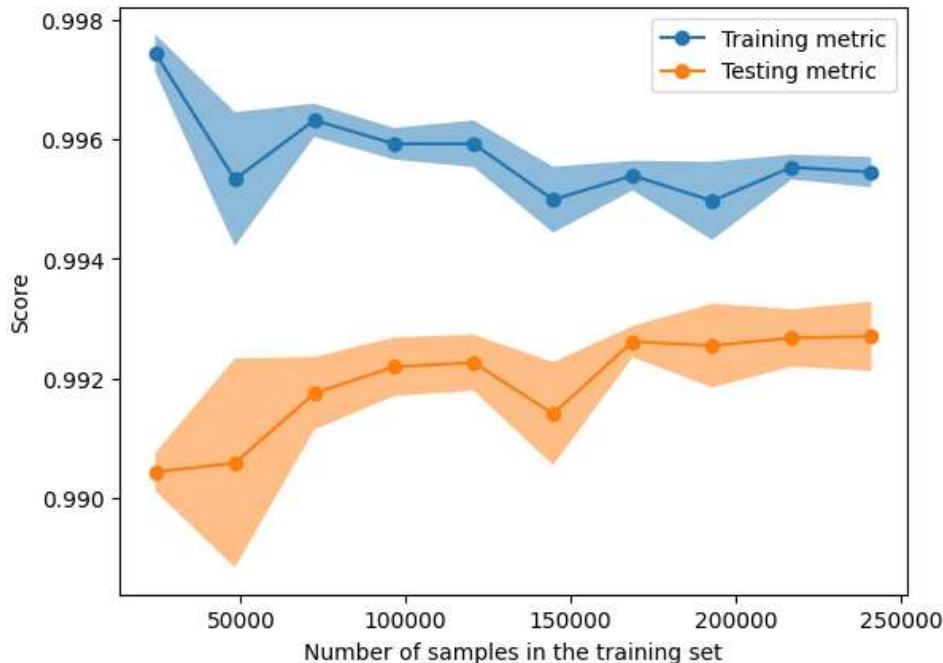
## Logistic Regression



Dapat dilihat Learning Curve disamping, data training dan data testing menunjukkan hasil yang bagus. Hasil dari data training dan data testing tidak terlalu jauh.

# Evaluation

## Decision Tree



Dapat dilihat Learning Curve disamping, data training dan data testing menunjukkan hasil yang tidak terlalu bagus (indikasi adanya overfitting). Hasil dari data training dan data testing jauh.

# Conclusion

## Logistic Regression

- Dengan precision 95%, model dapat memprediksi bahwa orang yang diprediksi layak mendapatkan pinjaman memang benar layak (true positive). Ini menunjukkan bahwa model sangat baik dalam meminimalkan Type I Error (False Positive) hingga 5%. Threshold yang digunakan sangat ketat, yaitu 0.16, sehingga hanya prediksi dengan probabilitas tinggi yang diklasifikasikan sebagai 'layak'.

## Decision Tree

- Dengan precision 98%, model dapat memprediksi bahwa orang yang diprediksi layak mendapatkan pinjaman memang benar layak (true positive). Ini menunjukkan bahwa model sangat baik dalam meminimalkan Type I Error (False Positive) hingga 2%. Threshold yang digunakan sangat ketat, yaitu 0.05, sehingga hanya prediksi dengan probabilitas tinggi yang diklasifikasikan sebagai 'layak'.



# Thank You

