

BAB II

LANDASAN TEORI

2.1 *Data Mining*

2.1.1 Definisi *Data Mining*

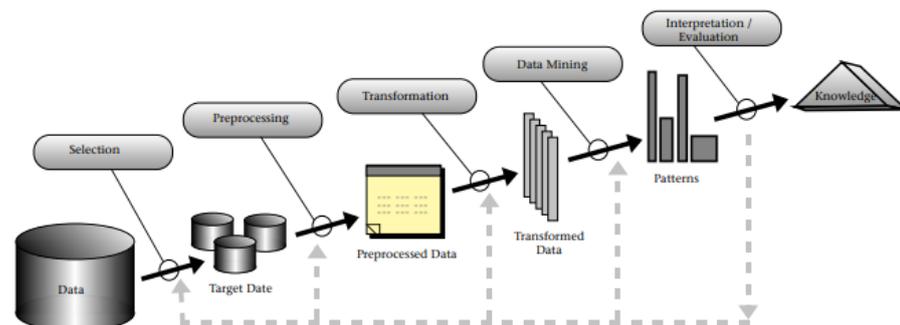
Data mining adalah proses mencari pola atau informasi menarik dalam data terpilih dengan menggunakan teknik atau metode tertentu. Teknik-teknik, metode-metode, atau algoritma dalam *data mining* sangat bervariasi. Pemilihan metode atau algoritma yang tepat sangat bergantung pada tujuan dan proses *Knowledge Discovery in Database* (KDD) secara keseluruhan. Menurut Gartner Group, *data mining* adalah proses menemukan hubungan baru yang mempunyai arti, pola dan kebiasaan dengan memilah-milah sebagian besar data yang disimpan dalam media penyimpanan dengan menggunakan teknologi pengenalan pola seperti teknik statistik dan matematika. *Data mining* merupakan gabungan dari beberapa disiplin ilmu yang menyatukan teknik dari pembelajaran mesin, pengenalan pola, statistik, database, dan visualisasi untuk penanganan permasalahan pengambilan informasi dari database yang besar. Berdasarkan fungsinya, *data mining* dikelompokkan menjadi 6 yaitu deskripsi, estimasi, prediksi, klasifikasi, *clustering* dan asosiasi (Marcoulides, 2005).

Sementara definisi lain dari *Data Mining* adalah proses yang menggunakan teknik statistik, matematika, kecerdasan buatan, dan *machine learning* untuk mengekstraksi dan mengidentifikasi informasi yang bermanfaat dan pengetahuan yang terkait dari berbagai database besar. *Data mining*

merupakan serangkaian proses untuk menggali nilai tambah dari suatu kumpulan data berupa pengetahuan yang selama ini tidak diketahui secara manual (O'Regan, 2011).

2.1.2 Tahapan *Data Mining*

Tahapan yang dilakukan pada proses *data mining* diawali dari seleksi data dari data sumber ke data target, tahap *preprocessing* untuk memperbaiki kualitas data, transformasi, *data mining* serta tahap interpretasi dan evaluasi yang menghasilkan output berupa pengetahuan baru yang diharapkan memberikan kontribusi yang lebih baik (Fayyad, Piatetsky-Shapiro, & Smyth, 2015). Secara detail dijelaskan pada gambar 2.1.



2.1 Gambar 2.1 Tahapan *Data Mining* (Fayyad et al., 2015)

1. *Data Selection*

Pemilihan (seleksi) data dari sekumpulan data operasional perlu dilakukan sebelum tahap penggalian informasi dalam KDD dimulai. Data hasil seleksi yang digunakan untuk proses *data mining*, disimpan dalam suatu berkas, terpisah dari basis data operasional.

2. *Pre-processing / Cleaning*

Sebelum proses *data mining* dapat dilaksanakan, perlu dilakukan proses *cleaning* pada data yang menjadi fokus KDD. Proses *cleaning* mencakup antara lain membuang duplikasi data, memeriksa data yang inkonsisten, dan memperbaiki kesalahan pada data.

3. *Transformation*

Coding adalah proses transformasi pada data yang telah dipilih, sehingga data tersebut sesuai untuk proses *data mining*. Proses coding dalam KDD merupakan proses kreatif dan sangat tergantung pada jenis atau pola informasi yang akan dicari dalam basis data.

4. *Data Mining*

Data mining adalah proses mencari pola atau informasi menarik dalam data terpilih dengan menggunakan teknik atau metode tertentu. Teknik, metode, atau algoritma dalam data mining sangat bervariasi. Pemilihan metode atau algoritma yang tepat sangat bergantung pada tujuan dan proses KDD secara keseluruhan.

5. *Interpretation / Evaluation*

Pola informasi yang dihasilkan dari proses *data mining* perlu ditampilkan dalam bentuk yang mudah dimengerti oleh pihak yang berkepentingan. Tahap ini merupakan bagian dari proses KDD yang disebut *interpretation*. Tahap ini mencakup pemeriksaan apakah pola atau informasi yang ditemukan bertentangan dengan fakta atau hipotesis yang ada sebelumnya.

2.2 *Machine Learning*

Machine Learning adalah bidang ilmu komputer yang melibatkan studi dan konstruksi teknik yang memungkinkan komputer untuk belajar mandiri berdasarkan data input untuk memecahkan masalah spesifik (Hoang & Nguyen, 2018). Jenis-jenis permasalahan yang umumnya diselesaikan dengan pendekatan *Machine Learning* adalah klasterisasi dan klasifikasi. Klasterisasi adalah aktivitas yang bertujuan mengelompokkan data berdasarkan kedekatan fitur yang dimilikinya, sedangkan klasifikasi bertujuan untuk memisahkan data menjadi kelas-kelas tertentu. Perbedaan yang mendasar antara 2 buah permasalahan ini adalah pada proses klasterisasi. Data-data dikelompokkan tanpa pelabelan, sedangkan klasifikasi mengelompokkan data-data menjadi label tertentu (Utama, 2018).

2.3 **Klasifikasi**

Klasifikasi (taksonomi) adalah proses menempatkan objek tertentu (konsep) dalam satu set kategori, berdasarkan masing-masing objek (konsep) *property*. Proses klasifikasi didasarkan pada empat komponen mendasar yaitu kelas, prediktor, *training set*, dan pengujian dataset. Di antara model klasifikasi yang paling populer adalah *Decision / Classification Trees*, *Bayesian Classifiers / Naïve Bayes Classifiers*, *Neural Networks*, *Statistical Analysis*, *Genetic Algorithms*, *Rough Sets*, *K-Nearest Neighbor Classifier*, *Rule based Methods*, *Memory Based Reasoning*, *Support Vector*. Klasifikasi adalah proses untuk menemukan model atau fungsi yang menjelaskan atau membedakan konsep atau kelas data, dengan tujuan untuk dapat memperkirakan kelas dari suatu objek yang

labelnya tidak diketahui. Dalam mencapai tujuan tersebut, proses klasifikasi membentuk suatu model yang mampu membedakan data ke dalam kelas-kelas yang berbeda berdasarkan aturan atau fungsi tertentu. Model itu sendiri bisa berupa aturan “jika maka”, berupa pohon keputusan, atau formula matematis (Fluorida Fibrianda & Bhawiyuga, 2018). Ada juga pendapat lain tentang klasifikasi yaitu merupakan suatu proses menemukan kumpulan pola atau fungsi yang mendeskripsikan serta memisahkan kelas data yang satu dengan yang lainnya untuk menyatakan objek tersebut masuk pada kategori tertentu yang sudah ditentukan atau dikenal juga sebagai *supervised learning*.

Metode klasifikasi terdiri dari dua proses, yaitu : *learning step* atau tahap *training phase* dimana algoritma klasifikasi membangun *classifier* dengan menganalisis atau “belajar dari” sebuah *training set* yang memiliki label dan telah tersedia sebelumnya, selanjutnya untuk mengetahui akurasi dari *classifier* yang telah dibentuk, pada tahap kedua dilakukan pengujian terhadap *classifier* tersebut dengan menggunakan *test set* yang merupakan kumpulan data baru yang dipilih secara acak dan bersifat independen dari *training set* yang artinya data yang digunakan pada *test set* tidak digunakan untuk membangun *classifier* (Zhao, 2014)

2.4 Naive Bayes

Teorema bayes adalah perhitungan statistik dengan menghitung probabilitas kemiripan kasus lama yang ada dibasis kasus dengan kasus baru. *Teorema bayes* memiliki tingkat akurasi yang tinggi dan kecepatan yang baik ketika diterapkan pada *database* yang besar. *Naive bayes* merupakan perhitungan

teorema bayes yang paling sederhana, karena mampu mengurangi kompleksitas komputasi menjadi multiplikasi sederhana dari probabilitas. Selain itu, algoritma *naive bayes* juga mampu menangani set data yang memiliki banyak atribut.

Keuntungan menggunakan *naive bayes* adalah metode ini hanya membutuhkan dua data yaitu data *training (training set)* dan data *testing (testing set)* untuk menguji suatu data yang ingin diperoleh (Jadhav & Channe, 2016).

Persamaan dari *naive bayes* sebagai berikut :

$P(C | X)$, dari $P(C)$, $P(X)$, dan $P(X | C)$

$$P(C | X) = \frac{P(X | C) \cdot P(C)}{P(X)}$$

Gambar 2.2 Persamaan Klasifikasi *Naive Bayes* (Jadhav & Channe, 2016)

Dimana :

$P(C | X)$ adalah probabilitas posterior dari kelas target.

$P(C)$ disebut probabilitas sebelumnya dari kelas.

$P(X | C)$ adalah kemungkinan yang merupakan probabilitas prediktor kelas yang diberikan.

$P(X)$ adalah probabilitas sebelumnya dari prediktor kelas.

2.5 *Support Vector Machine (SVM)*

Support Vector Machine (SVM) merupakan salah satu metode klasifikasi. Terdapat dua kategori SVM yaitu *Support Vector Machine Classification* dan *Support Vector Machine Regression*. SVM di perkenalkan pertama kali oleh Vapnik pada tahun 1992 sebagai konsep unggulan dalam bidang *pattern recognition*, algoritma ini dapat memilih model otomatis dan tidak memiliki

masalah *overfittin*. Metode SVM sangat baik untuk prediksi karena metode ini dapat meminimalkan kesalahan klasifikasi dan penyimpangan data pada data training. Cara kerja metode SVM yaitu dengan memisahkan beberapa kelompok data dengan garis. Garis ini dikenal dengan *hyperplane*, dengan teknik SVM bertujuan untuk mencari *hyperplane* yang optimal. Kernel merupakan fungsi yang digunakan untuk mendapatkan *hyperplane* yang optimum (Fadilah, Agfiannisa, & Azhar, 2020).

2.6 *Confusion Matrix*

Confusion matrix dapat diartikan sebagai suatu alat yang memiliki fungsi untuk melakukan analisis apakah *classifier* tersebut baik dalam mengenali tuple dari kelas yang berbeda. Nilai dari *True-Positive* dan *True-Negative* memberikan informasi ketika *classifier* dalam melakukan klasifikasi data bernilai benar, sedangkan *False-Positive* dan *False-Negative* memberikan informasi ketika *classifier* salah dalam melakukan klasifikasi data. TP (*True Positive*) merupakan jumlah data dengan nilai sebenarnya positif dan nilai prediksi positif, FP (*False Positive*) merupakan jumlah data dengan nilai sebenarnya negatif dan nilai prediksi positif, FN (*False Negative*) merupakan jumlah data dengan nilai sebenarnya positif dan nilai prediksi negatif, TN (*True Negative*) merupakan jumlah data dengan nilai sebenarnya negatif dan nilai prediksi negatif (Fibrianda & Bhawiyuga, 2018).

2.7 *Parameter Matric*

Accuracy merupakan tingkat keterhubungan antara suatu nilai yang diprediksi dengan nilai aktual yang ada (Devita et al., 2018).

$$Accuracy = \frac{True\ Positive + True\ Negative}{True\ Positive + False\ Negative + False\ Positive + True\ Negative} \text{ (Rumus 1)}$$

Precision merupakan pengukuran tingkat ketepatan antara informasi yang diminta oleh pemohon dengan jawaban yang diberikan oleh sistem.

$$Precision = \frac{True\ Positive}{True\ Positive + False\ Positive} \text{ (Rumus 2)}$$

Recall merupakan tingkat keberhasilan sistem dalam menemukan kembali sebuah informasi dalam suatu pemrosesan data.

$$Recall = \frac{True\ Positive}{True\ Positive + False\ Negative} \text{ (Rumus 3)}$$

2.8 WEKA (*Waikato Environment for Knowledge Analysis*)

WEKA merupakan sebuah *Tools* yang praktis. Kepanjangan dari WEKA adalah *Waikato Environment for Knowledge Analysis*, dibuat di Universitas Waikato, *New Zealand*. WEKA juga dapat digunakan pada beberapa tingkatan berbeda serta mengandung *Tools* untuk *pre-processing* data yaitu : *classification*, *regression*, *clustering*, *association rules*, dan *visualization* (Pujiono, Amborowati, Suyanto, & Kunci, 2013). Kelas yang paling penting disini adalah *classifier*, yang mendeklarasikan struktur umum dari skema klasifikasi dan prediksi.

Kelas ini memiliki 2 metoda, yaitu *build classifier* dan *classify instance*, yang harus diimplementasikan oleh kelas-kelas yang menginduk kelas ini. WEKA sangat membantu dalam mengolah data (*Data Mining*), karena dengan software WEKA ini dapat dengan mudah untuk mengolah setumpuk data dan mengambil informasi yang penting saja dari sekumpulan data tersebut. Selain itu juga, WEKA mampu menyelesaikan masalah-masalah *data mining* didunia nyata khususnya klasifikasi yang mendasari pendekatan *machine learning*.

2.9 State Of The Art (SOTA)

State Of The Art merupakan pengumpulan data penelitian sebelumnya dengan menentukan sumber atau referensi yang tertera pada tabel 2.1

Tabel 2.1 Referensi Penelitian

No	Peneliti	Tahun Terbit	Judul	Metode/Algoritma	Kesimpulan
1	Devi Nurul Anisa dan Jumanto (2022)	Dinamika Informatika, Vol.14, No.1, Maret 2022 : 33-42, E-ISSN : 2714-8769 P-ISSN : 2085-3343	Klasifikasi Penyakit Diabetes Menggunakan Algoritma <i>Naive Bayes</i>	<i>Naive Bayes</i>	Menerapkan metode klasifikasi untuk memprediksi apakah seseorang terkena diabetes atau tidak menggunakan algoritma <i>Naive Bayes</i> . Terdapat 390 data yang diklasifikasikan, dalam pengklasifikasiannya menggunakan 9 variabel yaitu <i>cholesterol, glucose, hdl cholesterol, age, gender, weight, systolic bp, distolic bp, diabetes</i> . Menghasilkan akurasi 92,3% untuk <i>Data Train</i> dan 91,6% untuk <i>Data Test</i> .
2	Hilda Apriyani dan Kurniati (2020)	Journal of Information Technology Ampera Vol. 1, No. 3, December 2020 e-ISSN: 2774-2121	Perbandingan Metode <i>Naive Bayes</i> Dan <i>Support Vector Machine</i> Dalam Klasifikasi Penyakit Diabetes Melitus	<i>Naive Bayes</i> dan <i>Support Vector Machine (SVM)</i>	Pengujian yang dilakukan menggunakan data training dan data testing kemudian diukur untuk mengetahui tingkat akurasi dengan evaluasi <i>confusion matrix</i> . Untuk algoritma SVM menggunakan dua kernel yaitu kernel polynomial dan RBF. Kernel Polynomial memiliki nilai akurasi yang tertinggi dengan nilai 96.2704% bisa dikatakan lebih akurat bila dibandingkan dengan algoritma <i>naive bayes</i> dengan tingkat akurasi 92.0746%.

Tabel 2.2 Referensi Penelitian (lanjutan)

3	Mercury Fluorida Fibrianda dan Adhitya Bhawiyuga (2018)	Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer, Vol. 2, No. 9, September 2018, hlm. 3112-3123, e-ISSN: 2548-964X	Analisis Perbandingan Akurasi Deteksi Serangan Pada Jaringan Komputer Dengan Metode <i>Naïve Bayes</i> Dan <i>Support Vector Machine</i> (SVM)	<i>Naive Bayes</i> , SVM <i>Linear</i> , SVM <i>Polynomial</i> , dan SVM <i>Sigmoid</i> .	Fitur yang digunakan dalam proses klasifikasi yaitu <i>totalSourceBytes</i> , <i>totalDestinationBytes</i> , <i>totalDestinationPacket</i> , <i>totalSourcePacket</i> , <i>direction</i> , <i>Source TCPFlagsDescription</i> , <i>Destination TCPFlagsDescription</i> , <i>protocolName</i> , <i>sourcePort</i> , <i>Destination</i> , <i>destinationPort</i> , <i>startDateTime</i> , dan <i>stopDateTime</i> . Performa yang dihasilkan dari <i>confusion matrix</i> pada masing-masing <i>classifier Naive Bayes</i> , SVM <i>Linear</i> , SVM <i>Polynomial</i> , dan SVM <i>Sigmoid</i> menghasilkan persentase akurasi berturut-turut sebesar 85,055%, 99,995%, 99,999% dan 99,995%. Performa kinerja klasifikasi yang dihasilkan dari kurva ROC pada <i>classifier Naive Bayes</i> yaitu baik, SVM <i>Linear</i> lemah, SVM <i>Polynomial</i> sangat lemah, dan SVM <i>Sigmoid</i> lemah. Sedangkan jika dilihat dari kurva ROC dengan <i>cross-validation</i> menunjukkan bahwa <i>classifier Naive Bayes</i> yaitu lemah dengan nilai AUC 0,5, SVM <i>Linear</i> baik dengan nilai AUC 0,75, SVM <i>Polynomial</i> sangat lemah dengan nilai AUC 0,33 dan SVM <i>Sigmoid</i> lemah dengan nilai AUC 0,5.
---	---	---	--	---	---

Tabel 2.3 Referensi Penelitian (lanjutan)

4	Hairani, Khurniawan Eko Saputro dan Sofiansyah Fadli (2020)	Jurnal Teknologi dan Sistem Komputer, 8(2), 2020, 89-93, DOI:10.14710/jtsiskom.8.2.2020.89-93	<i>K-Means-SMOTE</i> untuk menangani ketidakseimbangan kelas dalam klasifikasi penyakit diabetes dengan C4.5, SVM dan <i>Naive Bayes</i>	C4.5, SVM dan <i>Naive Bayes</i>	Penelitian ini dilakukan untuk menyelesaikan permasalahan ketidakseimbangan kelas pada dataset penyakit diabetes Pima Indian menggunakan <i>K-Means-SMOTE</i> . Dataset tersebut memiliki 268 data dari kelas positif (kelas <i>minoritas</i>) dan 500 data dari kelas negatif (kelas <i>mayoritas</i>). Kombinasi <i>K-Means-SMOTE</i> dengan metode klasifikasi SVM memiliki akurasi dan sensitivitas terbaik, yaitu sebesar 82% dan 77%. Sedangkan dengan metode <i>Naive Bayes</i> menghasilkan spesifisitas terbaik sebesar 89%.
5	Abu Wildan Mucholladin, Fitra Abdurrachman Bachtiar dan Muhammad Tanzil Furqon (2021)	Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer, Vol. 5, No. 2, Februari 2021, hlm. 622-633, e-ISSN: 2548-964X	Klasifikasi Penyakit Diabetes menggunakan Metode <i>Support Vector Machine</i>	<i>Support Vector Machine (SVM)</i>	Model SVM dibagi menjadi dua jenis yaitu model <i>benchmark</i> yang diimplementasikan menggunakan algoritma <i>Sequential Minimal Optimization (SMO)</i> dan model <i>scratch</i> yang diimplementasikan menggunakan algoritma <i>Sequential Learning</i> . Model yang sudah optimal diuji kembali pada beberapa metrik menggunakan <i>10-fold cross validation</i> . Hasil pengujian menunjukkan bahwa model <i>benchmark</i> memiliki nilai <i>mean accuracy</i> sebesar 0,87, <i>mean precision</i> sebesar 0,82, <i>mean sensitivity</i> sebesar 0,78, dan <i>mean specificity</i> sebesar 0,92. Model <i>scratch</i> memiliki nilai <i>mean accuracy</i> sebesar 0,78, <i>mean precision</i> sebesar 0,69, <i>mean sensitivity</i> sebesar 0,59, dan <i>mean specificity</i> sebesar 0,87

Tabel 2.4 Referensi Penelitian (lanjutan)

6	Erwin, Laras Azrisa Nurjanah, Dea Sellan Noviyanti dan Yurika (2018)	Prosiding Annual Research Seminar 2018 - Computer Science and ICT, Vol.4 No.1, ISBN : 978-979-587-813-1	Klasifikasi Penyakit <i>Diabetik Retinopathy</i> dengan Metode <i>Naive Bayes</i> pada Citra Retina	<i>Naive Bayes</i>	Dari data uji sebanyak 15 citra retina didapatkan hasil klasifikasi 5 citra Normal, 6 citra NPDR, dan 3 citra PDR. Dengan metode <i>Naive Bayes</i> mendapatkan akurasi sebesar 93%. Metode <i>Naive Bayes</i> yang digunakan mendapatkan hasil yang baik dalam pengklasifikasian.
7	Sayali D. Jadhav dan H. P. Channe (2013)	International Journal of Science and Research (IJSR) ISSN (Online): 2319-7064 Index Copernicus Value (2013): 6.14 Impact Factor (2014): 5.611	<i>Comparative Study of K-NN, Naive Bayes and Decision Tree Classification Techniques</i>	K-NN, <i>Naive Bayes</i> dan <i>Decision Tree</i>	Semua algoritma <i>Decision Tree</i> lebih akurat dan memiliki tingkat kesalahan yang lebih sedikit dan merupakan algoritma yang lebih mudah dibandingkan dengan <i>K-NN</i> dan <i>Bayesian</i> . Hasil dari implementasi WEKA pada dataset yang sama menunjukkan bahwa algoritma <i>Decision Tree</i> lebih unggul dan algoritma <i>Naive Bayes</i> memiliki tingkat akurasi yang hampir sama dengan algoritma <i>Decision Tree</i> .
8	Achmad Ridwan (2020)	Jurnal Sistem Komputer dan Kecerdasan Buatan Volume IV – Nomor 1 – September 2020	Penerapan Algoritma <i>Naive Bayes</i> Untuk Klasifikasi Penyakit Diabetes Mellitus	<i>Naive Bayes</i>	Data yang digunakan pada analisis ini merupakan data dari dataset <i>UCI Machine Learning</i> yaitu <i>Early Stage Diabetes Risk</i> tahun 2020 dan terdapat 16 atribut yang mempengaruhi Klasifikasi dataset <i>Early Stage Diabetes Risk prediction</i> yaitu <i>Age, Gender, Polyuria, Polydipsia, Sudden Weight Loss, Weakness, Polyphagia, Genital Thrush, Visual Blurring, Itching, Irritability, Delayed Healing, Partial Paresis, Muscle Stiffness, Alopecia, Obesity, Class</i> . Analisis yang dilakukan meliputi data <i>preprocessing</i> , model, dan evaluasi Hasil klasifikasi menunjukkan akurasi sebesar 90.20% dan nilai AUC nya yaitu 0,95.

Tabel 2.5 Referensi Penelitian (lanjutan)

9	Fatmawati (2016)	Jurnal Techno Nusa Mandiri Vol. XIII, No. 1 Maret 2016 50	Perbandingan Algoritma Klasifikasi <i>Data Mining</i> Model C4.5 Dan <i>Naive Bayes</i> Untuk Prediksi Penyakit Diabetes	C4.5 dan <i>Naive Bayes</i> .	Hasil klasifikasi data di evaluasi dengan menggunakan <i>Confusion Matrix</i> dan kurva <i>ROC</i> untuk mengetahui tingkat hasil akurasi menggunakan algoritma <i>Decision Tree</i> yaitu sebesar 73.30% dan nilai AUC dari kurva <i>ROC</i> adalah 0.733 sedangkan algoritma <i>Naive Bayes</i> sebesar 75.13% nilai AUC dari kurva <i>ROC</i> 0.810 sehingga dapat disimpulkan bahwa dengan menggunakan model <i>Naive Bayes</i> lebih tinggi tingkat akurasi, dengan peningkatan akurasi sebesar 1.83% dan peningkatan nilai AUC sebesar 0.077 sedangkan hasil pengujian dari prediksi diabetes hasilnya termasuk <i>Good Clasification</i> .
10	Dian Pramadhana (2021)	Jurnal Pendidikan Informatika Vol. 5 No. 1, Juni, 2021, Hal. 89-98, e-ISSN 2549-7472 DOI: 10.29408/edumatic.v5i1.3336	Klasifikasi Penyakit Diabetes Menggunakan Metode <i>CFS</i> dan <i>ROS</i> dengan Algoritma J48 Berbasis <i>Adaboost</i>	J48	Penggunaan seleksi fitur dengan metode <i>CFS (Correlation Feature Selection)</i> dan <i>Resample (Random Over Sampling)</i> yang didapatkan pada algoritma J48 menghasilkan nilai akurasi yang lebih baik yaitu sebesar 92,3%. Dengan penambahan metode <i>CFS</i> dan <i>Resample (Random Over Sampling)</i> serta <i>adaboost</i> dalam pengklasifikasian penyakit diabetes dapat meningkatkan kinerja Algoritma J48 sehingga membuat semakin lebih optimal dibandingkan tanpa menggunakan metode <i>CFS</i> dan <i>Resample (Random Over Sampling)</i> serta <i>adaboost</i> .

Tabel 2.6 Referensi Penelitian (lanjutan)

11	Hiya Nalatissifa, Windu Gata, Sri Diantika, Khoirun Nisa (2020)	Jurnal Informatika Universitas Pamulang, Vol. 5, No. 4, Desember 2020 (578-584), ISSN: 2541-1004, e-ISSN: 2622-4615	Perbandingan Kinerja Algoritma Klasifikasi <i>Naive Bayes</i> , <i>Support Vector Machine (SVM)</i> , dan <i>Random Forest</i> untuk Prediksi Ketidakhadiran di Tempat Kerja	<i>Naive Bayes</i> , <i>Support Vector Machine (SVM)</i> , dan <i>Random Forest</i>	Memprediksi ketidakhadiran di tempat kerja berdasarkan dataset <i>Absenteeism at work</i> menggunakan aplikasi Weka 3.8 dan algoritma <i>Naive Bayes</i> , <i>Support Vector Machine (SVM)</i> , dan <i>Random Forest</i> . Pada hasil penelitian, algoritma <i>Random Forest</i> memperoleh nilai akurasi, presisi, dan <i>recall</i> yang paling tinggi dibandingkan dengan algoritma <i>Naive Bayes</i> dan <i>SVM</i> , yaitu menghasilkan nilai akurasi sebesar 99.38%, presisi 99.42% dan <i>recall</i> 99.39%.
12	Yusra, Dhita Olivita, Yelfi Vitriani (2016)	Jurnal Sains, Teknologi dan Industri, Vol. 14, No. 1, Desember 2016, pp. 79 - 85 ISSN 1693-2390 print/ISSN 2407-0939 online	Perbandingan Klasifikasi Tugas Akhir Mahasiswa Jurusan Teknik Informatika Menggunakan Metode <i>Naive Bayes Classifier</i> dan <i>K-Nearest Neighbor</i>	<i>Naive Bayes Classifier</i> dan <i>K-Nearest Neighbor</i>	Pengujian akurasi metode pada penelitian ini dilakukan dengan <i>test option 10-fold cross validation</i> dan evaluasi data uji menggunakan <i>confusion matrix</i> . Dari penelitian yang telah dilakukan, didapatkan hasil pada seratus data tugas akhir dengan jumlah kelas acak, metode <i>Naive Bayes</i> menghasilkan nilai akurasi lebih baik, yaitu sebesar 87%. Pengujian pada metode <i>K-Nearest Neighbor</i> menghasilkan nilai akurasi 84% dengan nilai $k=3$, 85% dengan nilai $k=5$, 86% dengan nilai $k=7$ dan 84% dengan nilai $k=9$
13	Sri Diantika, Windu Gata, Hiya Nalatissifa, Mareanus Lase (2021)	Jurnal Ilmiah Elektronika dan Komputer, Vol.14, No.1, Juli 2021, pp. 10 – 15, p-ISSN : 1907-0012, e-ISSN : 2714-5417	Komparasi Algoritma SVM Dan <i>Naive Bayes</i> Untuk Klasifikasi Kestabilan Jaringan Listrik	<i>Support Vector Machine (SVM)</i> dan <i>Naive Bayes</i>	Dilakukan perbandingan penerapan algoritma klasifikasi SVM dan <i>Naive Bayes</i> terhadap dataset <i>Electrical Grid Stability Simulated</i> yang diambil dari UCI Machine Learning. Dari hasil pengujian klasifikasi kestabilan jaringan listrik yang telah dilakukan menggunakan aplikasi WEKA 3.8.2. Metode <i>Support Vector Machine (SVM)</i> menunjukkan tingkat <i>accuracy</i> yang lebih baik yaitu sebesar 98.9% jika

					dibandingkan dengan metode <i>Naive Bayes</i> yang menghasilkan nilai akurasi sebesar 97.64% Hasil akurasi ini akan menunjukkan hasil yang berbeda tergantung dengan jenis data, jumlah instance, label class dan Percentage split data yang digunakan.
--	--	--	--	--	---

2.10 Matriks Penelitian

Matriks penelitian menjelaskan tentang perbedaan antara penelitian yang dilakukan dengan penelitian terkait. Terdapat beberapa indikator yang menunjukkan perbedaan dan persamaan antara penelitian terdahulu dengan penelitian yang dilakukan.

Berikut merupakan tabel matriks penelitian :

Tabel 2.2 Matriks Penelitian

No	Peneliti	Judul	Algoritma							Parameter Metric						
			Naive Bayes	SVM	C4.5	KNN	Decision Tree	J48	Random Forest	Accuracy	Precision	Recall	Error Rate	F1 Score	Spesifisitas	ROC Curve
1	Devi Nurul, dkk (2022)	Klasifikasi Penyakit Diabetes Menggunakan Algoritma <i>Naive Bayes</i>	✓							✓						
2	Hilda Apriyani, dkk (2020)	Perbandingan Metode <i>Naive Bayes</i> Dan <i>Support Vector Machine</i> Dalam Klasifikasi Penyakit Diabetes Melitus	✓	✓						✓			✓			

Tabel 2.2 Matriks Penelitian (lanjutan)

3	Mercury Fluorida Fibrianda, dkk (2018)	Analisis Perbandingan Akurasi Deteksi Serangan Pada Jaringan Komputer Dengan Metode <i>Naïve Bayes</i> Dan <i>Support Vector Machine</i> (SVM)	✓	✓						✓	✓	✓		✓		
4	Hairani, dkk (2020)	<i>K-Means-SMOTE</i> untuk menangani ketidakseimbangan kelas dalam klasifikasi penyakit diabetes dengan C4.5, SVM dan <i>Naive Bayes</i>	✓	✓	✓					✓		✓				✓
5	Abu Wildan Mucholladin, dkk (2021)	Klasifikasi Penyakit Diabetes menggunakan Metode <i>Support Vector Machine</i>		✓						✓	✓	✓				✓
6	Erwin, dkk (2018)	Klasifikasi Penyakit <i>Diabetik Retinopathy</i> dengan Metode <i>Naive Bayes</i> pada Citra Retina	✓							✓						

Tabel 2.2 Matriks Penelitian (lanjutan)

7	Sayali D. Jadhav, dkk (2013)	<i>Comparative Study of K-NN, Naive Bayes and Decision Tree Classification Techniques</i>	✓			✓	✓			✓					✓
8	Achmad Ridwan (2020)	Penerapan Algoritma <i>Naïve Bayes</i> Untuk Klasifikasi Penyakit Diabetes Mellitus	✓							✓	✓	✓			
9	Fatmawati (2016)	Perbandingan Algoritma Klasifikasi <i>Data Mining</i> Model C4.5 Dan <i>Naive Bayes</i> Untuk Prediksi Penyakit Diabetes	✓		✓		✓			✓					
10	Dian Pramadhana (2021)	Klasifikasi Penyakit Diabetes Menggunakan Metode <i>CFS</i> dan <i>ROS</i> dengan Algoritma J48 Berbasis <i>Adaboost</i>						✓		✓	✓	✓			

Tabel 2.2 Matriks Penelitian (lanjutan)

11	Hiya Nalatissifa, dkk (2020)	Perbandingan Kinerja Algoritma Klasifikasi <i>Naive Bayes</i> , <i>Support Vector Machine (SVM)</i> , dan <i>Random Forest</i> untuk Prediksi Ketidakhadiran di Tempat Kerja	✓	✓						✓	✓	✓	✓				
12	Yusra, Dhita Olivita, dkk (2016)	Perbandingan Klasifikasi Tugas Akhir Mahasiswa Jurusan Teknik Informatika Menggunakan Metode <i>Naive Bayes Classifier</i> dan <i>K-Nearest Neighbor</i>	✓								✓						
13	Sri Diantika, dkk (2021)	Komparasi Algoritma <i>SVM</i> Dan <i>Naive Bayes</i> Untuk Klasifikasi Kestabilan Jaringan Listrik	✓	✓							✓	✓	✓				✓
14	Syifa Nazila Fauziah (2023)	Komparasi Algoritma <i>Naive Bayes</i> Dan <i>Support Vector Machine (SVM)</i> Untuk Klasifikasi Data Pada <i>Diabetes Prediction Dataset</i>	✓	✓							✓	✓	✓				