

BAB II

TINJAUAN PUSTAKA

2.1. Tinjauan Pustaka

Banyak penelitian terdahulu yang menggunakan algoritma *Random Forest* dalam berbagai konteks, terutama risiko kredit. Penerapan *Random Forest* pada klasifikasi risiko kredit menggunakan arsitektur yang standar (menggunakan parameter *default*) sehingga mendapatkan akurasi rendah. Namun ada beberapa penelitian yang menggunakan *feature engineering* dan *hyperparameter tuning* pada klasifikasi risiko kredit menunjukkan hasil yang lebih baik. Akan tetapi penelitian tersebut belum diterapkan pada data lain, seperti data risiko kredit atau “credit risk data” yang diunggah oleh Lao Tse [34]. Penelitian terkait *Random Forest* pada domain lain menunjukkan bahwa penerapan *feature engineering* dan *hyperparameter tuning* dapat meningkatkan keakuratan model *Random Forest*. Dalam beberapa penelitian, teknik-teknik tersebut berhasil mengoptimalkan performa model dengan memilih fitur-fitur yang paling informatif dan mengatur parameter-parameter yang tepat, sehingga meningkatkan prediksi dan keandalan model secara signifikan. Dengan memanfaatkan penelitian-penelitian sebelumnya, dapat diperoleh wawasan yang berharga dalam mengimplementasikan algoritma *Random Forest* dengan akurasi yang lebih tinggi. Berikut ini beberapa penelitian terdahulu yang menjadi rujukan penelitian ini.

Penelitian [10] tentang bagaimana menggunakan algoritma *Random Forest* untuk memprediksi kelayakan kredit secara dini. Keberagaman dan banyaknya fitur pada data serta adanya ketidak seimbangan kelas membuat algoritma *Random Forest* diterapkan pada penelitian ini. Penelitian ini bertujuan untuk mengetahui penerapan dan hasil terbaik dari algoritma *Random Forest* untuk klasifikasi risiko kredit pada data *German Credit Risk*. Data tersebut memiliki 1000 baris dengan 21 fitur. Pada penelitian ini mendapatkan akurasi model 83%.

Penelitian [13] yang bertujuan untuk membandingkan akurasi model *Random Forest*, *Deep Neural Network*, dan Adaboost. Penelitian ini menerapkan *hyperparameter tuning* pada model karena dianggap dapat membantu

meningkatkan akurasi, stabilitas, serta menangani masalah *underfitting* maupun *overfitting* pada data yang terdapat banyak fitur. Penelitian ini berhasil mendapatkan akurasi tertinggi pada *Random Forest* sebesar 90.63%. Namun pada penelitian ini kurang menjelajahi parameter lain pada *Random Forest* untuk dilakukan *hyperparameter tuning*.

Dilanjutkan dengan penelitian [9] yang melihat bahwa penelitian sebelum – sebelumnya yang tidak menerapkan *hyperparameter tuning* pada model *Random Forest*. Sehingga penelitian ini bertujuan untuk menerapkan *hyperparameter tuning* dan penyetelan strategi karena beberapa parameter di *Random Forest* dipercaya dapat mempengaruhi kinerjanya, parameter yang kurang tepat akan mempengaruhi hasil akurasi model. Penelitian ini memaparkan bahwa dengan melakukan *hyperparameter tuning* pada parameter *mtry*, *sample size*, *replacement*, *node size*, *number of trees*, dan *splitting rule* maka nilai *mtry* memiliki pengaruh besar terhadap keakuratan model. Sedangkan nilai *sample size* dan *node size* berpengaruh kecil terhadap performanya.

Didukung penelitian [14] yang menerapkan *hyperparameter tuning* model *Random Forest* pada data perusahaan di China dengan jumlah 3111 dan tidak seimbang (*imbalanced data*). Data yang berdimensi tinggi dan *imbalanced* membuat algoritma *Random Forest* diterapkan pada penelitian ini, namun jumlah parameter yang tidak tepat dianggap krusial pada hasil evaluasi model. Sehingga penelitian ini menerapkan *hyperparameter tuning* pada beberapa parameternya. Penerapan *hyperparameter tuning* ini berhasil meningkatkan akurasi model menjadi 99.67% dan nilai AUC yang didapat yaitu 0.9998, artinya model dapat membedakan kelas 0 dan 1 dengan baik.

Penelitian [11] tentang bagaimana algoritma optimasi pada *Random Forest* dapat meningkatkan akurasi model dan mampu mendeteksi kegagalan jantung pasien. Penelitian ini bertujuan untuk membandingkan 3 algoritma optimasi pada *Random Forest* yaitu *Grid Search*, *Random Search*, dan *Bayesian Search* untuk menghasilkan akurasi terbaik karena pada data yang mempunyai banyak fitur, akurasi model *Random Forest* yang didapat masih rendah. Penelitian ini menunjukkan hasil bahwa optimasi terbaik menggunakan *Random Search* dengan

akurasi tertinggi sebesar 85.63%.

Didukung penelitian [15] tentang bagaimana penerapan *hybrid* dari *Random Forest* dan *Grid Search* dalam menyetel *hyperparameter tuning* pada model. Pencarian akurasi menggunakan data *training* dapat menghasilkan *overfitting*, sehingga perlu dilakukan penyesuaian *hyperparameter*. Penelitian ini bertujuan melakukan pendekatan *hybrid* algoritma *Random Forest* dan *Grid Search* untuk menghasilkan akurasi model yang optimal. *Hyperparameter tuning* dilakukan pada jumlah *n_estimators* dan *max_features* sehingga berhasil meningkatkan akurasi model baseline menjadi 90.02%.

Didukung penelitian [16] yang memaparkan tentang bagaimana penerapan *feature engineering* dan *hyperparameter tuning* untuk meningkatkan akurasi model *Random Forest* dan menghasilkan model terbaik. Data yang berdimensi tinggi mempunyai kecenderungan *overfitting*, sehingga perlu dilakukan *feature selection*. Selain itu, akurasi model *baseline* yang masih rendah membuat penelitian ini menerapkan *hyperparameter tuning*. Penelitian ini melakukan *hyperparameter tuning* pada nilai *n_estimators*, *min_sample_leaf*, *max_features*, dan *max_depth* pada data berjumlah 12520 baris. Hasil akurasi model meningkat menjadi 71.88%.

Penelitian [17] tentang bagaimana penerapan hasil *hyperparameter tuning* menggunakan *Grid Search* untuk mendapatkan akurasi terbaik. Model yang terlalu kompleks membuat model tersebut *overfitting*, sehingga perlu dilakukan *hyperparameter tuning*. Penelitian ini bertujuan untuk mendapatkan akurasi model *Random Forest* yang paling tepat dan optimal. *Hyperparameter tuning* dilakukan pada parameter *max_depth*, *max_number_trees*, *weight parameter* dan *feature evaluation*. Penelitian ini menghasilkan akurasi yang sangat baik yaitu 97%.

Penelitian [18] memaparkan kondisi data yang tidak seimbang (*imbalanced*) dapat mempengaruhi proses *training* dan menyebabkan bias yang tidak menguntungkan terhadap kelas mayoritas. Algoritma yang digunakan yaitu *Random Forest* juga sensitif terhadap parameter – parameter yang digunakan. Untuk itu diperlukan penyeimbangan data ke kelas mayoritas atau *oversampling* dan *hyperparameter tuning* model *Random Forest* untuk mendapatkan hasil akurasi yang lebih baik. Hasil penelitian ini yaitu adanya peningkatan nilai akurasi dari 94%

menjadi 96% setelah menerapkan *oversampling* dan *hyperparameter tuning* menggunakan *Grid Search CV*.

Penelitian [19] yang menggunakan teknik *Random Oversampling* untuk mengatasi *imbalanced* data model *Random Forest* untuk klasifikasi risiko kehamilan. Penelitian ini mendapatkan hasil bahwa model yang dibuat setelah menerapkan *Random Over Sampling* mengalami kenaikan 1.09% pada akurasi model *Random Forest* menjadi 81.86%. Ringkasan 10 penelitian sebelumnya yang dijadikan referensi dalam penelitian ini terdapat pada Tabel 2.1.

Tabel 2.1 Penelitian Terdahulu

Literature Review		Masalah Penelitian/ Rumusan Masalah	Tujuan	Metode	Data	Hasil / Temuan/ Kesimpulan
Penulis, Tahun	Judul					
Budi Prasajo, Emy Haryatmi, (2020)	Analisa Prediksi Kelayakan Pemberian Kredit Pinjaman dengan Metode <i>Random Forest</i>	Keberagaman dan banyaknya fitur data pemberian kredit dan adanya kelas yang tidak seimbang (<i>imbalanced</i>), sehingga <i>Random Forest</i> tepat digunakan pada data yang ada.	Mengetahui penerapan dan hasil akurasi terbaik menggunakan algoritma <i>Random Forest</i> pada data <i>German Credit</i> .	Klasifikasi <i>Random Forest</i>	<i>German Credit</i> Data (1000 baris dengan 21 fitur) dan <i>imbalanced class</i>	Akurasi sebesar 83%.
Joseph Sanjaya, Erick Renata, Vincent Elbert Budiman, Francis Anderson, Mewati Ayub, (2020)	Prediksi Kelalaian Pinjaman Bank Menggunakan <i>Random Forest</i> dan <i>Adaptive Boosting</i>	Adanya fitur data yang beragam dan data yang tidak seimbang (<i>imbalanced</i>) sehingga penelitian ini menggunakan <i>Random Forest</i> , DNN, dan <i>Adaboost</i> .	Membandingkan akurasi dan performa model <i>Random Forest</i> , DNN, dan <i>Adaboost</i> pada data dengan fitur yang banyak.	<i>Random Forest</i> , <i>Adaboost</i> , dan DNN dengan melakukan <i>Hyperparameter Tuning</i> .	Data pinjaman (<i>loan</i>) yang diperoleh dari salah satu bank di Indonesia dengan 20 fitur.	Akurasi tertinggi sebesar 90.63% dengan model <i>Random Forest</i> dan dilakukan <i>undersampling</i> .

Literature Review		Masalah Penelitian/ Rumusan Masalah	Tujuan	Metode	Data	Hasil / Temuan/ Kesimpulan
Penulis, Tahun	Judul					
Philipp Probst, Marvin N. Wringht, Anne- Laure Boulesteix, (2018)	<i>Hyperparameters and Tuning Strategies for Random Forest</i>	Algoritma <i>Random Forest</i> memiliki beberapa parameter yang dapat mempengaruhi kinerjanya. Parameter yang kurang tepat dapat mempengaruhi model, sehingga perlu dilakukan <i>hyperparameter tuning</i>	Melakukan <i>hyperparameter tuning</i> dan penyetelan strategi untuk mendapatkan akurasi lebih baik pada <i>Random Forest</i> dengan melihat parameter berharga pada data yang fiturnya banyak.	<i>Random Forest</i> dengan melakukan tuning parameter pada <i>mtry</i> , <i>sample size</i> , <i>replacement</i> , <i>node size</i> , <i>number of trees</i> , dan <i>slitting rule</i> .	Data dari Open ML berjudul <i>monks-problem</i> .	<i>Mtry</i> dianggap berpengaruh sangat besar pada <i>tuning parameter</i> penelitian ini. Namun <i>sample size</i> dan <i>node size</i> berpengaruh kecil kepada performa.
Mohammad S. Uddin, Guotai Chi, Mazin A.M. Al Janabi, Tabassun Habib, (2020)	<i>Leveraging Random Forest in Micro- Enterprises Credit Risk Modelling for Accuracy and Interpretability</i>	Data yang digunakan memiliki ukuran besar, berdimensi tinggi dan <i>imbalanced</i> sehingga digunakan <i>Random Forest</i> . Namun jumlah pohon dianggap sangat krusial terhadap kebaikan model, sehingga dilakukan <i>hyperparameter tuning</i> .	Penerapan <i>Random Forest</i> pada pembuatan model klasifikasi risiko kredit pada data yang besar, berdimensi tinggi dan <i>imbalanced</i> .	<i>Random Forest</i> dengan menerapkan <i>hyperparameter tuning</i> pada <i>node</i> dan jumlah pohon.	Data perusahaan mikro di China dengan jumlah data 3111 dengan kondisi tidak seimbang.	Hasil klasifikasi algoritma <i>Random Forest</i> mendapatkan nilai AUC 0.9998 dan akurasi model 99.67%.

Literature Review		Masalah Penelitian/ Rumusan Masalah	Tujuan	Metode	Data	Hasil / Temuan/ Kesimpulan
Penulis, Tahun	Judul					
Unang Sunarya, Tita Haryanti, (2022)	Perbandingan Kinerja Algoritma Optimasi pada Metode <i>Random Forest</i> untuk Deteksi Kegagalan Jantung	Penerapan <i>Random Forest</i> memiliki akurasi yang rendah pada data yang mempunyai banyak fitur (kasus deteksi kegagalan jantung), sehingga diperlukan 3 teknik optimasi untuk dibandingkan hasil akurasi terbaiknya.	Membandingkan 3 algoritma optimasi pada algoritma <i>Random Forest</i> untuk mendapatkan akurasi terbaik.	<i>Grid Search</i> , <i>Random Search</i> , dan <i>Bayesian Search</i> .	Data publik berjudul ' <i>Heart Failure Clinical Records Data</i> ' yang terdiri dari 12 kolom fitur dan 299 baris.	Hasil dari <i>Random Forest</i> menggunakan optimasi <i>Random Search</i> mendapatkan akurasi tertinggi sebesar 85.63%.
Siji George C G, B. Sumathi, (2020)	<i>Grid Search Tuning of Hyperparameters in Random Forest Classifier for Customer Feedback Sentiment Prediction</i>	Mencari akurasi model dari data <i>training</i> dapat menghasilkan nilai <i>overfitting</i> , sehingga perlu dilakukan penyesuaian <i>hyperparameter</i> . <i>Grid Search</i> digunakan untuk mengidentifikasi parameter optimal dari klasifikasi sehingga model dapat secara akurat memprediksi data tidak berlabel.	Melakukan pendekatan <i>hybrid</i> algoritma <i>Random Forest</i> dan <i>Grid Search</i> untuk menghasilkan model yang optimal	<i>Random Forest</i> dengan menerapkan <i>tuning parameter</i> pada <i>n_estimators</i> dan <i>max_features</i> menggunakan <i>Grid Search</i> .	Data <i>sentiment analysis</i> dari UCI database berjumlah 1500 reviews pelanggan (<i>balanced data</i>)	Akurasi naik 5.5% menjadi 90.02% setelah menggunakan <i>proposed method</i> .

Literature Review		Masalah Penelitian/ Rumusan Masalah	Tujuan	Metode	Data	Hasil / Temuan/ Kesimpulan
Penulis, Tahun	Judul					
Phusanisa Charoen-Ung, Pradit Mittrapiyanuruk, (2019)	<i>Sugarcane Yield Grade Prediction Using Random Forest with Forward Feature Selection and Hyperparameter Tuning</i>	Data yang berdimensi tinggi mempunyai kecenderungan <i>overfitting</i> , sehingga perlu dilakukan <i>feature selection</i> . Di sisi lain, akurasi model <i>baseline</i> yang masih rendah membuat diterapkannya <i>hyperparameter tuning</i> .	Melakukan <i>feature selection</i> dan <i>hyperparameter tuning</i> pada model <i>Random Forest</i> untuk mendapatkan akurasi yang terbaik	<i>Random Forest</i> dengan tuning parameter <i>n_estimators</i> , <i>min_sample_leaf</i> , <i>max_features</i> , dan <i>max_depth</i>	Data berjumlah 12520 baris	Akurasi model yang dibuat mendapatkan hasil tertinggi sebesar 71.88%.
Wensu Zhao, Jie Hou, Qili Ran, (2022)	<i>Analysis of Corporate Credit Risk Based on Random Forest and TOPSIS Models</i>	Kekompleksan model dapat menimbulkan <i>overfitting</i> , sehingga penelitian ini menerapkan <i>hyperparameter tuning</i> untuk mendapatkan parameter terbaik dan optimal untuk pembangunan model pada data yang <i>imbalanced</i> .	Mendapatkan akurasi model <i>Random Forest</i> yang paling tepat dan optimal pada <i>data imbalanced</i> dengan <i>GridSearchCV</i> .	Menggunakan model <i>Random Forest</i> dengan <i>GridSearchCV</i> pada <i>max_depth</i> , <i>max number of trees</i> , <i>weight parameter</i> , dan <i>feature evaluation</i> .	Data milik perusahaan di China dengan jumlah 123 baris dengan 8 fitur.	Akurasi model yang didapat 97%

Literature Review		Masalah Penelitian/ Rumusan Masalah	Tujuan	Metode	Data	Hasil / Temuan/ Kesimpulan
Penulis, Tahun	Judul					
Anna Baita, Inggar Adi prasetyo, Nuri Cahyono (2023)	<i>Hyperparameter Tuning on Random Forest for Diagnose Covid-19</i>	Adanya ketidakseimbangan data yang mempengaruhi proses pelatihan dan menyebabkan bias yang tidak menguntungkan terhadap kelas mayoritas. <i>Random Forest</i> sendiri merupakan metode yang sensitif terhadap nilai parameternya. Sehingga diperlukan <i>balancing data</i> dan <i>hyperparameter tuning</i> .	Membandingkan penggunaan <i>oversampling</i> dan <i>hyperparameter tuning</i> pada model <i>Random Forest</i> khususnya pada kasus Diagnosa penyakit Covid-19.	<i>Oversampling</i> dan <i>hyperparameter tuning</i> model <i>Random Forest</i> .	Data penyakit Covid-19 dengan 20 fitur.	Penggunaan <i>oversampling</i> dan <i>hyperparameter tuning</i> khususnya <i>Grid Search</i> dapat meningkatkan akurasi model menjadi 96%.
Riska Aryanti, Titik Misriati, Rahmat Hidayat (2023)	Klasifikasi Risiko Kesehatan Ibu Hamil Menggunakan <i>Random Oversampling</i> Untuk Mengatasi Keseimbangan Data	Masalah ketidakseimbangan data yang terjadi pada klasifikasi risiko kesehatan ibu hamil yang menyebabkan model klasifikasi yang dibuat cenderung akan memprediksi kelas mayoritas saja.	Menerapkan <i>Random Oversampling</i> untuk mengatasi <i>imbalanced</i> pada data untuk mendapatkan nilai akurasi yang lebih baik.	Menggunakan teknik <i>Random Oversampling</i> untuk mengatasi <i>imbalanced data</i> model <i>Random Forest</i> untuk klasifikasi risiko kehamilan.	Data berasal dari UCI <i>Machine learning</i> yang terdiri dari 1014 baris dengan 6 fitur.	Kenaikan 1.09% pada akurasi model RF pada data hasil <i>Random Oversampling</i>

2.2. Landasan Teori

2.2.1. *Credit Risk* / Risiko Kredit

Surat Edaran Bank Indonesia No.13/24/DPNP/2011 menyatakan bahwa risiko kredit adalah risiko akibat kegagalan debitur dan atau pihak lain dalam memenuhi kewajiban kepada Bank. Risiko kredit terjadi karena ketidakmampuan nasabah atau peminjam dalam melunasi kredit di waktu yang ditentukan pada perjanjian awal [20]. Adanya keterlambatan pembayaran kredit pada suatu lembaga keuangan membuat nilai *Non-Performing Loan* (NPL) meningkat. Dari pernyataan tersebut dapat dikatakan bahwa semakin tinggi nilai *Non-Performing Loan* maka kualitas kredit di lembaga keuangan tersebut bisa disebut buruk.

Bank Indonesia menentukan batas maksimal nilai NPL adalah 5%. Jika nilai NPL di bawah 5% artinya lembaga keuangan dinyatakan sehat dalam peminjaman kreditnya. Namun jika NPL di atas 5% artinya lembaga keuangan tersebut mengalami kerugian atau likuiditas [21]. Perhitungan nilai *Non-Performing Loan* sesuai dengan persamaan 2.1.

$$NPL = \frac{\text{Kredit Bermasalah}}{\text{Total Kredit}} \quad (2.1)$$

Keterangan:

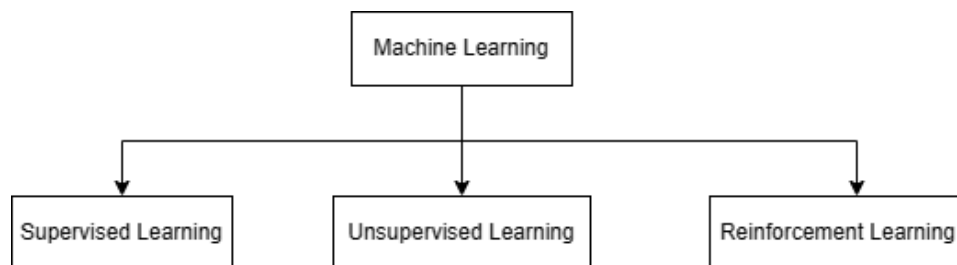
1. NPL (*Non-Performing Loan*) mencakup kredit yang mengalami tunggakan pembayaran atau tidak dapat dipenuhi oleh peminjam sesuai perjanjian.
2. Kredit bermasalah merupakan jumlah total kredit yang dianggap bermasalah atau tidak lancar.
3. Total kredit merupakan total keseluruhan kredit yang diberikan oleh lembaga keuangan kepada peminjamnya.

Untuk mengajukan kredit sendiri terdapat beberapa persyaratan yang harus diinputkan oleh kreditur kepada lembaga keuangan seperti identitas peminjam, penghasilan, jumlah pinjaman, dan lain sebagainya. Dari data tersebut risiko kredit bisa diklasifikasikan dan diprediksi.

2.2.2. *Machine Learning*

Machine learning atau biasa disebut dengan pembelajaran mesin adalah salah satu bagian dari *artificial intelligence* (AI) yang berfokus pada pengembangan

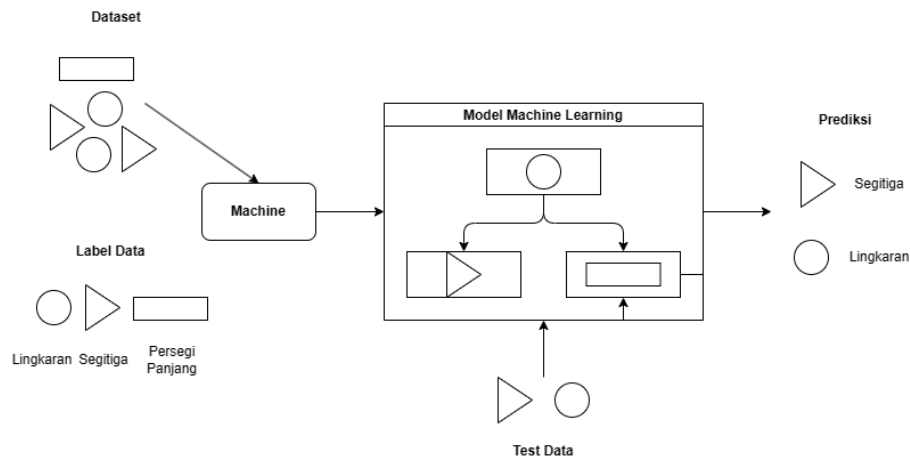
sebuah sistem untuk mampu belajar sendiri layaknya manusia tanpa harus diprogram berulang kali [22]. Menurut Mitchel [23], bisa dikatakan *machine learning* jika satu program komputer telah melakukan pembelajaran dari sebuah pengalaman (*experience*) terhadap tugas (*task*) dan mengukur hasil peningkatan kinerja (*performance measure*). Fokus utama dari *machine learning* adalah membangun aplikasi komputer yang dapat mempelajari data secara mandiri dan mampu membuat model yang siap digunakan untuk memecahkan sebuah kasus atau masalah [24]. *Machine learning* terdiri dari tiga sub bidang yang terlihat pada Gambar 2.1, yaitu *supervised learning*, *unsupervised learning*, dan *reinforcement learning*.



Gambar 2.1 Sub Pembagian *Machine learning*

Supervised learning merupakan sub bagian dari *machine learning* yang digunakan untuk mengolah kumpulan data berlabel [25]. Contoh penerapan dari *supervised learning* adalah klasifikasi dan regresi.

Klasifikasi merupakan proses pengelompokan objek yang mempunyai kesamaan karakteristik atau ciri kedalam beberapa kelas [3]. Proses klasifikasi dengan menerapkan *machine learning* dimulai dari inputan data berlabel yang akan diberikan pembelajaran menggunakan mesin dan mesin akan mendeteksi data – data baru (data uji / *testing*) ke dalam label atau kelas tertentu. Contoh pengaplikasian klasifikasi menggunakan *machine learning* terdapat pada Gambar 2.2.



Gambar 2.2 Alur Klasifikasi Menggunakan *Machine learning*

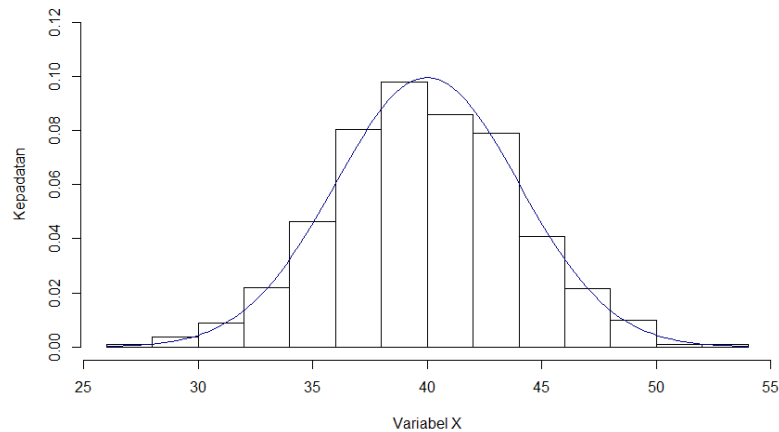
Unsupervised learning merupakan sub bagian pada *machine learning* yang digunakan untuk mengolah data tidak berlabel [24]. Contoh penerapannya pada kasus *clustering*. Sub bagian yang terakhir yaitu *reinforcement learning* yang pada prosesnya tanpa menggunakan inputan data, pada pembelajaran ini terdapat *agent* yang dapat mengamati lingkungan, memilih dan melakukan tindakan dan mendapatkan hasil untuk proses pembelajaran dan tindakan selanjutnya [24].

2.2.3. *Exploratory Data Analysis*

Exploratory Data Analysis (EDA) adalah salah satu tahapan pada analisis penelitian domain *Data Science*. *Exploratory data analysis* merupakan pendekatan analisis terhadap suatu data yang berguna untuk membuat dan mengetahui *summary* data. Tujuan utama *exploratory data analysis* adalah untuk memeriksa distribusi, *outliers*, dan tren data yang mengarah ke pengujian spesifik hipotesis yang dibuat [26]. untuk mengungkap wawasan yang dapat membantu memahami data dengan lebih baik. *Exploratory data analysis* berperan penting dalam mengidentifikasi pertanyaan yang relevan, merumuskan hipotesis, dan mempersiapkan data sebelum diterapkan pada model analisis yang lebih lanjut.

2.2.3.1. **Distribusi Data**

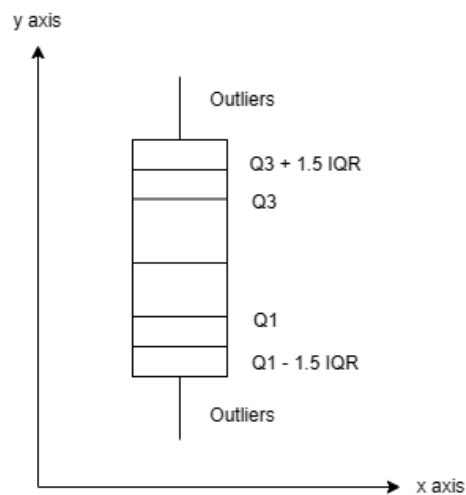
Terdapat bermacam – macam distribusi data pada statistika, seperti distribusi data normal dan tidak normal. Data dikatakan berdistribusi normal jika datanya simetris di sekitar nilai tengahnya. Gambar 2.3 memperlihatkan bahwa data berdistribusi normal.



Gambar 2.3 Contoh Data Berdistribusi Normal

2.2.3.2. Outliers Data

Sebuah data yang memiliki sifat dan karakteristik yang berbeda dari data pada umumnya dan kemunculan kejadian yang relatif sedikit dikatakan *outliers* [27]. *Outliers* bisa jauh lebih besar atau jauh lebih kecil daripada nilai-nilai lainnya, sehingga dapat menyebabkan distorsi pada analisis dan kesimpulan yang diambil dari data. Untuk melihat nilai *outliers*, biasanya menggunakan grafik boxplot.



Gambar 2.4 Data *Outliers*

Gambar 2.4 menunjukkan data yang terdeteksi *outliers* yaitu jika terletak di atas nilai *upper bound* dan di bawah nilai *lower bound*. Perhitungan nilai *upper bound* dan *lower bound* terdapat pada persamaan 2.2 dan 2.3.

$$\text{Lower Bound} = Q1 - (1.5 \times \text{IQR}) \quad (2.2)$$

$$\text{Upper Bound} = Q3 + (1.5 \times \text{IQR}) \quad (2.3)$$

Keterangan:

Q1 = Percentile ke-25

Q3 = Percentile ke-75

IQR = Q3 – Q1

2.2.3.2. Korelasi Fitur

Korelasi fitur adalah suatu ukuran statistik yang mengukur sejauh mana dua variabel atau fitur dalam dataset berhubungan satu sama lain. Terdapat beberapa jenis korelasi seperti korelasi *pearson*, *spearman*, *kendall*, dan lainnya. Korelasi *pearson* digunakan untuk digunakan untuk mengukur sejauh mana hubungan linier antara dua variabel kontinu. Perhitungan korelasi *pearson* pada persamaan 2.4.

$$r = \frac{n\sum xy - (\sum x)(\sum y)}{\sqrt{\{n\sum x^2 - (\sum x)^2\} \{n\sum y^2 - (\sum y)^2\}}} \quad (2.4)$$

Keterangan:

N = banyaknya pasangan x dan y

$\sum x$ = Total jumlah dari variabel x

$\sum y$ = Total Jumlah dari Variabel y

$\sum x^2$ = Kuadrat dari total jumlah variabel x

$\sum y^2$ = Kuadrat dari total jumlah variabel y

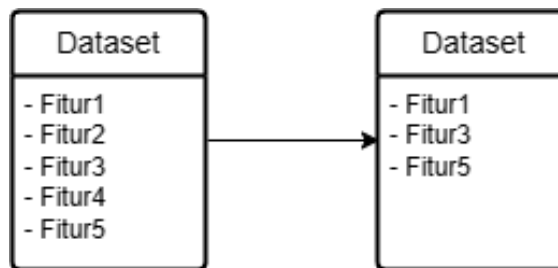
$\sum xy$ = Hasil perkalian dari total jumlah variabel x dan variabel y.

2.2.4. Feature Engineering

Feature atau fitur merupakan representasi numerik dari data mentah [20]. *Feature engineering* adalah proses mengubah, memodifikasi, atau membuat fitur baru dari data mentah (*raw data*) untuk meningkatkan kualitas dalam pemodelan. Tujuan utamanya adalah untuk meningkatkan performa dan akurasi model dengan membersihkan data dan memilih fitur data yang mempunyai pengaruh besar. *Feature engineering* merupakan langkah penting dalam membangun model *machine learning* karena fitur yang tepat dapat meringankan kesulitan pemodelan, sehingga menghasilkan akurasi yang lebih baik dan berkualitas [28]. Terdapat beberapa teknik yang termasuk *feature engineering* untuk data tabular, seperti *feature selection*, *standardization*, dan *oversampling*.

2.2.4.1. Feature Selection

Feature selection merupakan teknik menyeleksi fitur yang akan digunakan dalam membuat model *machine learning*. Contoh penerapan *feature selection* terdapat pada Gambar 2.5.

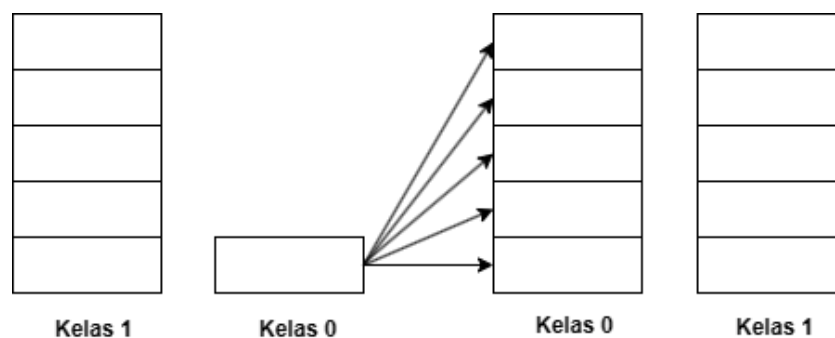


Gambar 2.5 Feature Selection

Untuk mendapatkan fitur – fitur yang relevan pada model *Random Forest*, dilihat dari nilai *feature importance*. *Feature importance* yaitu proses untuk melihat sejauh mana fitur-fitur berkontribusi terhadap pemodelan. *Random Forest* merupakan salah satu algoritma yang dapat dicari nilai *feature importancenya*. Semakin tinggi nilai *feature importance*, semakin besar kontribusinya fitur terhadap model. Dalam prediksi risiko kredit sendiri, mengabaikan fitur yang tidak relevan dapat meningkatkan akurasi model, mengurangi durasi waktu *training* atau pelatihan, dan menurunkan biaya komputasi saat menjalankan beberapa model *machine learning*.

2.2.4.2. Oversampling

Terdapat beberapa teknik *sampling* seperti *undersampling* dan *oversampling*. *Oversampling* merupakan teknik menyeimbangkan data dengan memperbanyak data pada kelas minoritas [31]. Contoh *oversampling* terlihat pada Gambar 2.6.



Gambar 2.6 Oversampling

Terdapat beberapa teknik *oversampling* seperti *Random Oversampling*, SMOTE, BorderLine SMOTE, KMeans SMOTE, SVM SMOTE, *Adasyn*, dan SMOTE-NC. SMOTE adalah singkatan dari *Synthetic Minority Over-sampling Technique*. SMOTE merupakan teknik *oversampling* yang menghasilkan data sintetik dengan menghasilkan *instance* baru di antara garis yang menghubungkan dua fitur [32]. *Random oversampling* merupakan teknik *sampling* dengan menambahkan jumlah data secara acak pada kelas minoritas [33]. Contoh penerapan *random oversampling* terdapat pada Gambar 2.7.

Index	Fitur A	Fitur B	Fitur C	Kelas
0	25	50	20	0
1	20	45	15	0
2	50	100	40	1
3	45	115	60	1
4	55	110	55	1
5	60	95	70	1
6	25	50	20	0
7	20	45	15	0

Gambar 2.7 Contoh *Random Ovesampling*

Gambar 2.7 menunjukkan bahwa index ke 0 dan 1 diduplikat secara random ke index 6 dan 7 sehingga jumlah kelas minoritas (0) menjadi sama dengan kelas mayoritas (1).

2.2.4.3. Standardisasi

Standardisasi atau *standardization* merupakan metode yang dilakukan untuk membuat beberapa variabel data memiliki rentang nilai yang sama, tidak ada yang terlalu besar maupun terlalu kecil [29]. Salah satu teknik standardisasi yaitu *Robust Scaler*, merupakan salah satu teknik optimasi untuk mentransformasikan suatu nilai pada data dengan menggunakan *median* dan *quartiles* [30]. Standardisasi menggunakan *Robust Scaler* dijelaskan pada persamaan

$$x' = \frac{x - \text{median}(x)}{Q3 - Q1} \quad (2.5)$$

Notasi x' merupakan hasil standardisasi menggunakan *Robust Scaler*, $Q3 - Q1$ merupakan hasil *Interquartile Range*. Menggunakan persamaan 2.5, maka contoh perhitungan *Robust Scaler* terdapat pada Tabel 2.2.

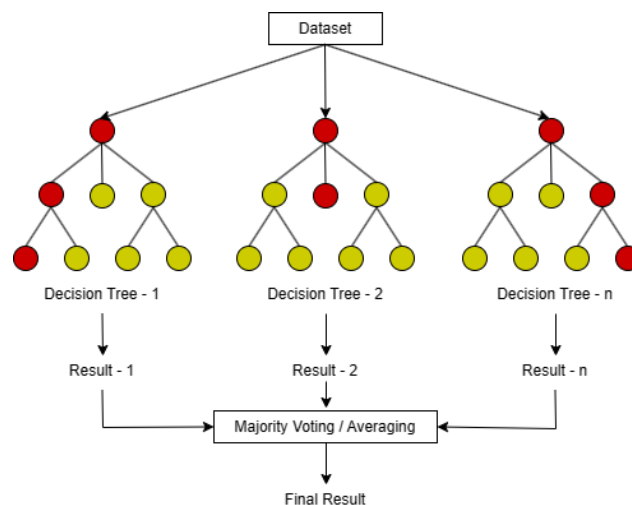
Tabel 2.2 Contoh Hasil Standardisasi Menggunakan *Robust Scaler*

Sebelum standardisasi		Setelah standardisasi	
Fitur A	Fitur B	Fitur A	Fitur B
12000	25	-1.14285714	-1.38461538
15000	50	-0.28571429	0.15384615
20000	75	1.14285714	1.69230769
17000	45	0.28571429	-0.15384615

Dari Tabel 2.2 menunjukkan bahwa sebelum dilakukan standardisasi, nilai pada kolom Fitur A (12000) dan Fitur B (25) memiliki rentang skala nilai yang berbeda. Setelah menjalankan proses standardisasi, Fitur A dan Fitur B akan memiliki rentang skala data yang serupa.

2.2.5. *Random Forest*

Random Forest adalah salah satu algoritma *machine learning* yang bekerja dengan cara mengkombinasikan beberapa algoritma *Decision Tree* [15]. *Random Forest* termasuk algoritma *ensemble* yang artinya algoritma ini membangun model di atas banyak model sehingga mendorong model tersebut mendapatkan hasil lebih baik [24]. Dengan menggunakan banyak *Decision Tree*, algoritma ini dapat memperoleh prediksi yang akurat dan meminimalkan *overfitting*. Selain kuat terhadap *overfitting*, *Random Forest* juga memiliki kinerja yang lebih baik jika dibandingkan dengan beberapa *machine learning* lain seperti *Support Vector Machine* (SVM) dan *Discriminant Analysis*. Gambar 2.8 menjelaskan alur kerja algoritma *Random Forest*.

Gambar 2.8 Cara Kerja Algoritma *Random Forest*

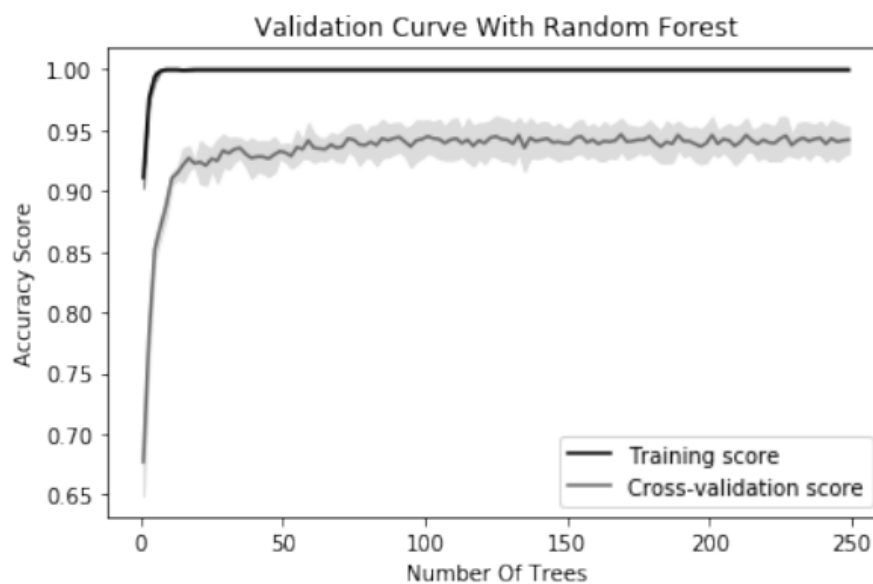
Cara kerja algoritma *Random Forest* dimulai dari pemilihan acak sampel data yang digunakan, kemudian membuat *Decision Tree* untuk setiap sampel data. Dilanjutkan dengan *voting* setiap hasil prediksi menggunakan nilai modus (pada model klasifikasi), sehingga hasil akhirnya atau hasil prediksi didapatkan dari hasil *voting* terbaik.

2.2.6. Hyperparameter Tuning

2.2.6.1. Parameter *Random Forest*

Hyperparameter Tuning adalah tugas menemukan parameter paling optimal untuk algoritma *machine learning* yang akan diterapkan pada data. Pada *machine learning*, optimasi menggunakan *hyperparameter tuning* sangat diperlukan untuk mendapatkan ukuran kinerja (*Accuracy* atau *AUC*) yang terbaik. Terdapat beberapa parameter yang bisa diubah dan diterapkan pada algoritma *Random Forest*, seperti jumlah pohon (*n_estimators*), kriteria (*criterion*), maksimal kedalaman pohon (*max_depth*), jumlah fitur yang dipertimbangkan (*max_features*), jumlah sampel minimum untuk pembentukan *node* baru (*min_samples_split*), dan jumlah sampel minimum pada daun (*min_samples_leaf*).

N_estimators atau jumlah pohon pada *Random Forest* dapat ditentukan untuk mendapatkan akurasi model paling optimal.



Gambar 2.9 Percobaan Jumlah *n_estimators*

Gambar 2.9 menunjukkan percobaan *tuning parameter* pada *n_estimators*

atau jumlah pohon mulai dari 0 sampai 250 yang menghasilkan nilai akurasi berbeda – beda. Nilai akurasi mengalami kenaikan dan penurunan di setiap jumlah pohon yang berbeda sehingga jumlah pohon keputusan pada model ini harus diperhatikan.

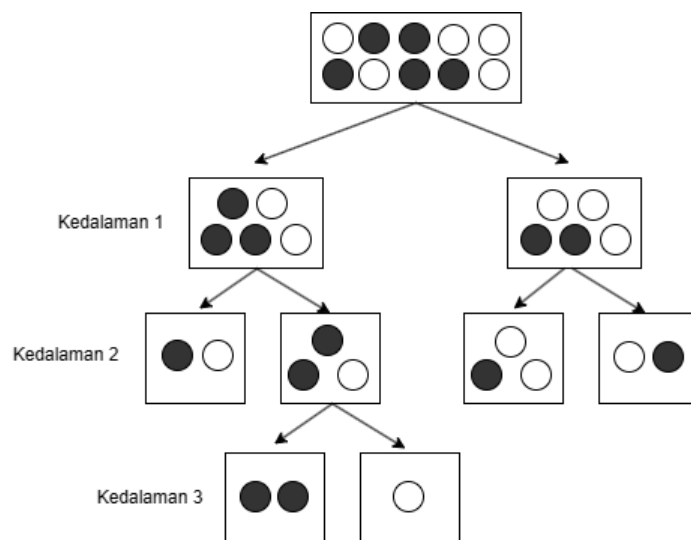
Criterion atau kriteria pembuatan model *Random Forest* terbagi menjadi dua yaitu *Gini* dan *Entropy*. Perhitungan menggunakan kriteria *gini* lebih cepat jika dibandingkan dengan *entropy* dikarenakan tidak menggunakan perhitungan logaritma. Rumus perhitungan *gini* dan *entropy* sesuai persamaan 2.2 dan 2.3

$$Gini = 1 - \sum_{i=1}^n (p_i)^2 \quad (2.2)$$

$$Entropy = \sum_{i=1}^c - p_i \log_2 p_i \quad (2.3)$$

Notasi p menunjukkan probabilitas *entropi* dan p_i merupakan probabilitas suatu objek dikategorikan ke kelas tertentu.

Max_depth atau kedalaman pohon dapat ditentukan berdasarkan data yang akan digunakan. Gambar 2.10 menunjukkan kedalaman pohon sebuah model *Random Forest* yang ditentukan.



Gambar 2.10 Kedalaman pohon

Gambar 2.10 menggambarkan kedalaman pohon yang diatur yaitu 3 sehingga terdapat 3 kedalaman pohon yang terbentuk. Jumlah kedalaman pohon atau *max_depth* dapat diatur sesuai kebutuhan.

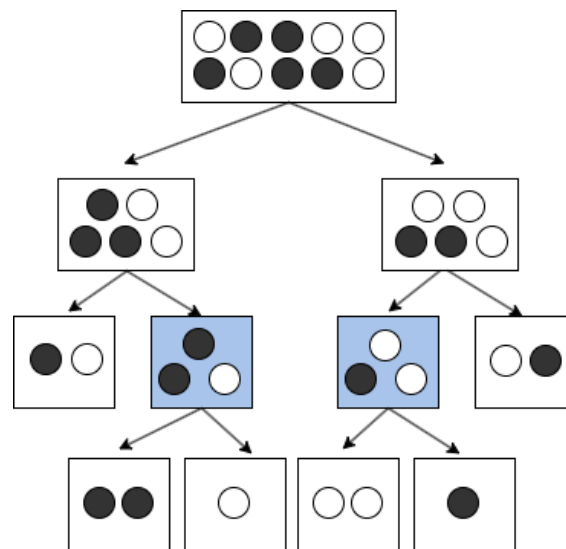
Selain itu terdapat *features* pada *Random Forest* yang bisa ditentukan jumlah maksimalnya atau bisa disebut *max_features*. Jumlah *max_features* yang digunakan pada umumnya yaitu berupa nilai *sqrt*, *log2*, dan *none* yang berupa nilai integer atau *float*. Secara *default* yang digunakan yaitu *sqrt* yang terdapat pada persamaan 2.4, 2.5, dan 2.6.

$$\max_{features}(sqrt) = sqrt(n_{features}) \quad (2.4)$$

$$\max_{features}(log2) = log2(n_{features}) \quad (2.5)$$

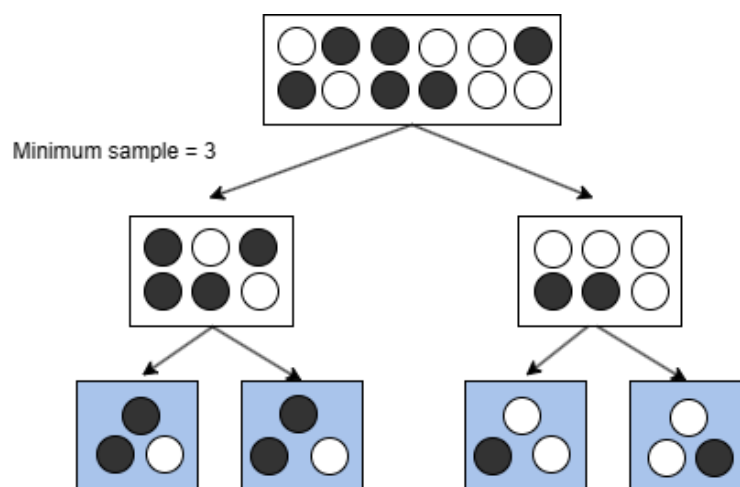
$$\max_{features}(none) = n_{features} \quad (2.6)$$

$n_{features}$ pada persamaan 2.4, 2.5, dan 2.6 merupakan jumlah fitur pada data. Parameter yang bisa ditentukan selanjutnya yaitu *min_sample_split* atau jumlah minimal *split* pada *node* yang ditunjukkan pada Gambar 2.11.



Gambar 2. 11 Contoh *Min_Samples_Split*

Gambar 2.11 menunjukkan bahwa jumlah *min_sample_split* yang ditetapkan yaitu 3 (kotak berwarna biru). Jika jumlah sampel di bawah 3 maka tidak dapat dilakukan *split*. Parameter lain yang bisa ditentukan pada model *Random Forest* yaitu *min_samples_leaf* atau jumlah minimal *node* pada daun.



Gambar 2.12 Contoh *Min_Samples_Leaf*

Pada Gambar 2.12 menunjukkan bahwa pada jumlah sampel pada daun telah memenuhi kondisi minimal yaitu 3 sampel (kotak berwarna biru). Dapat dikatakan bahwa di suatu daun tidak boleh terdapat sampel di bawah 3.

2.2.6.2. *Grid Search*

Untuk memudahkan dalam penggunaan *hyperparameter tuning*, terdapat algoritma optimasi seperti *Grid Search CV*, *Randomized Search CV*, dan *Bayesian CV* [11]. Algoritma *Grid Search* tepat untuk diterapkan saat ingin melakukan pencarian parameter secara eksplisit dan memeriksa semua kombinasi secara sistematis. Dibandingkan *Random Search*, *Grid Search* dianggap lebih sistematis karena mencoba semua kombinasi, sehingga cocok untuk ruang pencarian yang kecil dan terdefinisi dengan baik.

2.2.6.3. *Random Search*

Random Search tepat digunakan jika ingin mencari solusi dengan cepat karena algoritma ini secara acak mengambil parameter yang akan digunakan. Algoritma *Random Search* mempunyai kelebihan pada proses *training* dikarenakan secara acak memilih sekumpulan kombinasi parameter untuk dievaluasi. Namun kekurangan dari *Random Search* adalah kurang menjelajahi semua parameter dikarenakan hanya diambil beberapa parameter secara acak saja.

2.2.7. Matriks Konfusi

Terdapat beberapa cara pengukuran keakuratan model *machine learning*.

Beberapa metrik pengukuran yang digunakan seperti *accuracy*, *confusion matrix*, dan *classification report*. *Accuracy* merupakan sebuah matrik evaluasi yang digunakan untuk mengukur keakuratan model dalam melakukan prediksi atau klasifikasi yang secara perhitungan ditunjukkan persamaan (2.7).

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (2.7)$$

Keterangan:

- True Positive* (TP) merupakan jumlah observasi yang benar – benar termasuk ke dalam kelas positif dan diprediksi benar oleh model.
- True Negative* (TN) merupakan jumlah observasi yang benar – benar termasuk ke dalam kelas negatif dan diprediksi dengan benar oleh model.
- False Positive* (FP) merupakan jumlah observasi yang diprediksi sebagai bagian dari kelas positif oleh model, tetapi sebenarnya observasi tersebut termasuk dalam kelas negatif.
- False Negative* (FN) merupakan jumlah observasi yang diprediksi sebagai bagian dari kelas negatif oleh model, tetapi sebenarnya observasi tersebut termasuk dalam kelas positif.

Untuk membuktikan nilai akurasi model yang baik, dapat dilihat dari nilai *confusion matrix* seperti Gambar 2.13.

		Predicted Class	
		Normal	Attack
Actual Class	Normal	True Negative (TN)	False Positive (FP)
	Attack	False Negative (FN)	True Positive (TP)

Gambar 2.13 *Confusion Matrix*

Model dikatakan baik jika nilai TN dan TP > FN dan FP. Menggunakan tabel *confusion matrix* ini, akan terlihat jumlah hasil prediksi benar dan salah sehingga dapat membuktikan secara akurat kebaikan model.

2.2.8. Classification Report.

Classification Report merupakan salah satu metrik evaluasi kinerja model pada *machine learning* berbasis klasifikasi. Pada hasil *classification report* terdapat empat nilai yaitu *precision*, *recall*, *f1-score*, dan *support*.

- a. *Precision*, mengukur sejauh mana model memberikan hasil yang benar secara proporsional terhadap hasil yang dinyatakan sebagai positif, ditunjukkan persamaan (2.8).

$$Precision = \frac{TP}{TP + FP} \quad (2.8)$$

- b. *Recall*, mengukur sejauh mana model dapat mengidentifikasi secara benar semua kasus positif yang ada, seperti pada persamaan 2.9.

$$Recall = \frac{TP}{TP + FN} \quad (2.9)$$

- c. *F1-Score*, memberikan pengukuran yang seimbang antara *precision* dan *recall*. *F1-score* berguna ketika ingin memperhatikan *precision* dan *recall* secara bersamaan, serta mencari nilai rata-rata yang optimal dari keduanya. Perhitungan *F1-Score* terdapat pada persamaan 2.10.

$$F1 - Score = \frac{Precision \times Recall}{Precision + Recall} \quad (2.10)$$

- d. *Support* merupakan keterangan informasi tentang banyaknya data yang ada untuk setiap kelas dalam data yang dievaluasi.

2.2.9. Kurva AUC (Area Under the Curve)

Kurva AUC (*Area Under the Curve*) adalah alat evaluasi untuk mengukur kinerja model klasifikasi. Kurva AUC menggambarkan hubungan antara tingkat *true positive (sensitivity)* dan tingkat *false positive (1-specificity)* pada berbagai *threshold*. Nilai AUC berkisar antara 0 hingga 1, nilai 1 menunjukkan performa sempurna dan nilai 0,5 menunjukkan performa model setara dengan perbandingan acak. Semakin tinggi nilai AUC, semakin baik model dalam membedakan antara kelas positif dan negatif. Analisis kurva AUC untuk memperoleh wawasan tentang kinerja model dan memilih *threshold* yang tepat untuk pengambilan keputusan berdasarkan *trade-off* antara *true positive* dan *false positive*.