

BAB II

TINJAUAN PUSTAKA

2.1. Kajian Pustaka

Penelitian [10] membahas penggunaan metode SMOTE dengan *Random Forest* dan *Xgboost*. Fokus penelitian ini adalah penerapan metode *SMOTE* pada algoritma *Random Forest* dan *XGboost* dalam mengklasifikasikan tingkat penyakit Hepatitis C pada data dengan distribusi kelas yang tidak seimbang. Temuan dari penelitian ini menunjukkan bahwa kedua metode, yaitu *SMOTE Random Forest* dan *SMOTE XGboost* berhasil mencapai akurasi dan nilai *recall* yang melebihi 75%. Meskipun akurasi dari *Random Forest* dan *XGboost* mencapai sekitar 74%, nilai *recall* yang didapat tetap di bawah 2%.

Penelitian [11] meneliti mengenai pendekatan alternatif untuk memprediksi klasifikasi status vaksinasi Hepatitis-B menggunakan pendekatan *machine learning* seperti *Random Forest* dan *Naive Bayes*. Namun, untuk klasifikasi pada data yang tidak seimbang, algoritma cenderung bias terhadap kelas mayoritas. Untuk meningkatkan prediksi yang akurat, penelitian ini menggunakan *SMOTE*. Hasil penelitian menunjukkan bahwa penggunaan *SMOTE* pada *Random Forest* dan *Naive Bayes* meningkatkan akurasi identifikasi status non-vaksinasi Hepatitis-B sebesar 30,08% dan 26,09% secara berturut-turut, dibandingkan dengan *non-SMOTE*. *Random Forest* dengan *SMOTE* menjadi model terbaik untuk klasifikasi status vaksinasi Hepatitis-B.

Penelitian lain yang menerapkan *SMOTE* juga dilakukan pada penelitian [12] untuk mengatasi masalah ketidakseimbangan kelas. Penelitian ini menerapkan *Weighted Categorical Loss (WCL)* dan *Synthetic Minority Oversampling Technique (SMOTE)*. Enam model *CNN state-of-the-art* dilatih dengan *Transfer Learning (TL)* pada tiga data COVID-19 yang berbeda, dengan fokus pada klasifikasi multi kelas antara COVID-19, normal, dan kasus pneumonia virus. Hasil penelitian ini menunjukkan *SMOTE* secara signifikan meningkatkan nilai *AUC* pada model *DenseNet201* dan *VGG19* dibandingkan dengan penggunaan data yang tidak seimbang.

Penelitian [13] berfokus pada meningkatkan prediksi terhadap pasien penyakit gagal jantung kronis dengan menggunakan teknik *SMOTE* dan pendekatan *machine learning*. Delapan model klasifikasi digunakan dalam penelitian ini, yaitu *Random Forest (RF)*, *Extra Tree (ET)*, *Naive Bayes (NB)*, *k-Nearest Neighbor (KNN)*, *Decision Tree J48*, *Decision Table/Naive Bayes hybrid classifier (DTNB)*, *Optimized Forest*, dan *Alternating Decision Tree (ADTree)*, untuk memprediksi kelangsungan hidup pasien. Hasil percobaan menunjukkan teknik *SMOTE* mengubah akurasi keluaran dari metode pengklasifikasi yang dipilih dan *DTNB* mencapai akurasi tertinggi dengan 87,08% menggunakan *SMOTE* untuk memperkirakan kelangsungan hidup pasien jantung. Semua percobaan dilakukan dalam lingkungan simulasi menggunakan alat WEKA. dengan menggunakan *SMOTE* untuk memprediksi kelangsungan hidup pasien jantung.

Pada penelitian [14] dilakukan penelitian mengenai perbandingan dua algoritma klasifikasi yaitu *C4.5* dan *k-Nearest Neighbor (KNN)*, dan metode pra-pemrosesan *SMOTE* diterapkan dalam klasifikasi kinerja akademik mahasiswa. Hasil eksperimen menggunakan aplikasi Rapid Miner menunjukkan bahwa metode *C4.5 Decision Tree* menghasilkan performa prediksi yang lebih baik dalam hal akurasi, *recall*, dan presisi masing-masing sebesar 71,09%, 71,63%, dan 71,54% dibandingkan dengan algoritma *K-Nearest Neighbor*.

Penelitian [15] melakukan penelitian tentang penanganan pada data yang tidak seimbang pada penelitian ini menggunakan metode *oversampling* dengan teknik *SMOTE*, *Borderline-SMOTE*, dan *ADASYN*, hasil dari penelitian ini menunjukkan bahwa penggunaan teknik *oversampling* meningkatkan akurasi dari 2% hingga 11% pada beberapa data.

Penelitian [16] berfokus pada memprediksi kegagalan mesin di industri melalui algoritma pembelajaran mesin. Namun, tantangan ketidakseimbangan data pada setiap kelas menghambat kinerja model. Solusi yang diterapkan yaitu menggunakan teknik *oversampling* menggunakan *SMOTE*. Teknik *oversampling* tersebut menghasilkan sampel sintetis pada kelas minoritas. Hasil penelitian menunjukkan bahwa penggunaan *SMOTE* meningkatkan kinerja *Random Forest*

dalam mengenali kondisi kegagalan dan *non*-kegagalan mesin dengan peningkatan skor *AUC* 7,83%.

Penelitian [17] meneliti mengenai penanganan masalah data yang tidak seimbang pada data penipuan, hasil dari penelitian ini adalah teknik *oversampling* terbukti dapat meningkatkan kinerja model, dan teknik *oversampling* terbaik adalah menggunakan *SMOTE*.

Penelitian [18] mengadopsi teknik sampling *SMOTE* dan *Random Oversampling* untuk mengatasi data tidak seimbang. Teknik sampling *SMOTE* melibatkan mensintesis data baru yang mirip dengan data minoritas, sedangkan teknik sampling *ROS* melibatkan duplikasi data minoritas. Hasil penelitian menunjukkan bahwa model klasifikasi *SMOTE* dengan *Naive Bayes* memiliki akurasi yang lebih tinggi 0,61% daripada model klasifikasi *ROS* dengan *Naive Bayes*.

Penelitian [19] meneliti mengenai penggunaan metode resampling *oversampling SMOTE* pada dataset yang tidak seimbang menunjukkan bahwa model yang dilatih pada data yang *dioversampling* dapat memprediksi lebih, daripada model yang dilatih pada data yang tidak *dioversampling*. Ini menunjukkan bahwa *oversampling* dengan *SMOTE* dapat mengubah kinerja *XAI* menjadi lebih baik.

Tabel 2.1. Penelitian Terdahulu

No	Penulis	Judul Penelitian	Metode Penelitian	Masalah Penelitian	Hasil Penelitian
1	Muhamad Syukron, Rukun Santoso, Tatik Widiharih (2020)	Perbandingan Metode <i>SMOTE Random Forest</i> dan <i>SMOTE XGBoost</i> untuk Klasifikasi Tingkat Penyakit Hepatitis C pada <i>Imbalanced Class Data</i>	Menggunakan metode <i>SMOTE Random Forest</i> dan <i>SMOTE XGboost</i> dalam klasifikasi tingkat penyakit Hepatitis C pada data tidak seimbang.	Data yang digunakan dalam penelitian ini memiliki ketidakseimbangan yang dapat mempengaruhi akurasi prediksi. Oleh karena itu, diperlukan pendekatan untuk meningkatkan akurasi prediksi, salah satunya adalah menggunakan algoritma <i>SMOTE</i> .	Nilai akurasi dan <i>recall</i> dari kedua metode lebih dari 75%. Namun, akurasi <i>Random Forest</i> dan <i>XGBoost</i> sekitar 74% dengan nilai <i>recall</i> kurang dari 2%. Variabel hasil tes laboratorium memiliki pengaruh besar dalam menentukan tingkat penyakit Hepatitis C.
2	V M Putri, M Masjkur, C Suhaeni (2020)	<i>Performance of SMOTE in a Random Forest and Naive Bayes Classifier for Imbalanced Hepatitis-B Vaccination Status</i>	Menggunakan <i>SMOTE</i> dalam <i>Random Forest</i> dan <i>Naive Bayes Classifier</i> untuk klasifikasi status vaksinasi Hepatitis-B.	Masalah dalam penelitian ini adalah terkait dengan klasifikasi status vaksinasi Hepatitis-B yang tidak seimbang. Kondisi tidak seimbang ini mengacu pada perbedaan yang signifikan dalam jumlah sampel antara kelas mayoritas dan kelas minoritas. Dalam konteks penelitian ini, kelas mayoritas mungkin adalah status vaksinasi Hepatitis-B yang sudah divaksinasi, sedangkan kelas minoritas adalah status vaksinasi Hepatitis-B yang belum divaksinasi.	Penggunaan <i>SMOTE</i> meningkatkan akurasi identifikasi status <i>non-vaksinasi</i> Hepatitis-B sebesar 30,08% dan 26,09% pada <i>Random Forest</i> dan <i>Naive Bayes</i> . Model <i>Random Forest</i> dengan <i>SMOTE</i> menjadi yang terbaik.

No	Penulis	Judul Penelitian	Metode Penelitian	Masalah Penelitian	Hasil Penelitian
3	Ekram Chamseddine, Nesrine Mansouri, Makram Soui, Mourad Abed (2022)	<i>Handling class imbalance in COVID-19 chest X-ray images classification: Using SMOTE and weighted loss</i>	Penelitian ini menggunakan teknik <i>deep learning</i> , khususnya <i>Convolutional Neural Networks (CNNs)</i> , untuk mengklasifikasikan Gambar <i>X-ray</i> dada pasien COVID-19, dan untuk mengatasi masalah ketidakseimbangan kelas, penelitian ini menerapkan <i>Weighted Categorical Loss (WCL)</i> dan <i>SMOTE</i> .	Prediksi <i>churn</i> pelanggan bank dalam konteks data yang tidak seimbang. Dalam era <i>big data</i> , dengan perkembangan pesat keuangan internet dan meningkatnya tekanan pada bank untuk bersaing, bank beralih fokus pada mempertahankan pelanggan lama, sehingga penting untuk memprediksi <i>churn</i> pelanggan. Dalam penelitian ini, penulis merancang model prediksi <i>Borderline-SMOTE-Random Forest</i> untuk mengatasi masalah data yang tidak seimbang seperti data pelanggan bank.	Penggunaan <i>SMOTE</i> secara signifikan meningkatkan nilai <i>AUC</i> pada model <i>DenseNet201</i> dan <i>VGG19</i> dibandingkan dengan penggunaan data yang tidak seimbang.
4	S. Priyadarshinee, M. Panda (2022)	<i>Improving Prediction of Chronic Heart Failure using SMOTE and Machine Learning</i>	Menggunakan teknik <i>SMOTE</i> dan delapan model klasifikasi untuk memprediksi kelangsungan hidup pasien penyakit gagal jantung kronis.	Masalah penelitian ini adalah prediksi kelangsungan hidup pasien penyakit jantung menggunakan model <i>machine learning</i> . Penelitian ini bertujuan untuk meningkatkan prediksi kelangsungan hidup pasien penyakit jantung dengan memanfaatkan model <i>machine learning</i> .	Teknik <i>SMOTE</i> meningkatkan akurasi keluaran klasifikasi dan model <i>DTNB</i> mencapai akurasi tertinggi yaitu 87,08% untuk memprediksi kelangsungan hidup pasien jantung.
5	U. Pujianto, W. Agung Prasetyo, A.	<i>Students Academic Performance Prediction with k-</i>	Membandingkan <i>C4.5</i> dan <i>K-Nearest Neighbor (KNN)</i>	Masalah penelitian ini adalah prediksi kinerja akademik siswa dari usia dini akan memudahkan guru	Metode <i>C4.5 Decision Tree</i> menghasilkan performa prediksi yang lebih baik

No	Penulis	Judul Penelitian	Metode Penelitian	Masalah Penelitian	Hasil Penelitian
	Rakhmat Taufani (2020)	<i>Nearest Neighbor and C4.5 on SMOTE-balanced data</i>	menggunakan metode <i>SMOTE</i> dalam klasifikasi kinerja akademik mahasiswa.	untuk memberikan bantuan kepada siswa yang memiliki kemampuan akademik di bawah rata-rata kelas atau kesulitan dalam mengikuti proses pembelajaran di kelas.	dengan nilai akurasi, <i>recall</i> , dan presisi masing-masing sebesar 71,09%, 71,63%, dan 71,54% dibandingkan dengan <i>K-Nearest Neighbor</i> .
6	Nurheri Cahyana, Siti Khomsah, Agus Sasmito Aribowo (2019)	<i>Improving Imbalanced Dataset Classification Using Oversampling and Gradient Boosting</i>	Penelitian ini menggunakan metode <i>oversampling</i> dengan teknik <i>SMOTE</i> , <i>Borderline-SMOTE</i> , dan <i>ADASYN</i> untuk mengatasi ketidakseimbangan dataset. <i>Gradient Boosting</i> digunakan sebagai model klasifikasi.	Penelitian ini bertujuan untuk mengatasi masalah klasifikasi data yang tidak seimbang dengan menggunakan metode <i>oversampling</i> dan <i>Gradient Boosting</i> . Dalam penelitian ini, dicoba tiga teknik <i>oversampling</i> (<i>SMOTE</i> , <i>Borderline-SMOTE</i> , dan <i>ADASYN</i>) untuk mengatasi ketidakseimbangan data dan melihat dampaknya terhadap akurasi klasifikasi. Penelitian ini juga mencoba mengukur performa klasifikasi menggunakan metrik akurasi, <i>recall</i> , <i>precision</i> , <i>F1-Score</i> , dan AUC.	Eksperimen menunjukkan bahwa teknik <i>oversampling</i> meningkatkan akurasi dari 2% hingga 11% untuk beberapa data seperti <i>Mammography</i> , <i>Liver Disorders</i> , <i>Diabetes (Pima Indian)</i> , <i>Indian Liver</i> , <i>Habberman</i> , dan <i>Immunotherapy</i> .

No	Penulis	Judul Penelitian	Metode Penelitian	Masalah Penelitian	Hasil Penelitian
7	Sashank Sridhar, Sowmya Sanagavarapu (2021)	<i>Handling Data Imbalance in Predictive Maintenance for Machines using SMOTE-based Oversampling</i>	Penelitian ini fokus pada pengembangan algoritma pembelajaran mesin untuk memprediksi kegagalan mesin dalam lingkungan industri. Data sintetis digunakan dalam model pemeliharaan prediktif yang mencerminkan kegagalan nyata yang terjadi di industri. Masalah ketidakseimbangan kelas data ditangani dengan menerapkan teknik <i>oversampling</i> berbasis <i>SMOTE</i> .	Penelitian ini bertujuan untuk mengatasi masalah ketidakseimbangan data dalam model prediksi pemeliharaan prediktif untuk mesin industri. Ketidakseimbangan kelas data dapat menghambat kinerja algoritma pembelajaran mesin, dan masalah ini ditangani dengan menerapkan teknik <i>oversampling</i> berbasis <i>SMOTE</i> .	Hasil eksperimen menunjukkan bahwa penggunaan teknik <i>SMOTE</i> meningkatkan skor <i>AUC</i> sebesar 7,83%, yang menghasilkan peningkatan kinerja <i>Random Forest</i> dalam membedakan antara contoh kegagalan dan contoh <i>non-kegagalan</i> pada mesin. Dengan menerapkan teknik <i>SMOTE</i> , kinerja <i>Random Forest</i> dalam mengidentifikasi instansi kegagalan dan <i>non-kegagalan</i> meningkat.
8	Raneem Qaddoura and Mariam M. Biltawi (2022)	<i>Improving Fraud Detection in An Imbalanced Class Distribution Using Different Oversampling Techniques</i>	Ketidakseimbangan data ini dapat menyebabkan model pembelajaran mesin yang dilatih pada data ini menjadi bias	Menggunakan lima teknik <i>oversampling</i> yaitu: <i>SMOTE</i> , <i>ADASYN</i> , <i>Borderline1 SMOTE</i> , <i>Borderline2 SMOTE</i> , dan <i>Support Vector Machine SMOTE (SVM SMOTE)</i> untuk mengatasi masalah ketidak seimbangan data.	Teknik <i>oversampling</i> terbukti dapat meningkatkan kinerja model , dan Teknik <i>oversampling</i> terbaik adalah menggunakan <i>SMOTE</i>

No	Penulis	Judul Penelitian	Metode Penelitian	Masalah Penelitian	Hasil Penelitian
			terhadap transaksi <i>non-fraud</i> .		
9	Nur Silviyah Rahmi, Ni Wayan Surya Wardhani, Maria Bernadetha Mitakda, Regina Syahla Fauztina, dan Imelda Salsabila (2022)	<i>SMOTE Classification and Random Oversampling Naive Bayes in Imbalanced Data: (Case Study of Early Detection of Cervical Cancer in Indonesia)</i>	Mengadopsi teknik sampling <i>SMOTE</i> dan <i>Random Oversampling (ROS)</i> untuk mengatasi data tidak seimbang	Keberadaan data tidak seimbang membuat kinerja metode klasifikasi dalam <i>machine learning</i> menurun	Model klasifikasi <i>SMOTE Naive Bayes</i> memiliki akurasi yang lebih tinggi 0,61% daripada model klasifikasi <i>ROS</i> dengan <i>Naive Bayes</i> .
10	Aum Patil, Aman Framewala, dan Faruk Kazi (2022)	<i>Explainability of SMOTE Based Oversampling for Imbalanced Dataset Problems</i>	Menggunakan metode resampling <i>oversampling SMOTE</i> .	Keberadaan data tidak seimbang dapat membuat kinerja model pembelajaran mesin menurun	Model yang dilatih pada data yang diberlakukan <i>oversampling</i> dapat memprediksi lebih daripada model yang dilatih pada data yang tidak seimbang. Hasil penelitian menunjukkan bahwa <i>oversampling</i> menggunakan <i>SMOTE</i> dapat meningkatkan kinerja <i>XAI</i> .

Berdasarkan penelitian terdahulu pada Tabel 2.1 diketahui bahwa metode *SMOTE* dapat digunakan dalam penanganan *imbalanced data* yang kemudian menghasilkan nilai akurasi yang cukup bagus. Penelitian ini menerapkan metode *Random Forest* dan *C4.5* pada analisis kinerja model pada data yang tidak seimbang dan sebagai perbedaan penelitian tugas akhir ini terletak pada studi kasus yang diterapkan.

2.2. Landasan Teori

2.2.1. *Imbalanced Dataset*

Setiap kumpulan data dengan distribusi kelas yang tidak merata secara teknis dianggap mengalami ketidakseimbangan kelas [1]. Ketidakseimbangan kelas terjadi ketika jumlah data yang mewakili suatu kelas jauh lebih rendah dibandingkan dengan kelas lainnya. Kelas yang memiliki jumlah lebih sedikit dikenal sebagai kelas minoritas, sedangkan kelas lainnya disebut sebagai kelas mayoritas. Umumnya, kelas mayoritas memiliki dominasi yang lebih signifikan. Akibatnya, *classifier* cenderung membuat prediksi yang lebih condong ke arah kelas mayoritas sambil mengabaikan kelas minoritas.

Berbagai teknik telah diusulkan dan dibahas untuk mengatasi masalah data yang tidak seimbang. Beberapa solusi yang umum mencakup pendekatan pada tingkat data, pendekatan pada algoritma, sensitivitas terhadap biaya, dan penggunaan ansambel pembelajaran [20]. Dalam praktiknya, pendekatan pada tingkat data melibatkan tindakan seperti *undersampling*, *oversampling*, atau kombinasi keduanya.

2.2.2. Pendekatan Level Data

Pada pendekatan level data, ada berbagai teknik *resampling* dan sintesis data yang digunakan untuk mengatasi ketidakseimbangan distribusi kelas dalam data latih. Namun, di tingkat algoritma, pendekatan utamanya adalah melakukan penyesuaian pada operasi algoritma yang sudah ada, sehingga pengklasifikasi (*classifier*) lebih responsif terhadap klasifikasi kelas minoritas [21]. Teknik *resampling* adalah salah satu teknik *preprocessing* di mana distribusi data diseimbangkan kembali untuk mengurangi efek distribusi kelas tidak seimbang dalam proses pembelajaran [22] teknik *resampling* menyamakan distribusi kelas secara algoritmik untuk meningkatkan *imbalance ratio* [23] dan mengurangi efek distribusi kelas tidak seimbang dalam proses pembelajaran *machine learning*. Teknik *resampling* dapat dilakukan dengan metode *undersampling*, *oversampling*, dan gabungan keduanya (*hybrid*).

2.2.3. Oversampling

Oversampling adalah salah satu metode untuk mengatasi ketidakseimbangan kelas dalam dataset dengan cara menambahkan lebih banyak sampel dari kelas minoritas. *Oversampling* dapat meningkatkan jumlah data dari kelas minoritas sehingga distribusi antara kelas mayoritas dan minoritas menjadi lebih seimbang [5]. Teknik *oversampling* yang akan digunakan dalam penelitian ini adalah *SMOTE*.

Metode *Synthetic Minority Oversampling Technique (SMOTE)* adalah salah satu bentuk dari teknik *oversampling*. Diciptakan pertama kali oleh Nithes V. Chawla, pendekatan ini berfokus pada pembuatan replika data pada kelas minoritas. Replika data ini juga dikenal sebagai data sintetis.

SMOTE bekerja dengan mencari k tetangga terdekat (*K-Nearest Neighbors*) untuk setiap data pada kelas minoritas. Kemudian, data sintetis dibuat dengan jumlah persentase tertentu dari duplikasi data minoritas yang diinginkan. Pemilihan *K-Nearest Neighbors* ini dilakukan secara acak. *SMOTE* melakukan sintesis atau meng-*generate* data sampel baru berupa data vektor yang diambil dari kelas minoritas. Data yang digunakan untuk *SMOTE* adalah data yang sudah dilakukan *feature extraction*. Sintesis yang dibuat ditentukan dengan nilai k (tetangga terdekat). Berikut langkah-langkah atau algoritma yang dilakukan oleh *SMOTE* [5]:

Kode Program 2.1. Algoritma *SMOTE*(T , N , k)

Algoritma *SMOTE*(T , N , k)

Input: Jumlah sampel kelas minoritas T ; Persentase *SMOTE* $N\%$;

Jumlah tetangga terdekat k

Output: $(N/100) * T$ sampel kelas minoritas sintetis

1. (* Jika N kurang dari 100%, acak sampel kelas minoritas karena hanya sebagian persentase acak dari mereka yang akan di-*SMOTE*. *)
 2. If $N < 100$
 3. then acak sampel kelas minoritas sebanyak T

$$T = (N/100) * T$$

$$N = 100$$
 4. endif
-

```

5. N = (int)(N/100) (* Jumlah SMOTE diasumsikan dalam kelipatan integral dari 100.
   *)
6. k = Jumlah tetangga terdekat
7. numattrs = Jumlah atribut
8. Sampel[ ][ ]: array untuk sampel kelas minoritas asli
9. newindex: menyimpan hitungan jumlah sampel sintetis yang dihasilkan, diinisialisasi
   menjadi 0
10. Sintetis[ ][ ]: array untuk sampel sintetis
    (*Hitung k tetangga terdekat hanya untuk setiap sampel kelas minoritas*)
11. for i ← 1 hingga T
12.     Hitung k tetangga terdekat untuk i, dan simpan indeksinya dalam nnarray
13.     Populate(N, i, nnarray)
14. endfor
15. Populate(N, i, nnarray) (* Fungsi untuk menghasilkan sampel sintetis. *)
16. while N = 0
17.     (*Pilih angka acak antara 1 dan k, sebutlah nn. Langkah ini memilih salah
    satu dari k tetangga terdekat dari i*).
18.     for attr ← 1 hingga numattrs
19.         Hitung: dif = Sample[nnarray[nn]][attr] - Sample[i][attr]
20.         Hitung: gap = angka acak antara 0 dan
21.         Sintetis[newindex][attr] = Sample[i][attr] + gap * dif
22.     endfor
23.     newindex++
24.     N = N - 1
25. endwhile
26. return (* Akhir dari Populate. *)

```

Rumus untuk membuat data sintesis baru ditunjukkan pada formula 2.1 Di mana, x_{new} adalah data sintesis baru, x_i adalah data sampel asli kelas minoritas, \hat{x} adalah salah satu nilai dari K -Nearest Neighbors, δ merupakan nilai random dari 0 hingga 1. Lakukan sintesis ini berulang kali hingga jumlah kelas minoritas menyerupai kelas mayoritas.

$$x_{new} = x_i + (\hat{x} - x_i) \times \delta \quad (2.1)$$

2.2.4. *Imbalance Ratio*

Imbalance Ratio (IR) merujuk pada perbandingan antara jumlah sampel dalam kelas mayoritas dan jumlah sampel dalam kelas minoritas dalam sebuah dataset [1]. Beberapa peneliti menganggap bahwa kumpulan data mengalami masalah ketidakseimbangan jika IR lebih tinggi dari 3 [24]. Namun derajat ketidakseimbangannya dapat dibagi ke dalam kelompok-kelompok yang disajikan pada Tabel 2.2.

Klasifikasi optimal berhasil dicapai saat distribusi kelas seimbang dan jumlah sampel mencukupi untuk mewakili baik kelas mayoritas maupun minoritas. Hal ini penting karena sampel pelatihan yang lebih mewakili memberikan pemahaman yang lebih mendalam yang berkontribusi pada proses pembelajaran. Jika kondisi ini tidak terpenuhi, kinerja klasifikasi dapat menurun. Pada kasus dataset dengan dua kelas, perhitungan *imbalance ratio* bisa dijabarkan pada formula 2.2 [1].

$$IR = \frac{\text{Jumlah Sampel Kelas Mayoritas}}{\text{Jumlah Sampel Kelas Minoritas}} \quad (2.2)$$

Perhitungan proporsi kelas minoritas atas kelas mayoritas diberikan pada formula 2.3 [25].

$$\text{minoritas}\% = \frac{\text{Jumlah Sampel Kelas Minoritas}}{\text{Jumlah Sampel Kelas Mayoritas}} \times 100 \quad (2.3)$$

Proporsi kelas minoritas dapat digolongkan berdasarkan derajat ketidakseimbangan datanya. Derajat ketidakseimbangan datanya mengacu pada seberapa besar perbedaan antara jumlah sampel pada kelas mayoritas (kelas yang lebih banyak) dan kelas minoritas (kelas yang lebih sedikit) dalam dataset klasifikasi. Proporsi kelas minoritas menggambarkan persentase kelas minoritas dari total dataset[24].

Tabel 2. 2 Kategori Tingkat Ketidakseimbangan Data

Derajat Ketidak Seimbangan data	Proporsi Kelas Minoritas
Mild	20-40% dari dataset
Moderate	1-20% dari dataset
Extreme	<1% dari dataset

2.2.5. Pendekatan Level Data

Pendekatan level data mencakup berbagai teknik resampling dan sintesis data. teknik resampling merupakan salah satu teknik preprocessing di mana distribusi data diseimbangkan kembali untuk mengurangi efek distribusi kelas tidak seimbang dalam proses pembelajaran[23] teknik resampling menyamakan distribusi kelas secara algoritmik untuk meningkatkan *imbalance ratio* [24] dan mengurangi efek distribusi kelas tidak seimbang dalam proses pembelajaran *machine learning*. Teknik resampling dapat dilakukan dengan metode undersampling, oversampling, dan gabungan keduanya (hybrid).

2.2.6. Robust Scaller

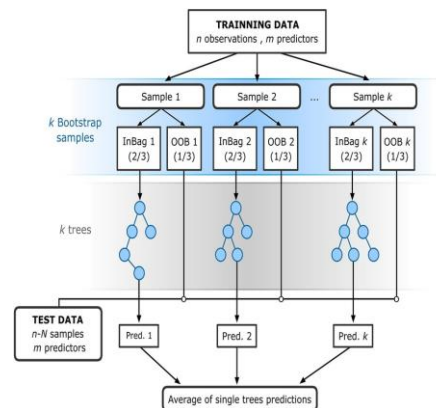
Robust Scaller adalah salah satu metode penskalaan fitur pada data yang digunakan dalam prapemrosesan data untuk analisis statistik atau pembelajaran mesin. Dengan menggunakan median dan kuartil rumus *robust scaller* dengan persamaan pada 2.7

$$X' = \frac{X - X_{Q1}}{X_{Q3} - X_{Q1}} \quad (2.7)$$

- Keterangan :
- X : Nilai awal
 - X_{Q1} : Kuartil pertama dari data
 - X_{Q3} : Kuartil ketiga dari data

2.2.7. Random Forest

Random Forest merupakan metode *ensemble* yang memanfaatkan kombinasi dari berbagai metode klasifikasi dari pemilah tunggal yang tidak stabil. Melalui proses pengambilan keputusan (*voting*), berbagai metode tersebut digabungkan untuk memberikan prediksi klasifikasi akhir yang lebih akurat [26]. Berikut merupakan alur algoritma *Random Forest* pada Gambar 2.1 [27].



Gambar 2. 1 Alur Algoritma Random Forest

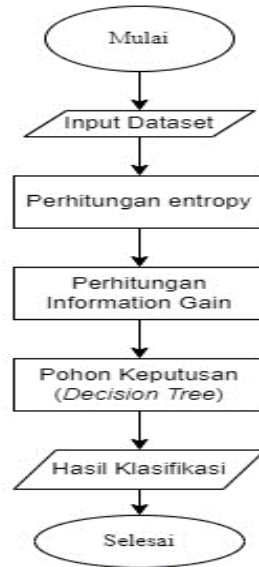
Random Forest terdiri atas kumpulan *Decision Tree* yang digunakan untuk mengklasifikasikan data ke dalam kelas-kelas tertentu[28]. *Random Forests* menggunakan pendekatan pemilihan atribut secara acak untuk membentuk beberapa pohon keputusan dengan atribut yang berbeda. Dalam *Decision Tree* konvensional, data uji diuji hanya pada satu pohon yang telah dibangun. Namun, pada *Random Forests*, data uji diuji pada semua pohon yang telah dibangun, dan klasifikasi dilakukan dengan memilih hasil mayoritas dari pemilihan yang dilakukan oleh setiap pohon [29].

2.2.8. Algoritma C4.5

Salah satu metode yang bisa dipakai untuk konstruksi pohon keputusan (*decision tree*) adalah menggunakan metode algoritma C4.5 [30]. Model algoritma C4.5 memiliki bentuk serupa pohon, dimana terdapat simpul internal yang merepresentasikan atribut-atribut, cabang-cabang menggambarkan hasil dari pengujian atribut, dan setiap ujung cabang (daun) mewakili kelas tertentu. Secara umum, data dalam struktur pohon keputusan direpresentasikan dalam format tabel yang memuat atribut-atribut. Atribut-atribut ini berperan sebagai parameter yang berfungsi sebagai kriteria dalam pembentukan pohon keputusan [31].

Dalam algoritma C4.5, dilakukan pemilihan solusi terbaik dengan menghitung dan membandingkan rasio keuntungan (*gain ratio*) hingga mencapai simpul-simpul pada level berikutnya. Algoritma C4.5 menerima masukan berupa

contoh latihan (*training samples*). Tahapan dalam pembuatan pohon keputusan menggunakan algoritma *C4.5* pada Gambar 2.2 [32].



Gambar 2. 2 Alur Algoritma *C4.5*

Berikut merupakan alur langkah-langkah dari cara kerja algoritma *C4.5*:

1. Persiapan data latih dapat dilakukan dengan memperoleh data historis sebelumnya yang telah dikelompokkan ke dalam kelas-kelas yang ditentukan.
2. Menentukan akar pohon, penentuan akar pohon dapat dilakukan dengan menghitung nilai *gain* yang tertinggi dari setiap atribut atau berdasarkan nilai *entropi* terendah. Dalam konteks *data mining*, *entropi* mengacu pada ukuran heterogenitas atau ketidakpastian dalam kumpulan data [33]. *Entropy* digunakan sebagai parameter untuk mengukur tingkat ketidakpastian dalam pengambilan keputusan. Semakin heterogen atau tidak teratur kumpulan data, semakin tinggi nilai *entropi*-nya. Sebelumnya, perlu dihitung nilai *entropi* dengan menggunakan rumus pada formula 2.8

$$Entropy (S) = \sum_{j=i}^m -p_i * \log_2 p_i \quad (2.8)$$

Keterangan	:	
i	:	Himpunan kasus
m	:	Jumlah partisi i
p_i	:	Atribut

3. Menghitung nilai *gain* pada formula 2.9

$$gain(S, A) = Entropy(s) - \sum_{I=1}^N \frac{|s_i|}{|S|} * Entropy(s_i) \quad (2.9)$$

Keterangan	:	
S	:	Himpunan kasus
A	:	Fitur
N	:	Jumlah atribut partisi A
s _i	:	Proporsi terhadap s
s	:	Jumlah kasus dalam s

4. Ulangi langkah ke-2 hingga semua *record* terpartisi. Proses partisi pohon keputusan akan berhenti disaat:

- Semua tupel dalam *record* dalam simpul *m* mendapat kelas yang sama
- Tidak ada atribut dalam *record* yang dipartisi lagi
- Tidak ada *record* didalam cabang yang kosong

2.2.9. Metrik Evaluasi

Metrik Evaluasi digunakan untuk mengevaluasi seberapa baik model bekerja pada data uji [36]. Hasil metrik evaluasi model klasifikasi berisikan berbagai metrik seperti *precision* (presisi), *recall* (sensitivitas), *F1-score*, dan lainnya. Berikut beberapa persamaan yang ada pada metrik evaluasi:

Akurasi berisi nilai yang menentukan jumlah prediksi yang benar atas jumlah prediksi oleh model. Rumus dari akurasi ditunjukkan pada formula 2.10

$$\text{Akurasi} = \frac{TP + TN}{TP + FP + TN + FN} \quad (2.10)$$

Presisi berisi nilai prediksi benar positif terhadap keseluruhan data yang diprediksi positif. Rumus dari Presisi ditunjukkan pada formula 2.10

$$\text{Presisi} = \frac{TP}{(TP + FP)} \quad (2.11)$$

Recall menggambarkan kapabilitas dalam mengidentifikasi seluruh contoh yang relevan dalam kumpulan data, sementara presisi mengindikasikan persentase titik data yang oleh model dianggap relevan dan sesuai dengan apa yang memang benar-benar relevan. Rumus *recall* ditunjukkan pada formula 2.12

$$Recall = \frac{TP}{(TP + FN)} \quad (2.12)$$

F1 Score merupakan nilai rata-rata dari presisi dan *recall* yang memasukkan perhitungan kedua metrik tersebut kedalam persamaan 2.13

$$F1 = 2 * \frac{precision * recall}{Precision + recall} \quad (2.13)$$

Metrik evaluasi, dalam bentuk tabel matriks seperti yang ditunjukkan pada Tabel 2.3, menunjukkan kinerja model klasifikasi pada sekumpulan data uji dengan nilai sebenarnya yang diketahui.

Tabel 2. 3 *Confusion Matriks*

<i>Confusion Matriks</i>		Actual	
		Positive	Negative
Predicted	Positive	TP	FP
	Negative	TN	FN

Keterangan :

- TP : True Positive
- TN : True Negative
- FP : False Positive
- FN : False Negative