

BAB II

TINJAUAN PUSTAKA DAN LANDASAN TEORI

2.1 Tinjauan Pustaka

Dalam penelitian ini, penelitian terkait dilakukan untuk mendapatkan pemahaman yang lebih mendalam tentang cara yang digunakan dan untuk memperjelas perbedaan topik penelitian. Berikut dilampirkan pada tabel 2.1 penelitian-penelitian terdahulu yang berkaitan dengan penelitian "Implementasi Algoritma *K-Nearest Neighbor* Berbasis *AdaBoost* Untuk Klasifikasi Penyakit Ginjal Kronis".

Tabel 2. 1 Penelitian Terdahulu

No.	Peneliti	Judul dan Tahun Penelitian	Metode	Masalah	Hasil
1.	Rezi Yuliani [11]	Penerapan Algoritma <i>C4.5</i> Berbasis <i>AdaBoost</i> Untuk Memprediksi <i>Financial Distress</i> Perusahaan, 2020[11]	Algoritma <i>C4.5</i> dan <i>AdaBoost</i> [11]	Meminimalisir terjadinya <i>financial distress</i> , dilakukan pengamatan dan menganalisa kondisi keuangan dari segi neraca dan laporan laba rugi yang terdapat dalam laporan keuangan[11].	Berdasarkan pengujian yang telah dilakukan dengan rasio data latih dan data uji 90% : 10% hasil <i>Accuracy</i> dari Algoritma <i>C4.5</i> adalah 72,97% dan setelah ditambah dengan <i>AdaBoost Accuracy</i> meningkat menjadi 86,49%.

No.	Peneliti	Judul dan Tahun Penelitian	Metode	Masalah	Hasil
					Algoritma <i>C4.5</i> dan <i>AdaBoost</i> baik digunakan dalam prediksi <i>Financial Distress</i> perusahaan [11].
2.	Komang Anggada Sugiarta[6]	Optimasi <i>K-Nearest Neighbor</i> Menggunakan <i>Bat Algorithm</i> Untuk Klasifikasi Penyakit Ginjal Kronis, 2019[6].	Algoritma <i>K-Nearest Neighbor</i> dan <i>Bat Algorithm</i> [6].	Tingginya angka kematian yang terus meningkat pada penderita penyakit ginjal kronis[6].	Dari hasil pengujian dengan mengukur optimasi <i>K-Nearest Neighbor</i> menggunakan metode <i>Bat Algorithm</i> didapatkan nilai akurasi 100% [6].
3.	Laila Qadrini [20]	<i>Decision Tree</i> Dan <i>AdaBoost</i> Pada Klasifikasi Penerima Program Bantuan Sosial, 2021[20]	<i>Decision Tree</i> dan <i>AdaBoost</i> [20]	Tingginya angka kemiskinan yang terjadi di Indonesia, pemerintah mengupayakan penerima program bantuan social[20].	Metode klasifikasi <i>Decision Tree</i> dan <i>AdaBoost</i> sama-sama memiliki hasil <i>accuracy</i> yang baik yaitu 94% dan 95% [20].

No.	Peneliti	Judul dan Tahun Penelitian	Metode	Masalah	Hasil
4.	Hartati[21]	Optimasi Analisis Sentimen Pada <i>Twitter</i> Olshop Tokopedia Menggunakan Textmining Dengan Algoritma <i>Naïve Bayes</i> dan <i>AdaBoost</i> , 2022[21].	<i>Naïve Bayes</i> dan <i>AdaBoost</i> [21].	Terdapat isu-isu terkait pelayanan maupun opini masyarakat di <i>twitter</i> Tokopedia yang bersifat positif dan negative[21].	Algoritma <i>Naïve Bayes</i> yang ditambahkan <i>feature Synthetic Minority Over-sampling Technique (SMOTE)</i> , yang dioptimasi dengan <i>AdaBoost</i> mendapatkan nilai <i>accuracy</i> 94,95%, <i>precision</i> 90,86%, <i>recall</i> 100% dan AUC 0,950[21].
5.	Amin Nur Rais[22]	Optimasi Akurasi Klasifikasi Pada Prediksi <i>Smokte Detection</i> dengan Menggunakan Algoritma <i>AdaBoost</i> , 2022[22].	Algoritma <i>AdaBoost</i> dan <i>Naïve Bayes</i> [22].	Permasalahan mengenai kebakaran yang menjadi ancaman bagi alam dan lingkungan[22].	Hasil penelitian menunjukkan algoritma <i>AdaBoost</i> meningkatkan nilai akurasi 98,70%, presisi 97,10%[22].

2.2 Landasan Teori

Penelitian terkait klasifikasi penyakit ginjal kronis, yang menerapkan metode *K-Nearest Neighbor* berbasis *AdaBoost* memerlukan berbagai bidang keilmuan untuk membantu penelitian ini, berikut landasan teori yang dibutuhkan.

2.2.1 Penyakit ginjal kronis

Persoalan penyakit ginjal kronis kerap kali mengalami peningkatan yang terjadi pada keadaan jasmani manusia secara global[1]. Penyakit ginjal kronis hadir lebih dari tiga bulan sebagai abnormalitas struktural ginjal dengan atau tanpa penurunan laju filtrasi glomerulus (GFR), didapati kelainan patologis dan ditandai kerusakan ginjal yang dapat berupa anomali laboratorium darah atau urin maupun anomali radiologis[23]. Dengan LPG <60 mL/menit/ $1,73 m^2$ selama kurang dari tiga bulan dan dapat disertai atau tanpa disertai kerusakan ginjal[23]. Ditemukan sebagian faktor risiko penyebab penyakit ginjal kronis seperti *hipertensi*[23], *diabetes melitus*[23], riwayat keluarga penyakit ginjal kronis[23] dan penyakit *kardiovaskular*[23]. Berikut data [19] yang digunakan dalam klasifikasi penyakit ginjal kronis dengan adanya pembagian kategori berdasarkan atribut untuk perhitungan *K-Nearest Neighbor* berbasis *AdaBoost*.

Tabel 2. 2 Kategori Berdasarkan Atribut Untuk Perhitungan Manual *K-Nearest Neighbor* berbasis *AdaBoost*

No.	Atribut	Keterangan Pembagian Kategori
1.	<i>Age</i> [19]	Usia 2 tahun \leq 25 tahun, 26 tahun \leq 49 tahun, 50 tahun \leq 82 tahun[19], [26].
2.	<i>Blood Pressure (BP)</i> [19]	Tekanan darah rendah \leq 90 mmHg, normal \leq 120 mmHg hipertensi stadium 2 \geq 140 mmHg, krisis hipertensi \geq 180 mmHg[19], [27], [28], [29], [30].

3.	<i>Specific Gravity</i> (SG)[19]	<i>Specific gravity</i> ≤ 1.010 yaitu mengencerkan urine, > 1.010 yaitu memekatkan urine[19], [31].
4.	<i>Albumin</i> (AL)[19]	Kadar <i>albumin</i> normal $3 \leq 5$ g/dL, kadar <i>albumin</i> rendah < 3 g/dL[19], [32].
5.	<i>Sugar</i> (SU)[19]	<i>Sugar</i> normal $4 \leq 5$ mmol/L, rendah ≤ 3 mmol/L[19], [33].
6.	<i>Red Blood Cells</i> (RBC)[19]	<i>Normal, abnormal</i> [19].
7.	<i>Pus Cell</i> (PC)[19]	<i>Normal, abnormal</i> [19].
8.	<i>Puss Cell Clumps</i> (PCC)[19]	<i>Present, notpresent</i> [19].
9.	<i>Bacteria</i> (BA)[19]	<i>Present, notpresent</i> [19].
10.	<i>Blood Glucose Random</i> (BGR)[19]	<i>Blood Glucose Random</i> rendah < 70 mg/dl, normal < 200 mg/dl, ≥ 200 mg/dl[19], [33].
11.	<i>Blood Urea</i> (BU)[19]	<i>Blood Urea</i> normal $20 \leq 40$ mg/dl, rendah < 20 mg/dl, > 40 mg/dl[19], [34].
12.	<i>Serum Creatinine</i> (SC)[19]	<i>Serum Creatinine</i> normal $0,6 \leq 1,3$ mg/dl, rendah $< 0,6$ mg/dl, tinggi $> 1,3$ mg/dl[19], [35].
13.	<i>Sodium</i> (SOD)[19]	<i>Hyponatremia</i> < 135 mmol/L, normal $135 \leq 145$ mmol/L, <i>hypertenatremia</i> > 145 mmol/L[19], [36].
14.	<i>Potassium</i> (POT)[19]	<i>Hypokalemia</i> $< 3,5$ mEq/L, normal $3,5 \leq 5,0$ mEq/L, <i>hyperkalemia</i> $5,1 \leq 6,0$ mEq/L, <i>dangerously high</i> $> 6,0$ mEq/L[19], [37].

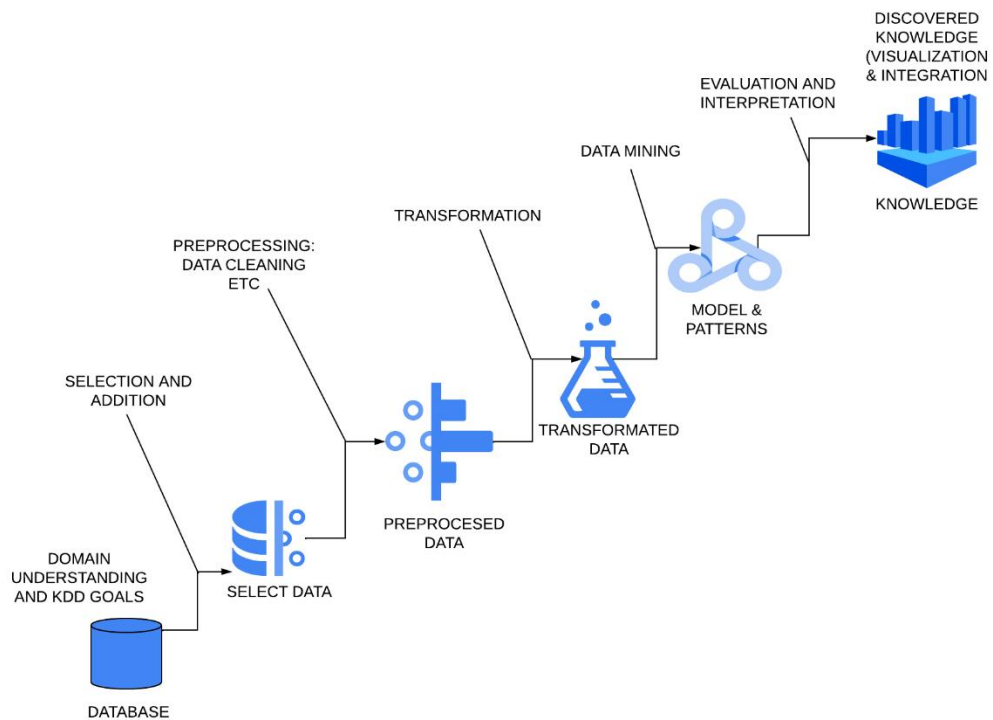
15.	<i>Hemoglobin (HEMO)</i> [19]	<i>Hemoglobin</i> $12 \leq 18$ g/dl, <i>anemia</i> < 12 g/dl, <i>polisitemia</i> > 18 g/dl[19], [38].
16.	<i>Packed Cell Volume (PCV)</i> [19]	<i>Packed Cell Volume</i> rendah < 35 , normal $35 \leq 47$, tinggi > 47 [19], [39].
17.	<i>White Blood Cell Count (WC)</i> [19]	<i>Leukopenia</i> < 4000 , normal $4000 \leq 11000$, <i>leukositosis</i> > 11000 [19], [40].
18.	<i>Red Blood Cell Count (RC)</i> [19]	<i>Mikrositik</i> $< 4,5$, normal $4,5 \leq 6$, <i>makrositik</i> > 6 [19], [41].
19.	<i>Hypertension (HTN)</i> [19]	<i>Yes, no</i> [19]
20.	<i>Diabetes Mellitus (DM)</i> [19]	<i>Yes, no</i> [19]
21.	<i>Coronary Artery Disease (CAD)</i> [19]	<i>Yes, no</i> [19]
22.	<i>Appetite (APPET)</i> [19]	<i>Good, poor</i> [19].
23.	<i>Pedal Edema (PE)</i> [19]	<i>Good, poor</i> [19]
24.	<i>Anemia (ANE)</i> [19]	<i>Yes, no</i> [19]

2.2.2 Data mining

Tahap perencanaan *data mining* memuat informasi yang akurat untuk perkiraan berdasarkan tren masa lalu dan kondisi saat ini[42]. *Data mining* merupakan analisis melihat gabungan data untuk mendeteksi kesinambungan yang tidak terduga dan meringkas data dengan cara yang berbeda dari sebelumnya untuk mudah dimengerti dan berguna bagi pemilik data[43]. Dalam beberapa kasus, banyak teknik yang termasuk dalam *data mining* dapat dieksploitasi, salah satunya adalah teknik klasifikasi, klasifikasi itu sendiri adalah bentuk dasar dari analisis data[44].

Berdasarkan kategori kegunaan *data mining* meliputi deskriptif dan prediktif, deskriptif didefinisikan sebagai menemukan pola dan kesinambungan yang tidak

diketahui dalam kumpulan data, dan *data mining* yang dimanfaatkan untuk membentuk sebuah model pengetahuan dalam klasifikasi adalah prediktif, beberapa operasi *data mining* yang termasuk deskriptif yaitu *clustering*, *summarization*, dan *sequence discover*, pada kegunaan *data mining* yang termasuk prediktif adalah klasifikasi, regresi, *time series analysis*, dan *prediction*[45]. *Data mining* memiliki beberapa cara yang dimanfaatkan pada alur aplikasi penambangan data, salah satu strategi yang umum dimanfaatkan adalah *Knowledge Discovery Databases* (KDD) yang terdiri dari 9 fase yang ditampilkan pada Gambar 2.1[46].



Gambar 2. 1 Tahapan *Knowledge Discovery Databases* (KDD)[46]

Setiap langkah *Knowledge Discovery Databases* (KDD), dijelaskan sebagai berikut[46]:

1. *Domain understanding and KDD Goals*

Sasaran ditentukan dari sudut pandang pengguna dalam membantu mengembangkan dan memahami aplikasi dan informasi sebelumnya[46].

2. *Selection and additions*

Langkah kedua, berpusat mendefinisikan data target dan himpunan bagian dari sampel atau data atribut, tahap pembagian data latih dan uji yang akan diolah menggunakan rasio 70:30[47].

3. *Preprocessing*

Prapemrosesan data adalah operasi mendasar untuk menyelesaikan data yang konsisten tanpa gangguan *missing value*[46]. Menangani *missing value* dengan metode imputasi[48]. Mengganti nilai kosong dengan nilai rata-rata, dari nilai atribut berdasarkan masing-masing jenis atribut[49]. Kedua, menangani dengan menghapus data atribut, dengan syarat data kosong lebih dari 20%, karena semakin tinggi data kosong yang diisi mempengaruhi keakuratan data, menggunakan rumus persentase untuk mengetahui berapa persen *missing value* yang ada pada atribut, dihitung menggunakan rumus berikut[50].

$$P = \frac{f}{n} \times 100\% \quad (2.1)$$

Keterangan:

P = persentase

f = frekuensi data *missing value*

n = jumlah keseluruhan data

100% = angka tetap

Rumus untuk mengetahui nilai rata-rata, yaitu[51]:

$$\bar{X} = \frac{\sum X_i}{n} \quad (2.2)$$

Keterangan :

\bar{X} = nilai rata – rata

X = nilai data

n = banyak data

4. *Transformation*

Mengkonversi nilai nominal ke tipe data numerik, sehingga informasi mudah diimplementasikan[46].

5. *Data mining (choosing the suitable data mining task)*

Pilih algoritma penambangan data yang sesuai berdasarkan tujuan spesifik yang ditentukan di fase awal, contohnya teknik *data mining* adalah klasifikasi, regresi, pengelompokan dan peringkasan[46].

6. *Data mining (choosing the suitable data mining algorithm)*

Saat menentukan algoritma yang tepat untuk menemukan pola data, algoritma yang dipilih didasarkan pada konsistensi metode *data mining*[46].

7. *Data mining (implying data mining algorithm)*

Mengaplikasikan algoritma yang telah dipilih[46].

8. *Evaluation and interpretation*

Fase ini berfokus pada interpretasi dan evaluasi, memeriksa pola atau informasi yang ditemukan yang mungkin atau mungkin tidak bertentangan dengan hipotesis yang sudah ada sebelumnya[46].

9. *Discovered knowledge*

Dalam proses KDD, informasi yang ditemukan dapat dimanfaatkan untuk memutuskan apa yang harus dilakukan dengan informasi yang dihasilkan[46].

2.2.3 Klasifikasi

Klasifikasi adalah bagian metode pembelajaran terawasi, yang membutuhkan data pembelajaran berlabel pembelajar untuk memperkirakan kelas dari objek yang tidak diketahui[52]. Klasifikasi memiliki nilai data akurasi, atau kesamaan hasil pengukuran dengan bilangan real atau data[53], hasil dari jumlah data uji yang benar dan data uji yang salah pada algoritma adalah nilai dari *accuracy*, *recall* dan *precision* yang diperoleh menggunakan *confusion matrix*[54]. Akurasi yang tinggi

diartikan sebagai klasifikasi yang berhasil dalam menyelesaikan suatu kasus dan dapat diklasifikasikan sebagai klasifikasi yang akurat dan efisien[52].

2.2.4 Algoritma *K-Nearest Neighbor*

Algoritma *K-Nearest Neighbor* (K-NN) termasuk algoritma *supervised learning*[55]. Sistem kinerja dari *supervised learning* adalah dari data pelatihan yang diberikan, sistem akan mempelajari data yang telah ada, sistem akan mempelajari pola dari data sebagai acuan untuk data selanjutnya[56]. Algoritma *K-Nearest Neighbor* adalah unsur pengklasifikasian objek berdasarkan data latih yang paling dekat dengan objek tersebut, data latih tersebut dijelaskan dengan atribut numerik n-dimensi, pada setiap data latih memiliki titik di n-dimensi mewakili ruang, diberikan data kueri yang pengidentifikasinya tidak diketahui, dengan *K-Nearest Neighbor* mencari k data pelatihan yang paling dekat dengan data kueri dalam ruang n-dimensi, jarak antara data kueri dan data pelatihan dihitung dengan mengukur jarak antara titik yang mewakili data kueri dan semua titik yang mewakili data pelatihan[57]. Jumlah tetangga yang digunakan sama dengan k tetangga, ambil k tetangga terdekat terlebih dahulu, hitung data terdekat dari k tetangga terdekat dari kelas yang ada, kelas dengan data terbaru menjadi kelas pemenang yang ditetapkan sebagai kelas pengenalan dalam data X[58].

Data numerik mampu ditangani oleh algoritma *K-Nearest Neighbor*, namun tipe non-numerik akan dikonversi atau dimasukkan ke dalam *preprocessing* data[59]. Sebelum melakukan perhitungan, langkah awal yang diambil yaitu dengan menentukan jumlah tetangga (k) dalam pertimbangan penentuan kelas, untuk menentukan nilai k tidak memiliki rumus pasti, menentukan nilai k dengan mempertimbangkan jumlah data yang ada, jumlah data yang genap menggunakan nilai k ganjil dan jumlah data yang ganjil menggunakan nilai k genap, dari *ranking* data yang dilakukan menentukan kelas data baru dari hasil kelas mayoritas yang diperoleh[60]. Tahapan dari algoritma K-Nearest Neighbor, pertama siapkan data latih, tahap kedua tentukan nilai k, ketiga menghitung jarak terdekat dapat menggunakan jarak *euclidean*, lalu meranking hasil perhitungan *euclidean*, tahap

selanjutnya menghitung jumlah titik pada k untuk setiap kelas dan tahap akhir mengklasifikasikan data baru ke kelas k terbanyak[61].

Metode jarak *euclidean* merupakan perhitungan jarak pendekatan dalam mengukur dua titik yang berbeda[62]. Metode pendekatan *K-Nearest Neighbor* yang digunakan pada penelitian ini, yaitu metode pendekatan jarak *euclidean*[63]. Persamaan perhitungan untuk mencari jarak *euclidean* adalah sebagai berikut[64]:

$$d_i = \sqrt{\sum_{1=1}^p (X_{1i} - X_{2i})^2} \quad (2.3)$$

Keterangan:

X_1 : sample data uji

X_2 : data uji

d : jarak antara X dan Y

p : dimensi data

i : setiap data

2.2.5 AdaBoost

Adaptive Boost (AdaBoost) adalah unsur varian dari algoritma *boosting*, *boosting (AdaBoost)* didefinisikan sebagai pendekatan pembelajaran mesin untuk meningkatkan aturan perkiraan yang akurat dengan menggabungkan aturan yang relatif lemah dan tidak tepat[65]. Penerapan algoritma *AdaBoost* dapat dikombinasikan dengan algoritma klasifikasi lainnya untuk memperoleh hasil kinerja klasifikasi yang baik, yaitu dengan kinerja bobot dilampirkan pada setiap pelatihan dan kemudian pengklasifikasi yang berbeda dilatih menggunakan bobot tersebut[11]. Di bawah ini adalah tahapan dari algoritma *AdaBoost*[20]:

- 1) Input: Suatu kumpulan sample penelitian dengan label $\{(x_i, y_i), \dots, (x_N, y_N)\}$ suatu *component learn* algoritma, jumlah perputaran T.
- 2) *Initialize*: Bobot suatu sampel pelatihan $w_i^1 = 1/N$, untuk semua $i = 1, \dots, N$
- 3) *Do for* $t = 1, \dots, T$

- 4) Gunakan *component learn* algoritma untuk melatih suatu komponen klasifikasi, pada sample bobot pelatihan.
- 5) Hitung kesalahan pelatihnannya pada $h_t: \varepsilon_t = \sum_i^N w_i^t, y_i \neq h_i(x_i)$
- 6) Tetapkan bobot untuk *component classifier* $h_t = \alpha_t = \frac{1}{2} \ln \left(\frac{1-\varepsilon_t}{\varepsilon_t} \right)$
- 7) *Update* bobot sample pelatihan $w_i^{t+1} = \frac{w_i^t \exp \{-\alpha_t h_t(x_i)\}}{C_t}$, $i = 1, \dots, N$ C_t adalah suatu konstanta normalisasi.
- 8) *Output*: $f(X) = \text{sign} (w_{t=1}^T \alpha_t h_t(x))$

2.2.6 Label encoding

Label-encoding bekerja dengan mengubah data menjadi bentuk numerik, yang diubah menjadi bentuk yang dapat dibaca komputer atau diproses dengan memberikan setiap kategori data nomor unik, yang dimulai dari 0[16]. Dan fungsi *One-hot encoding* mengklasifikasikan pengenalan kelas yang ada menjadi bilangan biner, angka 1 berarti pengenalan kelas ada dan angka 0 berarti pengenalan kelas tidak ada[66]. *Label encoding* merupakan bagian dari *pre-processing*[67].

2.2.7 Entropi

Entropi digunakan untuk menentukan nilai perolehan informasi, entropi menjelaskan jumlah informasi yang dibutuhkan untuk mengkodekan suatu kelas [68]. Semakin rendah nilai entropi, semakin baik digunakan untuk mengekstrak kelas[69]. Nilai entropi dapat dihitung dengan menggunakan persamaan berikut [70]:

$$\text{Entropi } S = \sum_{i=1}^n -p_i \times \log_2 p_i \quad (2.4)$$

Keterangan:

S : Himpunan kasus

n : Jumlah partisi S

p_i : Proporsi dari S_i terhadap S

2.2.8 Information gain

Information gain (IG) didefinisikan sebagai teknik pemilihan fitur yang menggunakan metode pembobotan atribut nominal atau kontinu yang didiskritisasi menggunakan entropi maksimum, istilah perolehan informasi diukur dengan menghitung jumlah bit informasi dari kelas yang diklasifikasi, dengan dokumen tersebut mengandung term atau tidak, teknik pemilihan fitur *information gain* berarti pemilihan simpul fitur dari pohon keputusan berdasarkan nilai *information gain*, nilai *information gain* suatu fitur ditentukan oleh efek dari fitur tersebut, pada satuan kelas yang diukur dalam data dibagi menjadi sub-data dengan nilai karakteristik tertentu[68]. Pemilihan atribut pada *root* didasarkan pada nilai validasi tertinggi dari atribut yang ada, *gain* (S.A) adalah kumpulan informasi tentang atribut A dalam kaitannya dengan informasi awal S, berikut adalah persamaan untuk menghitung nilai gain[70]:

$$Gain S, A = \sum_{i=1}^n \frac{S_i}{S} \times Entropy S_i \quad (2.5)$$

Keterangan:

S : Himpunan kasus

A : Atribut

n : Jumlah partisi atribut A

|S_i|: Jumlah kasus pada partisi ke-i

|S| : Jumlah kasus dalam S

2.2.9 Split Data

Proses klasifikasi membutuhkan data yang telah dibagi terlebih dahulu atau disebut juga *split data*, terdiri dari dua jenis, yaitu data latih (*data training*) dan data uji (*data testing*)[71]. Pembagian data memiliki rasio perbandingan, pada penelitian ini akan menggunakan 70% data latih dan 30% data uji[72]. Menggunakan jumlah data latih lebih banyak, untuk mempermudah pengujian kebenaran dan keakuratan pada data uji[73].

2.2.10 Confusion matrix

Teknik yang digunakan untuk mencari nilai akurasi dengan menggunakan tabel yang menjelaskan jumlah data uji yang benar dan salah disebut *confusion matrix* [74]. *Confusion matrix* menjadi alat visualisasi yang umum digunakan dalam pembelajaran terawasi, setiap kolom matriks adalah contoh kelas yang diklasifikasi dan setiap baris mewakili kejadian kelas yang sebenarnya [75]. Setiap sel berisi angka menunjukkan berapa banyak kasus yang sebenarnya dari kelas yang diamati untuk diperkirakan, TP sebagai jumlah *record positive* yang diklasifikasikan sebagai positif, FP sebagai jumlah *record negative* yang diklasifikasikan sebagai positif, TN sebagai jumlah *record positive* yang diklasifikasikan sebagai negatif, dan FN sebagai jumlah *record negative* yang diklasifikasikan sebagai negatif [76].

		PREDICTED CLASS		
		YES	NO	TOTAL
ACTUAL CLASS	YES	TP	FN	P
	NO	FP	TN	N
TOTAL		P'	N'	P + N

Gambar II.1 *Confusion Matrix* [76]

2.2.11 Accuracy, recall, precision

Ada tiga nilai yang dimanfaatkan untuk mengukur kemampuan suatu sistem klasifikasi yaitu *precision*, *recall* dan *accuracy* [77]. Berikut adalah prosedur *confusion matrix* dengan beberapa parameter untuk mengevaluasi kinerja model klasifikasi yang dihasilkan [78][79][80]:

1. Precision

Precision adalah tingkat ketepatan antara informasi yang diminta oleh pengguna dengan jawaban yang diberikan oleh system menggunakan persamaan (2.8), pada persamaan (2.6) yaitu *precision (yes)* merupakan

nilai evaluasi dari *precision* yang menghitung kelas positif dan pada persamaan (2.7) yaitu *precision (no)* merupakan nilai evaluasi yang menghitung kelas negatif. Berikut rumus persamaan (2.6), (2.7) dan (2.8) ada di bawah ini:

$$Precision (yes) = \frac{TP}{TP+FP} \quad (2.6)$$

$$Precision (no) = \frac{TN}{TN+FN} \quad (2.7)$$

$$Precision = \frac{Precision(yes)+Precision(no)}{2} \quad (2.8)$$

2. Recall

Recall merupakan tingkat keberhasilan sistem dalam menemukan kembali sebuah informasi dihitung menggunakan rumus persamaan (2.11), pada persamaan (2.9) merupakan sensitivitas yang merupakan nama lain dari *recall* untuk menghitung kelas positif dan pada persamaan (2.10) merupakan spesifisitas yang merupakan nama lain dari *recall* untuk menghitung kelas negatif. Berikut persamaan (2.9), (2.10) dan (2.11) ada dibawah ini:

$$Sensivitas = \frac{TP}{TP+FP} \quad (2.9)$$

$$Spesifisitas = \frac{TN}{TN+FP} \quad (2.10)$$

$$Recall = \frac{Sensivitas+Spesifisitas}{2} \quad (2.11)$$

3. *Accuracy* didefinisikan sebagai tingkat kedekatan antara nilai klasifikasi dengan nilai aktual dengan menggunakan rumus persamaan di bawah ini:

$$Accuracy = \frac{TP+TN}{TP+FP+TN+FN} \quad (2.12)$$

2.2.12 *Rapidminer*

Rapidminer adalah perangkat lunak yang dikembangkan oleh Dr. Markus Hofmann dari *Institute of Technology Blanchardstown* dan Raif Klinkenberg dari *rapid-i.com* yang menggunakan GUI (*Graphical User Interface*) yang membuat pengguna tidak sukar dalam menggunakan *software* tersebut[81]. *Rapidminer* menjadi bagian dari salah satu *software* dalam pengolahan *data mining*[82]. Pengolahan pada tahap awal yaitu dengan merakit dan memodelkan operator dalam suatu proses sesuai dengan teknik penambangan data yang disediakan oleh alat *Rapidminer*[83].