

## BAB 2

### TINJAUAN PUSTAKA

#### 2.1. Kinerja Mahasiswa

Menurut Naomi & Nindyati (2008) dalam buku Indra et al., (2021) Kinerja akademik merupakan hasil akhir yang dicapai oleh siswa/mahasiswa sebagai keberhasilan selama mengikuti pendidikan dalam sebuah institusi pendidikan (I Made Indra P., S.KM., MPH., QRG.P. et al., 2021).

#### 2.2. *Data Mining*

*Data Mining* adalah suatu proses yang menggunakan teknik statistik, matematika, kecerdasan buatan, dan *machine learning* untuk memperoleh dan mengidentifikasi informasi yang berguna dan pengetahuan yang terkait dari berbagai *database* besar (Turban, dkk. 2005). Berikut merupakan beberapa teknik dari *data mining* antara lain (Osman, 2019):

a. Asosiasi

Asosiasi merupakan salah satu teknik *data mining* yang cukup dikenal dalam menemukan pola berdasarkan hubungan antar variabel dalam transaksi yang sama. Selain itu, asosiasi juga dikenal sebagai teknik relasi karena menggunakan relasi antara item dan menemukan seringnya kemunculan item berbeda yang muncul dengan frekuensi tertinggi dalam dataset.

b. Klasifikasi

Klasifikasi merupakan teknik yang digunakan untuk mengklasifikasikan kumpulan data ke dalam kelompok atau kelas untuk mendapatkan prediksi dan analisis yang akurat dalam kumpulan data yang besar.

c. Klustering

Klustering merupakan salah satu teknik pertama yang digunakan dalam *data mining*. Dalam prosesnya melibatkan satu atau lebih atribut untuk mengidentifikasi antara data yang mirip satu sama lain untuk mengetahui perbedaan dan persamaan antara kumpulan data.

d. Prediksi

Prediksi merupakan topik yang komprehensif, mulai dari memprediksi kegagalan komponen hingga memahami dan memprediksi keuntungan perusahaan. Hal ini digunakan dengan kombinasi teknik *data mining* lainnya, teknik prediksi ini termasuk analisis tren, klasifikasi, pencocokan pola, dan hubungan. Selain itu, prediksi juga dibuat dengan menganalisis peristiwa masa lalu.

### 2.2.1. Algoritma Klasifikasi

Klasifikasi adalah jenis analisis data yang dapat membantu orang memprediksi label kelas sampel harus diklasifikasikan. Berbagai macam teknik klasifikasi telah diusulkan dalam bidang-bidang seperti pembelajaran mesin, sistem pakar dan statistik (Yu et al., 2008). Menurut (Ardiyansyah et al., 2018), Algoritma Klasifikasi merupakan suatu teknik pembelajaran untuk memprediksi dari beberapa atribut yang menggambarkan dan membedakan kelas data atau konsep dengan tujuan memprediksi dari objek kelasnya yang tidak diketahui.

### 2.2.2. Random Forest

*Random Forest* adalah suatu teknik yang dikembangkan dari algoritma *Classification and Regression Tree* (CART), yang pada prosesnya menggunakan pendekatan *bootstrap aggregating (bagging)* dan *random feature selection* (Breiman, 2001). *Random Forest* adalah salah satu algoritma *machine learning* yang merupakan pengembangan dari algoritma *Decision Tree*, *RF* bisa dilihat sebagai gabungan beberapa buah *decision tree* (Primartha, 2021). Rumus dari *RF* yang terdiri dari  $N$  trees dinyatakan sebagai berikut (Liparas et al., 2014):

$$l(y) = \operatorname{argmax}_c \left( \sum_{n=1}^N I_{hn}(y) = c \right) \quad (2.1)$$

Dimana variabel  $l$  adalah fungsi indikator dan  $hn$  merupakan *tree* ke- $n$  dari *RF*.

Metode *CART* yang digunakan untuk membangun pohon pada algoritma *Random Forest Classifier* menggunakan aturan *Gini Impurity* untuk menentukan pecahan dari pohon keputusan (Daniya et al, 2020). Perhitungan dimulai dengan

penentuan nilai *Gini Index* untuk menentukan distribusi probabilitas atribut terhadap kelas target dan dilanjutkan pada perhitungan *Gini Impurity*. Berikut adalah rumus perhitungan ***Gini*** (Daniya et al, 2020):

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

Dimana:

$n$  = Merupakan jumlah kelas target

$l$  = Merupakan kelas target

$p$  = Merupakan rasio kelas target

Berikut ini merupakan rumus perhitungan *Gini Impurity*:

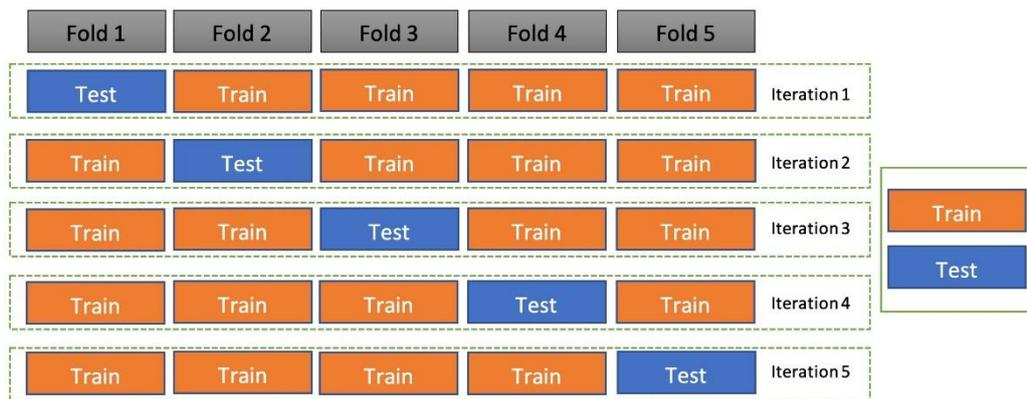
$$Gini\ impurity = 1 - \sum_{i=1}^n P_i^2 \quad (2.3)$$

### **2.2.3. Data Preprocessing**

*Data Preprocessing* merupakan proses mengubah data mentah menjadi data yang lebih mudah dipahami serta mudah diproses. Hal ini diperlukan untuk memperbaiki kesalahan yang ada pada data mentah seperti data tidak lengkap, kosong, tidak sesuai, dan formatnya salah. Tujuan dari *data preprocessing* untuk membuat kualitas data menjadi lebih baik, termasuk kelengkapan, konsistensi, ketepatan waktu, dan meningkatkan hasil dari akurasi (rifqi dharma, 2022). Tahapan-tahapan dari *data preprocessing* yaitu, *data selection*, *data cleaning*, *data transformation*, dan *data reduction* (rifqi dharma, 2022):

### **2.2.4. K-Fold Cross Validation**

*Cross Validation* merupakan metode yang dapat digunakan untuk mengevaluasi kinerja dari suatu model atau algoritma yang dimana memisahkan data menjadi 2 subset yaitu data latih dan data uji (Daqiqil, 2021). Subset pembelajaran melatih model atau algoritma dan subset validasi memvalidasinya lalu pemilihan jenis CV berdasarkan ukuran dari datasetnya. *K-fold* adalah sebuah metode yang memecah dataset menjadi dua bagian yaitu data latih dan data uji sebanyak  $k$  kelompok yang dimana jumlah data latih dan data uji pada tiap kelompok sama.



**Gambar 2. 1 K-Fold Cross Validation**

*Sumber: (Gopal Krishna Ranjan, 2021)*

### 2.2.5. Confusion Matrix

*Confusion Matrix* merupakan matrik yang membandingkan hasil aktual dengan hasil prediksi, terutama digunakan dalam *supervised learning* untuk mengevaluasi akurasi prediksi dari suatu klasifikasi (Pu et al., 2021). Dalam *Confusion Matrix* terdapat istilah TP, TN, FP, dan FN yang antara lain (Gde Agung Brahma Suryanegara et al., 2021):

**Tabel 2. 1 Confusion Matrix**

		Actual Value	
		TRUE	FALSE
Prediction value	TRUE	True Positive (TP)	False Positive (FP)
	FALSE	False Negative (FN)	True Negative (TN)

TP = *True Positive* artinya data positif yang diprediksi secara positif

FP = *False Positive* artinya data negatif yang diprediksi secara positif

FN = *False Negative* artinya data positif yang diprediksi secara negatif

TN = *True Negative* artinya data negatif yang diprediksi secara negatif

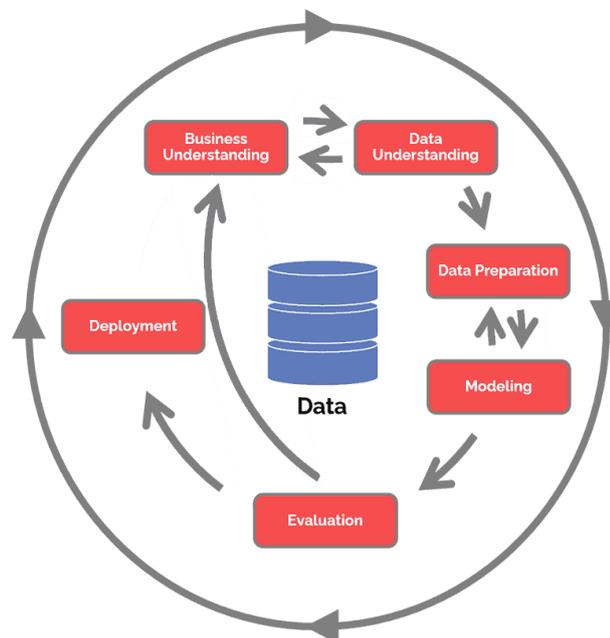
Dari tabel diatas, *confusion matrix* dapat diukur dan dievaluasi tingkat performa model klasifikasinya dengan menghitung akurasi. Akurasi adalah persentase dari data uji yang diklasifikasikan dengan benar oleh model klasifikasi yang dibangun. Akurasi ini dapat dihitung dengan rumus berikut (Gde Agung Brahma Suryanegara et al., 2021):

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \times 100\%$$

(2.4)

### 2.2.6. CRISP-DM

*Cross-Industry Standard Process for Data Mining* atau yang disingkat dengan CRISP-DM merupakan suatu model atau tahapan dalam mengelola dan menyempurnakan *data mining* (Huber et al., 2019). Ada beberapa tahapan yang ada di dalam CRISP-DM diantaranya:



**Gambar 2. 2 CRISP-DM**

Sumber: (Hotz, 2018)

1. *Business Understanding*

Tahap awal dari CRISP-DM adalah menentukan tujuan bisnis lalu menilai situasi saat ini dan setelahnya menetapkan tujuan dilakukannya *data mining*.

2. *Data Understanding*

Tahap ini merupakan persiapan, mengevaluasi persyaratan data, dan mengumpulkan data. Data yang dikumpulkan kemudian dideskripsikan untuk mengetahui mana yang menjadi atribut, kelas, dan jenis tipe datanya.

### 3. *Data Preparation*

Selanjutnya pada tahap ini data-data perlu diidentifikasi, dipilih, dibersihkan dan diubah ke dalam format yang diinginkan. Pada tahapan ini juga bisa disebut sebagai data *preprocessing*.

### 4. *Modeling*

Pada tahap ini dilakukan penerapan algoritma dalam mencari, mengidentifikasi, dan menemukan pola. Pemilihan algoritma yang digunakan menyesuaikan dengan jenis tipe data karena dari tipe datanya, dapat diketahui apakah data tersebut akan diprediksi, diklasifikasi, diklaster atau dengan melihat hubungan asosiasi.

### 5. *Evaluation*

Di tahap ini evaluasi dilakukan pada model atau algoritma yang dilakukan sebelumnya, Pada algoritma klasifikasi, evaluasi yang sering digunakan adalah akurasi, *sensitivity*, *F-Measure*, dan lain-lain.

### 6. *Deployment*

Tahap *Deployment* merupakan tahap akhir dalam CRISP-DM, yang bertujuan untuk melakukan otomatisasi model atau pengembangan aplikasi yang terintegrasi dengan sistem informasi manajemen atau operasional yang ada.

## **2.2.7. Analysis of Variance (ANOVA)**

*Analysis of Variance* merupakan teknik standar untuk mengukur signifikansi statistik dari suatu set variabel independen dalam memprediksi variabel dependen (Wenda, 2022). *ANOVA* digunakan untuk memangkas fitur tanpa memengaruhi keakuratan prediktor (Chen et al., 2022). Dalam penyelesaian menggunakan *ANOVA*, berikut langkah-langkahnya:

- a. Semua fitur dipilih dari dataset.
- b. Fungsi fitur target dari *scikit-learn* dihitung menggunakan *ANOVA F-Score* untuk setiap fitur. Dibawah ini merupakan rumus untuk menghitung *ANOVA*.

$$F = \frac{\text{variance between groups}}{\text{variance within groups}}$$

$$\text{Variance between groups} = \frac{\sum_i^n n_i (\bar{Y}_i - \bar{Y})^2}{(k - 1)}$$

$$\text{Variance within groups} = \frac{\sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2}{(n - k)}$$

(2.5)

- c. Hasil dari pengujian digunakan untuk melakukan pemilihan fitur yang memungkinkan pembuangan fitur yang tidak terkait dengan variabel target. Fitur yang memiliki pengaruh paling tinggi dengan varian terendah dipilih dalam eksperimen ini dan diuji dengan *SelectKBest()*; *K* mewakili jumlah fitur yang ada untuk dataset akhir.
- d. Jumlah fitur dengan peringkat tertinggi digunakan untuk membuat berbagai subset fitur.

### 2.3. Penelitian Terdahulu

Tabel 2. 2 Penelitian Terdahulu

No.	Nama Peneliti	Judul Penelitian	Hasil Penelitian
1	(Hendrawan et al., 2021)	Implementasi Pemilihan Fitur Metode <i>Wrapper</i> dan <i>Embedded</i> dalam Prediksi Ketepatan Kelulusan Mahasiswa	Hasil dari pengujian menggunakan algoritma <i>Random Forest</i> tanpa pemilihan fitur memiliki nilai akurasi sebesar 73% dengan nilai <i>AUC</i> 0.5. Sedangkan hasil pengujian menggunakan algoritma <i>Random Forest</i> dengan pemilihan fitur metode <i>wrapper</i> dan metode <i>embedded</i> memiliki nilai akurasi sebesar 100% dengan nilai <i>AUC</i> 1.
2	(Oon Wira Yuda et al., 2022)	Penerapan <i>Data Mining</i> Untuk Klasifikasi Kelulusan Mahasiswa Tepat Waktu Menggunakan Metode <i>Random Forest</i>	Berdasarkan dari hasil pengujian menggunakan algoritma <i>Random Forest</i> , akurasi yang diperoleh sebesar 0.98 atau 98%. Dengan diperolehnya tingkat akurasi tersebut, dapat disimpulkan bahwa pengujian dengan algoritma <i>Random Forest</i> memiliki hasil yang akurat

			dalam klasifikasi kelulusan mahasiswa.
3	(Rachmatika & Bisri, 2020)	Perbandingan Model Klasifikasi untuk Evaluasi Kinerja Akademik Mahasiswa	Dalam hasil uji Friedman, algoritma <i>Random Forest</i> memiliki peringkat tertinggi dengan nilai rata-rata 8,38. Berdasarkan hasil tersebut, dapat disimpulkan bahwa perbandingan pada Sembilan algoritma klasifikasi yang di evaluasi pada empat dataset, algoritma <i>Random Forest</i> lebih unggul dan dapat diandalkan dalam mengevaluasi kinerja akademik mahasiswa
4	(Linawati et al., 2020)	Prediksi Prestasi Akademik Mahasiswa menggunakan Algoritma <i>Random Forest</i> dan C4.5	Algoritma <i>Decision Tree</i> C4.5 menghasilkan akurasi sebesar 87.1%, presisi sebesar 85.4%, dan <i>recall</i> sebesar 87.1%. Sedangkan <i>Random Forest</i> menghasilkan akurasi sebesar 92.4%, presisi sebesar 91.4%, dan <i>recall</i> sebesar 92.4%. Dalam hal ini <i>Random Forest</i> memiliki akurasi, presisi, dan <i>recall</i> yang lebih baik dibandingkan algoritma <i>Decision Tree</i> C4.5 yang menjadikannya sebagai solusi untuk memprediksi prestasi akademik mahasiswa.
5	(Hasan et al., 2022)	Perbandingan K- <i>Nearest Neighbor</i> dan <i>Random Forest</i> dengan Seleksi Fitur <i>Information Gain</i> untuk Klasifikasi Lama Studi Mahasiswa	Hasil menunjukkan bahwa 15 atribut dan salah satu diantaranya merupakan atribut target, didapatkan 4 atribut terbaik berdasarkan seleksi fitur <i>Information Gain</i> . Algoritma <i>Random Forest</i> memiliki akurasi sebesar 100% dibandingkan algoritma K-NN yang memiliki akurasi sebesar 86.67% yang menjadikan algoritma <i>Random Forest</i> bekerja lebih baik dalam mengklasifikasi lama studi mahasiswa.

6	(Alhassan et al., 2020)	<i>Predict Students' Academic Performance based on their Assessment Grades and Online Activity Data</i>	Dalam model prediksi, hasil menunjukkan bahwa model dasar dan sub model bekerja dengan baik dalam memprediksi prestasi akademik siswa. <i>Random Forest</i> mengungguli algoritma lain dalam memprediksi kinerja siswa dengan akurasi tertinggi untuk model dasar dan sub model, diikuti oleh <i>decision tree</i> .
---	-------------------------	---	--

Berdasarkan dari penelitian-penelitian terdahulu di atas tentang klasifikasi nilai mahasiswa yang telah dilakukan, ada banyak metode atau teknik *data mining* yang digunakan dalam melakukan klasifikasi dan prediksi. Dalam penelitiannya (Rachmatika & Bisri, 2020), (Linawati et al., 2020), dan (Alhassan et al., 2020) memiliki pembahasan yang berkaitan dengan nilai mahasiswa menggunakan algoritma *Random Forest*. Namun, yang membedakan penelitian ini dengan penelitian sebelumnya adalah indikator, dan seleksi fitur yang digunakan dalam klasifikasi nilai mahasiswa. Penelitian ini mengangkat studi kasus pada Universitas Muhammadiyah Kalimantan Timur dengan indikator atribut data yang diperoleh dari MKDU dan BAA UMKT dan penerapan seleksi fitur *Analysis of Variance (ANOVA)* pada algoritma *Random Forest Classifier* dalam mengklasifikasi nilai mahasiswa.