BAB II TINJAUAN PUSTAKA

2.1 Landasan Teori

2.1.1 Data Mining

Data mining adalah sebuah ilmu komputer dalam proses pencarian dan analisis dari kumpulan data yang besar untuk menentukan pola tertentu yang memiliki makna serta aturan (Barhate dkk., 2018). Data mining memiliki tujuan untuk merancang pola tertentu dari sejumlah kumpulan data yang besar (database) dengan secara efisien untuk menghasilkan suatu informasi atau pengetahuan sehingga dikenal dengan istilah knowledge discovery, sedangkan untuk pengenalan pola yang tepat pada data mining yang digunakan untuk menemukan pola tersembunyi dalam kumpulan data dikenal dengan istilah pattern recognition (Nabila dkk, 2021).

Proses data mining yaitu dengan mengekstraksi pengetahuan dari sejumlah data yang besar, data mining sebagai proses mengeksplorasi dan menganalisis kumpulan data untuk menemukan pola dan aturan untuk memecahkan suatu masalah, serta dalam penggunaan data mining didorong dari kebutuhan akan teknik mengambil keputusan dengan cara menganalisis, memahami, dan melihat data yang dikumpulkan dari data warehouse (PASCU, 2018). Proses pengerjaan pada data mining memiliki metodologi yang dapat digunakan, metodologi dapat dengan mudah digunakan oleh seorang pemula dalam bidang data mining karena setiap fase pengerjaan terstruktur, terdokumentasi dengan baik serta terdefinisi dengan jelas. Metodologi penelitian ini terdiri dari beberapa fase utama, yaitu fase

pemahaman bisnis, fase pemahaman data, fase pengolahan data, fase pemodelan, fase evaluasi, serta fase penyebaran hasil. (Fadillah, 2015).

Data mining memiliki beberapa tahapan dalam pengerjaannya yaitu:

1. Pembersihan Data (*Data Cleaning*)

Pembersihan data merupakan proses menghilangkan noise dan data yang tidak relevan. Pada umumnya data yang didapatkan memiliki isi yang tidak sempurna, seperti data yang hilang, data yang tidak *valid*, atau data yang salah dalam pengetikan atau *human error*, sehingga baik untuk dihilangkan karena akan mempengaruhi hasil ataupun proses dari teknik *data mining*.

2. Integrasi Data (*Data Integration*)

Proses integrasi data merupakan penggabungan data dari berbagai *database* kedalam satu *database* yang baru. Proses integrasi harus dikerjakan dengan teliti karena jika terjadi kesalahan saat proses integrasi data dapat menghasilkan hasil yang salah dan akan mempengaruhi saat pengambilan keputusan.

3. Seleksi Data (Data Selection)

Data yang terdapat pada *database* untuk dilakukan analisis tidak selalu semuanya dipakai, akan tetapi hanya data yang sesuai yang digunakan oleh karena itu dilakukan proses seleksi data.

4. Transformasi Data (*Data Transformation*)

Transformasi data merupakan perubahan data baik itu data yang digabungkan maupun data yang diubah sesuai dengan format untuk diproses.

5. Proses *Mining*

Proses *mining* merupakan proses utama saat metode diterapkan untuk menemukan pengetahuan baru yang tersembunyi dari suatu data.

6. Evaluasi Pola (*Pattern evaluation*)

Proses evaluasi pola digunakan untuk mengidentifikasi pola-pola menarik yang ditemukan dalam *knowledge based*, sehingga tahap ini akan menghasilkan pola-pola yang baru dan unik maupun model prediksi dievaluasi untuk menilai tercapainya hipotesa yang ada.

2.1.2 Klasifikasi

Klasifikasi merupakan salah satu teknik pada *data mining* yang berfungsi untuk mengelompokan data ke dalam kelas tertentu dari sejumlah kelas yang ada berdasarkan atribut-atribut tertentu (Leomongga Oktaria Sihombing, Hannie, 2021). Dasar yang diturunkan dari model klasifikasi yaitu pada analisis dari *training data*. Menurut (Ulfatul *et al.*, 2022) proses klasifikasi terbagi menjadi dua fase yatu *fase learning* dan *testing*. *Fase training* merupakan fase untuk membangun sebuah model dari data yang akan digunakan, sedangkan *fase testing* merupakan pengujian dari model yang telah dibuat dengan data lainnya untuk mengetahui akurasi dari model tersebut.

Klasifikasi memiliki 3 tahapan (Nasution, Khotimah and Chamidah, 2019).

1. Pembangunan Model

Pembangunan model merupakan data latih yang telah memiliki atribut dan kelas digunakan untuk membangun sebuah model.

2. Penerapan Model

Penerapan model merupakan penerapan model yang telah dibangun untuk menentukan kelas dari data atau objek baru.

3. Evaluasi

Evaluasi merupakan proses yang dilakukan untuk melihat akurasi terhadap data baru dari pembangunan dan penerapan model yang telah dilakukan pada tahap sebelumnya.

2.1.3 Adaptive boosting (Adaboost)

Algoritma boosting adalah konsep dari machine learning yang dapat digunakan untuk mengkombinasikan algoritma klasifikasi lain untuk meningkatkan performa klasifikasi. Adaptive boosting adalah salah satu algoritma boosting yang mampu menyelesaikan nilai error secara adaptif yang dihasilkan dari klasifikasi lemah untuk dijadikan acuan proses pelatihan klasifikasi berikutnya (Latief, Subekti and Gata, 2021). Algoritma adaboost termasuk algoritma yang populer dan cukup banyak digunakan, karena adaboost ini mudah digunakan dan diimplementasikan juga fleksibel sehingga mudah dikombinasikan dengan algoritma lain. Adaboost banyak berhasil digunakan pada beberapa bidang karena memiliki teori yang kuat, prediksi yang besar, juga sederhana (Prasetio and Susanti, 2019). Weak learner yang digunakan pada algoritma ini yaitu algoritma Decision Tree.

Algoritma Decision Tree merupakan suatu model berbasis diagram alur yang memiliki struktur menyerupai pohon, setiap *internal node* merepresentasikan proses pengujian terhadap suatu atribut, setiap cabang (*branch*) menunjukkan hasil

dari pengujian tersebut, dan *leaf node* merepresentasikan kategori kelas atau distribusi kelas yang dihasilkan. Node yang terletak di bagian paling atas disebut sebagai *root node* atau simpul akar, yang memiliki sejumlah *edge* keluar tanpa *edge* masuk. Sementara itu, *internal node* memiliki satu *edge* masuk dan beberapa *edge* keluar, sedangkan *leaf node* hanya memiliki satu *edge* masuk tanpa *edge* keluar. Decision Tree umumnya digunakan dalam proses klasifikasi untuk menentukan kelas dari suatu sampel data yang belum diketahui kategorinya, berdasarkan kelaskelas yang telah ditetapkan sebelumnya (Septhya *et al.*, 2023).

2.1.4 K-Nearest Neighbor (KNN)

Algoritma *K-Nearest Neighbor (KNN)* merupakan algoritma tertua dan paling sederhana dalam penerapannya dan efektif, (Gamadarenda and Waspada, 2020), serta akurat untuk melakukan klasifikasi pola dan regresi. Algoritma ini sederhana namun memiliki kinerja yang dapat bersaing dengan klasifikasi yang kompleks (Suwirmayanti, 2017), algoritma ini memiliki kemampuan dalam mendeteksi dan menganalisa permasalahan yang cukup kompleks dan *non-linier*, serta perhitungan dilakukan secara paralel sehingga waktu komputasi akan lebih cepat (Abu Alfeilat *et al.*, 2019).

Menurut (Arifin, 2019) Algoritma *K-Nearest Neighbor (KNN))* merupakan pengklasifikasian objek berdasarkan data pembelajaran terdekat berdasarkan dengan data sebelumnya yang dimiliki sebagai sampel untuk menemukan suatu hasil akhir. Algoritma *K-NN* bersifat nonparametrik yang berarti tidak memiliki jumlah parameter, namun parameter akan ditentukan oleh ukuran kumpulan data

pelatihan, meskipun tidak ada asumsi yang perlu dibuat untuk distribusi data yang mendasarinya (Setiyorini and Asmono, 2019).

Terdapat beberapa langkah-langkah dalam melakukan perhitungan KNN (K-Nearest Neighbor) (Arifin, 2019), berikut adalah langkah-langkah tersebut :

- a. Menentukan parameter k atau jumlah tetangga paling dekat.
- b. Menghitung jarak *euclidean* atau *query instance* pada masing-masing objek terhadap data sampel yang diberikan.
- c. Mengurutkan objek-objek ke dalam kelompok
- d. yang mempunyai jarak euclidean terkecil.
- e. Mengumpulkan kategori Y atau klasifikasi KNN (K-Nearest Neighbor).
- f. Dengan menggunakan kategori *nearest neighbor* yang paling banyak maka dapat dilakukan prediksi nilai *query instance* yang telah dihitung.

Klasifikasi algoritma *K-NN* dilakukan berdasarkan *data training* dilihat dari jarak paling dekat dengan objek berdasarkan nilai k (Setianto, Kusrini and Henderi, 2019). Jarak terdekat dapat dilakukan dengan membagi data menjadi *data training* dan *data testing*, jika *data training* dan *data testing* telah dibuat maka bisa menggunakan metode perhitungan jarak untuk menghitung jarak masing-masing data *testing* terhadap data *training*. Terdapat beberapa metode perhitungan jarak pada algoritma *KNN* (*K-Nearest Neighbor*) yaitu:

1. Euclidean Distance

Metode *euclidean* merupakan salah satu metode perhitungan jarak yang digunakan untuk mengitung jarak antara dua titik dalam ruang *euclidean*. Pada metode ini, semakin kecil jarak maka akan semakin kecil

jarak antar kedua titik. Metode ini merupakan metode yang paling sederhana dan umum digunakan (Yudhana, Sunardi and Hartanta, 2020), dan untuk mengkur tingkat kemiripan data dengan *euclidean* mengunakan rumus sebagai berikut:

$$d(x,y) = \sqrt{\sum_{i=1}^{n} \left(x_{training}^{i} - y_{testing}^{i}\right)^{2}}....(2.1)$$

Keterangan:

$$d(x,y) = jarak$$

i = variabel data

n = dimensi data

 $x^{i}_{training} = data training$

 $y^{i}_{testing} = data \ testing$

2. Chebyshev Distance

Chebyshev disebut sebagai maximum distance karena menerapkan pencarian nilai maksimum dari selisih jarak titk data x dan y di demensi ruang d (Hidayati et al., 2021). Chebyshev merupakan metode pengukuran jarak yang dihitung dari nilai absolute atau nilai mutlak dari selisih sepasang koordinat (Rani, Aziz and Sulistyono, 2022). Jarak yang diukur dengan chebyshev berdasarkan nilai mutlak dari perbedaan anta elemen-elemen pada vektor dan jumlah data yang secara otomatis harus sama. Rumus yang dapat digunakan untuk menghitung jarak dua buah objek berdasarkan chebyshev sebagai berikut:

$$d(x, y) = max_{i=1}^{n} |x_i - y_i|$$
(2.2)

3. Manhattan Distance

Manhattan (city distance) diartikan sebagai jumlah jarak dari semua atribut. Manhattan merupakan jarak geometri dari dua buah objek, digunakan untuk mengambil kasus yang cocok dari basis kasus dengan menghitung jumlah bobot absolute dari perbedaan antara kasus sekarang dan kasus lain dalam basis kasus. Untuk melakukan perhitungan jarak dapat menggunakan rumus berikut:

$$d(x,y) = \sum_{i=1}^{n} |x_i - y_i|...........(2.3)$$

2.1.5 Synthetic Minority Oversampling Technique (SMOTE)

Metode sytnthetic minority oversampling technique atau sering dikenal dengan SMOTE merupakan model pengembangan dari metode oversampling, yang diterapkan dalam rangka menangani ketidakseimbangan kelas yang ada pada dataset. Masalah ketidakseimbangan data banyak sekali muncul pada proses klasifikasi, ketidakseimbangan dapat terjadi apabila jumlah objek suatu kelas data lebih banyak (mayor) dari kelas lainnya (minor). Perbedaan jumlah data yang besar akan berdampak pada model klasifikasi, sehingga model klasifikasi tidak dapat memprediksikan kelas minor dengan tepat sehingga banyak data tes yang seharusnya ada pada kelas minor diprediksikan salah oleh model klasifikasi (Prianggi, Nilogiri and Umilasari, 2021).

Cara kerja pada metode *SMOTE* ini yaitu dengan membuat data sintesis pada kelas data minor sehingga data menjadi seimbang (Arifiyanti and Wahyuni, 2020). Beberapa algoritma yang dikombinasikan dengan metode *SMOTE* terbukti

mengatasi ketidakseimbangan data dan mendapatkan hasil klasifikasi yang lebih baik (Sutoyo and Fadlurrahman, 2020).

2.1.6 Confusion Matrix

Confusion matrix merupakan proses untuk melakukan perhitungan akurasi pada konsep data mining (Ernawati and Wati, 2018). Terdapat 4 kombinasi yang berbeda pada confusion matrix yaitu terdapat nilai prediksi dan aktual. Adapun tabel evaluasi model confusion matrix (Said et al., 2022), sebagai berikut:

Tabel 2. 1 Confusion Matrix

Confusion Matrix		Nilai Aktual	
		Positif	Negatif
Nilai	Positif	True Positive	False Negative
Prediksi	Negatif	False Positive	True Negative

Keterangan:

- True Positive (TP): jumlah data bernilai positif yang diklasifikasikan sebagai positif.
- 2. False Positive (FP): jumlah data bernilai negatif yang diklasifikasikan positif.
- 3. False Negative (FN): jumlah data bernilai positif yang diklasifikasikan negatif,
- 4. *True Negative* (TN): jumlah data bernilai negatif yang diklasifikasikan negatif.

2.1.7 *Recall*

Recall merupakan rasio antara jumlah dengan prediksi benar pada kelas positif atau true positif (TP) terhadap total data yang sebenarnya masuk dalam kelas positif, dengan menjumlahkan antara true positif dan false negatif (TP + FN). Nilai ini merepresentasikan sejauh mana model klasifikasi mampu mengidentifikasi dan menemukan informasi yang benar-benar relevan (Clara et al., 2021)

2.2 Penelitian Terkait dan Kebaruan Penelitian

2.2.1 State Of The Art (SOTA)

Tabel 2. Membahas mengenai penelitian sebelumnya yang telah dilakukan sebagai perbandingan dengan fokus pada penelitian yang akan dilakukan yaitu pada peningkatan akurasi algoritma klasifikasi. Terdapat beberapa kesamaan serta perbedaan dari masing-masing penelitian yang dapat dilihat dari penggunaan metode dan algoritmanya.

Tabel 2. 2 State Of The Art

No.	Nama	Tahun	Judul	Isi Ringkasan	Hasil
	Pengarang				
1.	Tessy	2020	Machine Learning	Stroke adalah penyebab utama	Sebelum dilakukan klasifikasi,
	Badriyah, Nur		Algorithm for	kematian dan obesitas nomor satu di	dilakukan pengolahan citra dan
	Sakinah, Iwan		Stroke Disease	banyak negara. Penelitian ini	ekstraksi ciri pada data citra. Dan
	Syarif, Daisy		Classification	melakukan preprocessing untuk	setelah itu digunakan perbandingan 8
	Rahmania			meningkatkan kualitas citra dan	(delapan) algoritma untuk melakukan
	Syarif			mengurangi noise, serta menerapkan	klasifikasi yaitu : KNN (K-Nearest
				algoritma machine learning untuk	Neighbor)s, naive bayes, logistic
				mengklasifikasikan pasien menjadi dua	regression, decision tree, random
				sub tipe penyakit stroke, yaitu stroke	forest, multi-layer perceptron (mlp-
				iskemik dan stroke pendarahan.	nn), deep learning and support vector
				Delapan algoritma digunakan dalam	machine. Dari hasil perbandingan ini
				penelitian ini salah satunya yaitu KNN	algoritma K-NN memiliki nilai
				(K-Nearest Neighbor).	akurasi 96%.

No.	Nama	Tahun	Judul	Isi Ringkasan	Hasil
	Pengarang				
2.	Iswanto,	2021	Comparison of	Stroke adalah penyakit kardiovaskular	Berdasarkan hasil pengujian, model
	Tulus, Poltak		Distance Models on	(CVD) yang disebabkan oleh	chebyshev memiliki tingkat akurasi
	Sihombing		KNN (K-Nearest	kegagalan sel-sel otak untuk	tertinggi dibandingkan ketiga model
			Neighbor)	mendapatkan suplai oksigen sehingga	jarak lainnya dengan nilai akurasi
			Algorithm in Stroke	menimbulkan risiko kerusakan iskemik	rata-rata 95,49%, akurasi tertinggi
			Disease Detection	dan mengakibatkan kematian.	96,03%, pada K = 10. Model jarak
				Mendeteksi stroke membutuhkan	<i>euclidean</i> dan <i>minkowski</i> memiliki
				metode pembelajaran mesin. Dalam	tingkat akurasi yang sama pada setiap
				penelitian ini, menggunakan salah satu	nilai K dengan nilai akurasi rata-rata
				dari metode klasifikasi pembelajaran	95,45%, akurasi tertinggi 95,93%
				terbimbing yaitu KNN (K-Nearest	pada K = 10. Sedangkan <i>Manhattan</i>
				Neighbor) (K NN). Penelitian ini	memiliki rata-rata terendah
				membandingkan model jarak	dibandingkan model jarak lainnya,
				euclidean, minkowski, manhattan ,	yaitu 95,42% tetapi memiliki akurasi
				chebyshev untuk mendapatkan hasil	tertinggi 96,03% pada nilai K = 6.
				yang optimal.	Dengan demikian, model jarak yang
					paling optimal yang digunakan untuk

No.	Nama	Tahun	Judul	Isi Ringkasan	Hasil
	Pengarang				
					dataset penyakit stroke adalah model
					jarak <i>chebyshev</i> dengan nilai K
					optimal 10. Selanjutnya model
					perhitungan jarak manhattan dengan
					nilai K optimal th adalah 6
3.	Kartika	2018	Detection of Atrial	Seseorang yang tidak pernah memiliki	Struktur sistem diperoleh dengan
	Resiandi,		Fibrillation	riwayat penyakit jantung bahkan	menggunakan metode klasifikasi K-
	Adiwijaya,		Disease Based on	berpeluang menderita AF. Risiko yang	NN. Penelitian ini menggunakan
	Dody Qori		Electrocardiogram	ditimbulkan oleh AF, yaitu	parameter K=1,3,5,7,9, dan 11
	Utama		Signal	kemungkinan stroke, gagal jantung,	dengan train/test split untuk
			Classification	dan kematian. Bagi seseorang yang	mendapatkan akurasi tertinggi.
			Using RR Interval	sudah memiliki gejala AF sebaiknya	Kinerja algoritma <i>K-NN</i> diukur
			and KNN (K-	segera memeriksakan diri salah satunya	dengan menggunakan confusion
			Nearest Neighbor)	dengan menggunakan alat	matrix. Berdasarkan akurasi skema
				elektrokardiogram (EKG). Karena	keseluruhan, hasil terbaik adalah k =
				dengan adanya deteksi dini dapat	1 dengan akurasi rata-rata sebesar
				menurunkan jumlah persentase	91,75% dan tingkat akurasi,

No.	Nama	Tahun	Judul	Isi Ringkasan	Hasil
	Pengarang				
				populasi AF, dan prognosis penyakit	sensitivitas, spesifisitas tertinggi
				AF juga lebih baik. Ada tiga tahapan	sebesar 95,45%, 91,67%, dan 100%
				dalam penelitian ini; yaitu pra-	dengan data train/test split sebesar 60
				pemrosesan sebagai proses	:40 persen. Disimpulkan bahwa
				penyeragaman dimensi data, ekstraksi	diagnosis AF jantung dengan metode
				ciri, dan klasifikasi K-NN.	K-NN sudah memadai untuk
					pemeriksaan medis.
4.	Agus Byna,	2020	Penerapan Metode	Stroke adalah penyakit paling	Hasil penelitian untuk nilai akurasi
	Muhammad		Adaboost Untuk	mematikan nomor dua di dunia	algoritma Naïve Bayes memiliki nilai
	Basit		Mengoptimasi	menurut WHO. Saat ini perkembangan	0.976 dengan Split data 80/20,
			Prediksi Penyakit	Era Revolusi Industri 4.0 yang	sedangkan untuk nilai akurasi
			Stroke Dengan	berkolaborasi di bidang teknologi dan	optimasi <i>adaboost</i> dengan <i>naïve</i>
			Algoritma Naïve	ilmu kesehatan sehingga menjadi	bayes senilai 0.981 split data 70/30
			Bayes	sesuatu yang bermanfaat dengan	kedua model tersebut memiliki
				menggunakan Machine Learning.	diagnosa Excellent Classification,
				Banyak sekali manfaat yang digunakan	dalam pengujian prediksi penyakit
				dalam memprediksi beberapa penyakit	stroke dengan dengan 11 variabel

No.	Nama	Tahun	Judul	Isi Ringkasan	Hasil
	Pengarang				
				yang dapat diantisipasi. Khususnya	dengan jumlah data 28,500 untuk data
				penyakit stroke dengan menggunakan	training dan 572 untuk data testing
				algoritma adaboost yang mempunyai	(hasil dari kuisioner dan aplikasi di
				kelebihan bisa digabung dengan	rumah sakit). Algoritma adaboost
				algoritma naïve Bayes. Penerapan	dapat meningkatkan dan
				metode ini menggunakan split data	mengoptimasi yang dapat digabung
				yaitu data training dan data testing	dengan naïve bayes sebagai algoritma
				dibuat porsi dalam melakukan	estimator sehingga menghasilkan
				pengujian	akurasi yang dapat meningkatkan
					hasil yang terbaik dengan memiliki
					margin <i>error</i> yang kecil.
5.	Nia Novianti,	2022	Penerapan	Data – data penyakit penderita	Berdasarkan hasil penelitian dan
	Muhammad		Algoritma	diabetes terdahulu sudah tersimpan dan	pengujian yang telah dilakukan maka
	Zarlis, Poltak		Adaboost Untuk	tersusun data sebuah gudang	dapat ditarik kesimpulan
	Sihombing		Peningkatan	penyimpanan data atau yang biasa	bahwasannya algoritma adaboost
			Kinerja Klasifikasi	disebut dengan dataset. Maka dari itu	dapat dipergunakan dengan baik
			Data Mining Pada	perlu dilakukan proses untuk	untuk membantu kinerja dari pada

No.	Nama	Tahun	Judul	Isi Ringkasan	Hasil
	Pengarang				
			Imbalance Dataset	melakukan pengolahan data yang	algoritma KNN (K-Nearest Neighbor)
			Diabetes	terdapat pada <i>dataset</i> . Tetapi	untuk proses klasifikasi pada dataset
				penggunaan dari pada teknik data	penyakit diabetes. Hal ini dapat
				mining sendiri harus dibantu dengan	dilihat pada hasil pengujian dengan
				menggunakan teknik yang terdapat	nilai K = 7, 13, 19, 25 dan 31 terdapat
				pada data mining tersebut yaitu	peningkatan hasil akurasi yang
				teknik klasifikasi. Pada penelitian ini,	didapatkan setelah menggunakan
				menggunakan algoritma KNN (K-	algoritma adaboost. Akurasi tertinggi
				Nearest Neighbor) dengan menerapkan	terdapat pada nilai K=7, sebelum
				algoritma adaboost untuk	menerapkan algoritma adaboost
				meningkatkan akurasi metode	memiliki akurasi 92,70%, sedangkan
				klasifikasi.	setelah menerapkan algoritma
					adaboost meningkat menjadi 95,40%.
6.	Aah Sumiah,	2020	Perbandingan	Membandingkan algoritma KNN (K-	Hasil dari implementasi ini algoritma
	Nita Mirantika		Metode KNN (K-	Nearest Neighbor) dan algoritma naïve	KNN (K-Nearest Neighbor) diperoleh
			Nearest Neighbor)	bayes untuk rekomendasi penerimaan	akurasi 100%, sedangkan algoritma
			dan Naive Bayes	beasiswa pada Universitas Kuningan.	naïve bayes akurasinya 99,89%. Hal

No.	Nama	Tahun	Judul	Isi Ringkasan	Hasil
	Pengarang				
			untuk Rekomendasi	Algoritma ini digunakan karena kedua	ini menunjukan bahwa pada data
			Penentuan	algoritma tersebut biasa digunakan	penerimaan beasiswa, algoritma KNN
			Mahasiswa	dalam klasifikasi data. Dari penelitian	(K-Nearest Neighbor) lebih baik
			Penerima Beasiswa	ini mengimplementasikan menjadi	karena memiliki nilai akurasi lebih
			pada Universitas	sistem informasi yang menggunakan	tinggi dibandingkan dengan <i>naïve</i>
			Kuningan	visual basic.net dan sql server.	bayes. Banyaknya data latih akan
					mempengaruhi keakuratan data.
7.	Dita Noviana,	2019	Analisis	Beasiswa merupakan bantuan biaya	Dari penelitian ini, dalam
	Yuliana		Rekomendasi	pendidikan yang sangat membantu	merekomendasikan penerima
	Susanti, Irwan		Penerima Beasiswa	prestasi mahasiswa. Seiring dengan	beasiswa PPA dapat dilakukan
	Susanto		Menggunakan	meningkatnya jumlah mahasiswa yang	menggunakan algoritma K-NN dan
			Algoritma KNN (K-	mengajukan beasiswa, maka	C.45 dengan akurasi masing-masing
			Nearest Neighbor)	dibutuhkan suatu metode klasifikasi	sebesar 90.7% dan 88.3%. Hasil
			(K-Nn) Dan	yang dapat membantu menentukan	penelitian menunjukkan bahwa
			Algoritma C4.5	siapa yang layak mendapatkan	variabel memiliki pengaruh adalah
				beasiswa PPA dari Universitas Sebelas	IPK, prestasi, penghasilan, jumlah
				Maret. Metode yang digunakan dalam	orang tua bekerja, dan jumlah

No.	Nama	Tahun	Judul	Isi Ringkasan	Hasil
	Pengarang				
				analisis ini adalah KNN (K-Nearest	tanggungan. Variabel yang paling
				Neighbor) dan C4.5. Hasil klasifikasi	berpengaruh adalah variabel IPK.
				digunakan sebagai keputusan dalam	Dalam mengukur kedua fungsi
				rekomendasi penerima beasiswa.	algoritma tersebut menggunakan
					confusion matrix dengan hasil bahwa
					algoritma K-NN memiliki tingkat
					akurasi yang lebih tinggi
					dibandingkan algoritma C4.5.
8.	Riski Tri	2019	Prediksi Harapan	Operasi bedah toraks menjadi salah	Untuk melakukan prediksi ini yaitu
	Prasetio, dan		Hidup Pasien	satu solusi utama untuk kanker paru-	dengan mengkombinasikan teknik
	Sari Susanti		Kanker Paru Pasca	paru. Akan tetapi,terdapat banyak	boosting Adaboost sebagai optimasi
			Operasi Bedah	resiko dan komplikasi pasca operasi	level algoritma KNN (K-Nearest
			Toraks	bedah toraks hingga berujung pada	<i>Neighbor)</i> dengan didapatkan akurasi
			Menggunakan	kematian. Maka prediksi ini dilakukan	85.11% dengan menggunakan
			Boosted KNN (K-	dengan menganalisis kondisi pasien	validasi 10 fold cross validation
			Nearest Neighbor)	sebelum dan sesudah operasi. Adaptive	dengan parameter nilai k pada
				boost digunakan sebagai optimasi	algoritma K-NN bernilai 5 untuk

No.	Nama	Tahun	Judul	Isi Ringkasan	Hasil
	Pengarang				
				level algoritma pada algoritma KNN	prediksi harapan hidup pasien pasca
				(K-Nearest Neighbor).	operasi bedah toraks.
9.	Ikhsan	2020	Implementasi Data	Penyakit ginjal kronis (PGK)	Hasil pemodelan data mining terbaik
	Wisnuadji		Mining Untuk	merupakan masalah kesehatan yang	dari sistem dibuat menggunakan
	Gamadarenda,		Deteksi Penyakit	dihadapi oleh dunia dengan angka	Backward Elimination ($\alpha = 0.05$) dan
	dan Indra		Ginjal Kronis (Pgk)	kejadian yang terus meningkat. Deteksi	kNN (k = 3) dengan menghitung
	Waspada		Menggunakan KNN	PGK membutuhkan banyak atribut,	kenaikan biaya inspeksi dan
			(K-Nearest	sehingga membutuhkan biaya yang	sensitivitas tertinggi. Rekomendasi
			Neighbor) Dengan	cukup mahal. Penelitian ini	sistem menghasilkan 10 atribut
			Backward	menggunakan metode data mining	terpilih dari 24 atribut awal yang
			Elimination	dengan mengimplementasikan	digunakan yaitu : berat jenis (sg),
				algoritma KNN (K-Nearest Neighbor)	albumin (al), ureum darah (bu),
				dengan menambahkan algoritma	kreatinin serum (sc), natrium (tanah),
				backward elimination pada tahap	hemoglobin (hemo), sel darah merah
				preprocessing sehingga lebih optimal	(rbc), hipertensi (htn), diabetes
				dan dapat menekan biaya pemeriksaan	mellitus (dm), dan nafsu makan
				laboratorium.	(nafsu makan). Penggunaan atribut

No.	Nama	Tahun	Judul	Isi Ringkasan	Hasil
	Pengarang				
					terpilih berhasil mencapai biaya
					pemeriksaan hingga 73,36%.
					Selanjutnya pendeteksian penyakit
					menggunakan algoritma KNN (K-
					Nearest Neighbor) menghasilkan
					nilai akurasi sebesar 99,25%,
					sensitivitas 99,5%, dan spesifisitas
					98,745%.
10	Irene Lishania,	2019	Perbandingan	Penelitian ini akan membandingkan	Hasil ketepatan klasifikasi penyakit
	Rito		Klasifikasi Metode	hasil akurasi klasifikasi yaitu naïve	stroke pada data pasien di RSUD
	Geojantoro,		Naïve Bayes dan	bayes dan decision tree (J48) pada	Abdul Wahab Sjahranie bulan
	dan Yuki		Metode Decision	pasien stroke. Artinya, seseorang yang	November dan Desember 2017
	Novia		Tree Algoritma	mengalami stroke akan	dengan metode naive bayes adalah
	Nasution		(J48) pada Pasien	diklasifikasikan dengan menggunakan	81,25% dan metode decision tree
			Penderita Penyakit	data pasien di RSUD Abdul Wahab	algoritma (J48) diperoleh tingkat
			Stroke di RSUD	Sjahranie Samarinda dengan 7 faktor	akurasi sebesar 87,5%. Hal ini
			Abdul Wahab	yaitu usia, jenis kelamin, tekanan	menunjukkan bahwa pada penelitian

No.	Nama	Tahun	Judul	Isi Ringkasan	Hasil
	Pengarang				
			Sjahranie	darah, diabetes melitus, dislipidemia,	ini, metode decision tree algoritma
			Samarinda	kadar asam urat dan penyakit jantung.	(J48) memberikan ketepatan prediksi
					klasifikasi yang lebih baik.

Pada *state of the art* tersebut belum ada penelitian yang membandingkan dua algoritma klasifikasi *k-nearst neighbor* dan *adaptive boosting* pada kasus prediksi penyakit stroke. Algoritma *k-nearst neighbor* dan *adaptive boosting* memiliki kekurangan dan kelebihan masing-masing, oleh karena itu penelitian ini akan membandingkan kedua algoritma tersebut untuk meperoleh hasil algoritma yang terbaik dan dengan menerapkan teknik *SMOTE* untuk mengatasi *imbalance data* pada data prediksi penyakit stroke ini.

2.2.2 Matriks Penelitian

Tabel 2. 3 Matriks Penelitian

	Penulis/ Tahun	Judul	Ruang Lingkup							
No.			Algoritma/Metode				Tujuan			
			Naïve Bayes	KNN	C.45	Ada Boost	Klasifik asi	Mengatasi <i>Imbalance</i> Data	Prediksi	
1.	(Badriyah,dkk,	Machine Learning Algorithm for								
	2020)	Stroke Disease	$\sqrt{}$	√	√	_	-	$\sqrt{}$	-	
		Classification								
2.	(Tulus dan	Comparison of Distance Models on								
	Sihombing, 2021)	KNN (K-Nearest Neighbor) Algorithm in Stroke Disease Detection.	-	V	-	-	-	√	-	

	Penulis/ Tahun	Judul	Ruang Lingkup							
No.			Algoritma/Metode				Tujuan			
			Naïve Bayes	KNN	C.45	Ada Boost	Klasifik asi	Mengatasi <i>Imbalance</i> Data	Prediksi	
3.	(Novianti,dkk,	Penerapan Algoritma Adaboost								
	2022)	Untuk Peningkatan Kinerja	_	1	_	1	_	$\sqrt{}$	_	
		Klasifikasi <i>Data Mining</i> Pada		,		,		·		
		Imbalance Dataset Diabetes								
4.	(Sumiah dan	Perbandingan Metode K-Nearest								
	Mirantika, 2020)	Neighbor dan Naive Bayes untuk								
		Rekomendasi Penentuan Mahasiswa	\checkmark	√	-	-	-	\checkmark	-	
		Penerima Beasiswa pada Universitas								
		Kuningan								
5.	(Prasetio dan	Prediksi Harapan Hidup Pasien	-	V	_	√	-	V	-	
	Susanti, 2019)	Kanker Paru Pasca Operasi Bedah								

No.	Penulis/ Tahun	Judul	Ruang Lingkup							
			Alg	oritma	/Meto	de	Tujuan			
			Naïve Bayes	KNN	C.45	Ada Boost	Klasifik asi	Mengatasi <i>Imbalance</i> Data	Prediksi	
		Toraks Menggunakan Boosted KNN								
		(K-Nearest Neighbor)								
6.	(Gamadarendra	Implementasi Data Mining Untuk								
	dan Waspada,	Deteksi Penyakit Ginjal Kronis (Pgk)								
	2020)	Menggunakan K-Nearest Neighbor	-	√	-	-	$\sqrt{}$	$\sqrt{}$	-	
		(K-NN) Dengan Backward								
		Elimination.								
7.	(Nopi, 2022)	Penerapan Metode Adaboost Untuk								
		Optimalisasi Kerja Algoritma K-	_	$\sqrt{}$	-	√	-	$\sqrt{}$	V	
		Nearest Neighbor Untuk Prediksi	_							
		Penyakit Stroke								