

Nama Dosen : Teguh Iman Hermanto, M.Kom  
 Mata Kuliah : Machine Learning 1  
 Pembahasan : Exploratory Data Analysis (EDA)  
 Pokok Pemb :  
     - Mengenal Library ML  
     - Mengenal Statistik Deskriptif  
     - Mengenal EDA Data Numerik  
     - Mengenal EDA Data Kategori  
     - Mengenal EDA data Multivariabel  
     - Mengenal Univariate Analysis  
     - Mengenal Bivariate Analysis

## 1. Mengenal Library Machine Learning

```

● ● ●
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns

```

```

● ● ●
1 df = pd.read_csv("adult.data.csv")

```

```

● ● ●
1 df.head()

```

```

● ● ●
1 df.info()

```

## 2. Mengenal Statistik Deskriptif

```

● ● ●
1 # Menghitung statistik deskriptif untuk kolom numerik
2 numerical_stats = df.describe()
3 print("Statistik Deskriptif untuk Kolom Numerik:\n", numerical_stats)

```

```

● ● ●
1 # Menghitung frekuensi untuk kolom kategorikal
2 categorical_stats = df.select_dtypes(include=['category']).apply(lambda x: x.value_counts(normalize=True) * 100)
3 print("\nPersentase Frekuensi untuk Kolom Kategorikal:\n", categorical_stats)

```

```
● ● ●  
1 # Menampilkan jumlah missing values pada tiap kolom  
2 missing_values = df.isnull().sum()  
3 print("\nJumlah Missing Values pada tiap kolom:\n", missing_values)
```

### 3. Mengenal Bentuk EDA dengan Data Numerik

```
● ● ●  
1 # Histogram  
2 plt.figure(figsize=(8, 6))  
3 sns.histplot(df['age'], kde=True)  
4 plt.title('Histogram of Age')  
5 plt.xlabel('Age')  
6 plt.ylabel('Frequency')  
7 plt.show()
```

```
● ● ●  
1 # Box Plot  
2 plt.figure(figsize=(8, 6))  
3 sns.boxplot(x='sex', y='hours-per-week', data=df)  
4 plt.title('Box Plot of Hours per Week by Sex')  
5 plt.xlabel('Sex')  
6 plt.ylabel('Hours per Week')  
7 plt.show()
```

```
● ● ●  
1 # Scatter Plot  
2 plt.figure(figsize=(8, 6))  
3 sns.scatterplot(x='age', y='capital-gain', data=df)  
4 plt.title('Scatter Plot of Age vs Capital Gain')  
5 plt.xlabel('Age')  
6 plt.ylabel('Capital Gain')  
7 plt.show()
```

#### 4. Mengenal Bentuk EDA dengan Data Kategori

```
● ● ●
1 # Line Plot (Example: Average capital gain by age)
2 avg_capital_gain_by_age = df.groupby('age')['capital-gain'].mean()
3 plt.figure(figsize=(8, 6))
4 plt.plot(avg_capital_gain_by_age.index, avg_capital_gain_by_age.values)
5 plt.title('Line Plot of Average Capital Gain by Age')
6 plt.xlabel('Age')
7 plt.ylabel('Average Capital Gain')
8 plt.show()
```

```
● ● ●
1 # Bar Chart
2 # Contoh: Melihat jumlah orang berdasarkan tingkat pendidikan
3 education_counts = df['education'].value_counts()
4 plt.figure(figsize=(10, 5))
5 plt.bar(education_counts.index, education_counts.values)
6 plt.xlabel('Tingkat Pendidikan')
7 plt.ylabel('Jumlah Orang')
8 plt.title('Jumlah Orang Berdasarkan Tingkat Pendidikan')
9 plt.xticks(rotation=45, ha='right')
10 plt.show()
```

```
● ● ●
1 # Pie Chart
2 # Contoh: Melihat persentase jenis kelamin
3 sex_counts = df['sex'].value_counts()
4 plt.figure(figsize=(6, 6))
5 plt.pie(sex_counts.values, labels=sex_counts.index, autopct='%.1f%%', startangle=90)
6 plt.title('Persentase Jenis Kelamin')
7 plt.show()
```

```
● ● ●
1 # Donut Chart
2 # Contoh: Melihat persentase ras
3 race_counts = df['race'].value_counts()
4 plt.figure(figsize=(6, 6))
5 plt.pie(race_counts.values, labels=race_counts.index, autopct='%.1f%%', startangle=90, wedgeprops=dict(width=0.4))
6 plt.title('Persentase Ras')
7 plt.show()
```

## 5. Mengenal Bentuk EDA dengan Data Multivariabel

```
● ● ●  
1 # Create a heatmap  
2 # Select numerical features for the heatmap  
3 numerical_features = ['age', 'fnlwgt', 'education-num', 'capital-gain', 'capital-loss', 'hours-per-week']  
4 heatmap_data = df[numerical_features].corr()  
5  
6 plt.figure(figsize=(10, 8))  
7 sns.heatmap(heatmap_data, annot=True, cmap='coolwarm', fmt=".2f")  
8 plt.title('Correlation Heatmap')  
9 plt.show()
```

```
● ● ●  
1 # Create a pair plot  
2 sns.pairplot(df[['age', 'education-num', 'hours-per-week', 'capital-gain', 'capital-loss']])  
3 plt.show()
```

```
● ● ●  
1 # Create a bubble chart  
2 # Select features for the bubble chart  
3 bubble_data = df[['age', 'hours-per-week', 'capital-gain']]  
4  
5 plt.figure(figsize=(10, 8))  
6 plt.scatter(x='age', y='hours-per-week', s='capital-gain', data=bubble_data, alpha=0.5)  
7 plt.xlabel('Age')  
8 plt.ylabel('Hours per Week')  
9 plt.title('Bubble Chart: Age, Hours per Week, and Capital Gain')  
10 plt.show()
```

## 6. Mengenal Univariate Analysis

```
● ● ●  
1 # Univariate Analysis for Age  
2 plt.figure(figsize=(8, 6))  
3 sns.histplot(df['age'], kde=True)  
4 plt.title('Distribution of Age')  
5 plt.xlabel('Age')  
6 plt.ylabel('Frequency')  
7 plt.show()
```

```
● ● ●
1 # Univariate Analysis for Education (Categorical)
2 plt.figure(figsize=(10, 6))
3 sns.countplot(x='education', data=df)
4 plt.title('Count of Individuals by Education Level')
5 plt.xlabel('Education Level')
6 plt.ylabel('Count')
7 plt.xticks(rotation=45, ha='right')
8 plt.show()
```

```
● ● ●
1 # Univariate Analysis for Salary (Categorical)
2 plt.figure(figsize=(6, 6))
3 sns.countplot(x='salary', data=df)
4 plt.title('Count of Individuals by Salary')
5 plt.xlabel('Salary')
6 plt.ylabel('Count')
7 plt.show()
```

## 7. Mengenal Bivariate Analysis

```
● ● ●
1 # 1. Age vs. Salary
2 # Deskripsi: Melihat hubungan antara umur dan tingkat pendapatan.
3 # Kita mengharapkan bahwa orang yang lebih tua mungkin memiliki pendapatan yang lebih tinggi.
4 plt.figure(figsize=(8, 6))
5 sns.boxplot(x='salary', y='age', data=df)
6 plt.title('Hubungan antara Umur dan Tingkat Pendapatan')
7 plt.show()
```

```
● ● ●
1 # 2. Education vs. Salary
2 # Deskripsi: Melihat hubungan antara tingkat pendidikan dan pendapatan.
3 # Kita mengharapkan bahwa orang dengan pendidikan lebih tinggi cenderung memiliki pendapatan lebih tinggi.
4 plt.figure(figsize=(12, 6))
5 sns.countplot(x='education', hue='salary', data=df)
6 plt.title('Hubungan antara Tingkat Pendidikan dan Tingkat Pendapatan')
7 plt.xticks(rotation=45, ha='right')
8 plt.show()
```

```
● ● ●  
1 # 3. Hours-per-week vs. Salary  
2 # Deskripsi: Melihat hubungan antara jumlah jam kerja per minggu dan pendapatan.  
3 # Kita mengharapkan bahwa orang yang bekerja lebih banyak jam cenderung memiliki pendapatan lebih tinggi.  
4 plt.figure(figsize=(8, 6))  
5 sns.boxplot(x='salary', y='hours-per-week', data=df)  
6 plt.title('Hubungan antara Jam Kerja per Minggu dan Tingkat Pendapatan')  
7 plt.show()
```

```
● ● ●  
1 # 4. Sex vs. Salary  
2 # Deskripsi: Melihat perbedaan pendapatan antara pria dan wanita.  
3 plt.figure(figsize=(8, 6))  
4 sns.countplot(x='sex', hue='salary', data=df)  
5 plt.title('Perbedaan Tingkat Pendapatan Berdasarkan Jenis Kelamin')  
6 plt.show()
```

```
● ● ●  
1 # 5. Occupation vs. Salary  
2 # Deskripsi: Melihat hubungan antara jenis pekerjaan dan pendapatan.  
3 # Kita mengharapkan bahwa beberapa jenis pekerjaan memiliki pendapatan yang lebih tinggi daripada yang lain.  
4 plt.figure(figsize=(12, 6))  
5 sns.countplot(x='occupation', hue='salary', data=df)  
6 plt.title('Hubungan antara Jenis Pekerjaan dan Tingkat Pendapatan')  
7 plt.xticks(rotation=45, ha='right')  
8 plt.show()
```

```
● ● ●  
1 # 6. Correlation Matrix  
2 # Deskripsi: Melihat korelasi antara variabel-variabel numerik.  
3 # Kita dapat melihat variabel mana yang memiliki korelasi kuat dengan variabel lainnya.  
4 numeric_cols = df.select_dtypes(include=['number']).columns  
5 correlation_matrix = df[numeric_cols].corr()  
6 plt.figure(figsize=(10, 8))  
7 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")  
8 plt.title('Correlation Matrix')  
9 plt.show()
```

### Latihan

1. Tampilkan Plot Education berdasarkan Salarynya
2. Tampilkan jumlah pendidikan Bachelors dengan gaji <=50k
3. Tampilkan data occupation dengan gaji diatas 50k
4. Tampilkan Plot occupation dengan gaji diatas 50k
5. Tampilkan rata-rata hours per week untuk occupation sales?
6. Tampilkan rata-rata hours per week untuk occupation Prof-Speciality?

Nama Dosen : Teguh Iman Hermanto, M.Kom  
 Mata Kuliah : Machine Learning 1  
 Pembahasan : Simple Linear Regresion  
 Pokok Pemb :  
   - Mengenal Regresi Linier  
   - Membangun Model Regresi Linier  
   - Simulasi Algoritma Regresi Linier  
   - Evaluasi Algoritma Regresi Linier

## 1. Mengenal Regresi Linier

### a. Konsep Regresi Linier

Regresi linier berusaha menemukan garis lurus yang terbaik (disebut line of best fit) untuk mendekati data dan memprediksi nilai target. Garis ini dinyatakan dengan persamaan:

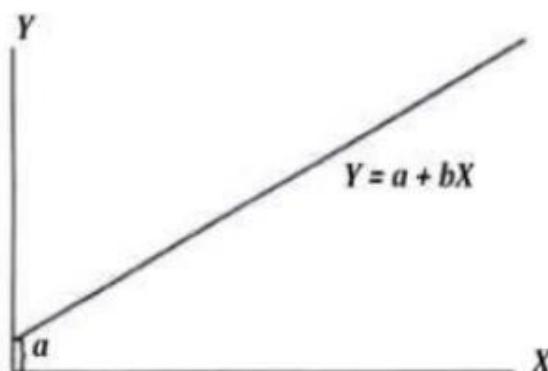
$$Y = a + bX$$

$Y$  = Variabel Response atau Variabel Akibat (Dependent)

$X$  = Variabel Predictor atau Variabel Faktor Penyebab (Independent)

$a$  = Konstanta / intercept

$b$  = Koefisien regresi (kemiringan/slope); besaran response yang ditimbulkan oleh predictor



### b. Langkah-langkah regresi linier

- Membuat dataset

Nama	Kalori/hari(X)	Berat Badan (Y)
Adi	530	89
Rudi	300	48
Didi	358	56
Budi	510	72
Intan	302	54
Putu	300	42
Parta	387	60
Esti	527	85
Wike	415	63
Elis	512	74

- Menghitung nilai variabel

Nama	Kalori/ hari(X)	X2	Berat Badan (Y)	Y2	XY
Adi	530	280900	89	7921	47170
Rudi	300	90000	48	2304	14400
Didi	358	128164	56	3136	20048
Budi	510	260100	72	5184	36720
Intan	302	91204	54	2916	16308
Putu	300	90000	42	1764	12600
Parta	387	149769	60	3600	23220
Esti	527	277729	85	7225	44795
Wike	415	172225	63	3969	26145
Elis	512	262144	74	5476	37888
<b>Σ</b>	<b>4141</b>	<b>1802235</b>	<b>643</b>	<b>43495</b>	<b>279294</b>

- Perhitungan koefisien a dan b

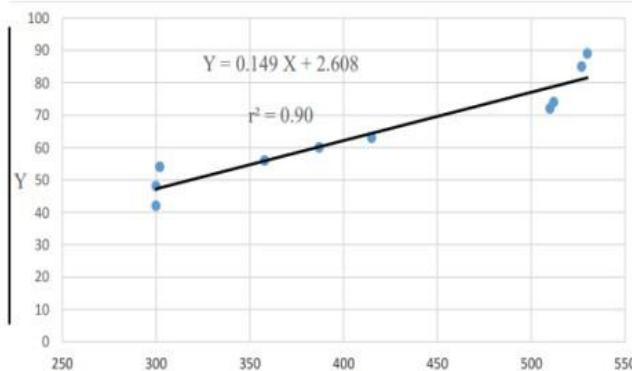
$$a = \frac{(\sum Y_i)(\sum X_i^2) - (\sum X_i)(\sum X_i Y_i)}{n \sum X_i^2 - (\sum X_i)^2}$$

$$= \frac{(643)(1802235) - (4141)(279294)}{10(1802235) - (4141)^2} = \frac{2280651}{874469} \cong 2,608$$

$$b = \frac{n (\sum X_i Y_i) - (\sum X_i)(\sum Y_i)}{n \sum X_i^2 - (\sum X_i)^2}$$

$$= \frac{10(279294) - (4141)(643)}{10(1802235) - (4141)^2} = \frac{130227}{874469} \cong 0,14892 \approx 0,149$$

- Persamaan garis regresi



Setelah didapat koefisien a dan b, maka persamaan garisnya adalah  $Y = 2,608 + 0.149X$

- Hitung estimasi

Melakukan Estimasi terhadap variable predictor atau response

**Estimasi berat badan mahasiswa jika asupan kalori adalah 600 kalori/hari :**

$$Y = 2,608 + 0,149 X$$

$$Y = 2,608 + (0,149 * 600)$$

Estimasi Y = **92 Kilogram**

**Estimasi asupan kalori mahasiswa, jika berat badan mahasiswa adalah 40 kilogram :**

$$40 = 2,608 + 0,149 X$$

$$40 - 2,608 = 0,149 X$$

$$37,392 = 0,149 X$$

Estimasi X = **250.59 kalori/hari**

## 2. Membangun model regresi linier

- Import Library yang dibutuhkan

```
● ● ●
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_squared_error, r2_score
5 import matplotlib.pyplot as plt
6 import seaborn as sns
```

- Deskripsi Dataset

```
● ● ●
1 df = pd.read_csv('konsumsi-bbm.csv')
```

```
● ● ●
```

```
1 df.head()
```

```
● ● ●
```

```
1 df.info()
```

```
● ● ●
```

```
1 df.describe()
```

- Exploratory Data Analysis

```
● ● ●  
1 plt.figure(figsize=(8, 6))  
2 sns.scatterplot(x='Jarak_Tempuh', y='Konsumsi_BBM', data=df)  
3 plt.title('Hubungan Jarak Tempuh dan Konsumsi BBM')  
4 plt.xlabel('Jarak Tempuh')  
5 plt.ylabel('Konsumsi BBM')  
6 plt.show()
```

```
● ● ●  
1 plt.figure(figsize=(8, 6))  
2 sns.boxplot(data=df[['Jarak_Tempuh', 'Konsumsi_BBM']])  
3 plt.title('Box Plot Jarak Tempuh dan Konsumsi BBM')  
4 plt.show()
```

```
● ● ●  
1 plt.figure(figsize=(8, 6))  
2 sns.histplot(df['Jarak_Tempuh'], kde=True)  
3 plt.title('Distribusi jarak tempuh')  
4 plt.xlabel('jarak tempuh')  
5 plt.ylabel('Frekuensi')  
6 plt.show()
```

- Membangun Model Regresi

```
● ● ●  
1 X = df[['Jarak_Tempuh']]  
2 y = df['Konsumsi_BBM']
```

```
● ● ●
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
● ● ●
1 model = LinearRegression()
```

```
● ● ●
1 model.fit(X_train, y_train)
```

```
● ● ●
1 y_pred = model.predict(X_test)
```

### 3. Simulasi algoritma regresi linier

```
● ● ●
1 jarak_baru = [60, 70, 80] # Contoh jarak tempuh baru
2 konsumsi_prediksi = model.predict([[jarak] for jarak in jarak_baru])
```

```
● ● ●
1 print('\nSimulasi Konsumsi BBM:')
2 for i in range(len(jarak_baru)):
3     print(f'Jarak Tempuh: {jarak_baru[i]}, Konsumsi BBM Prediksi: {konsumsi_prediksi[i]}')
```

```
● ● ●
1 plt.scatter(X_train, y_train, color='blue', label='Data Pelatihan')
2 plt.plot(X_train, model.predict(X_train), color='red', label='Garis Regresi')
3 plt.scatter(jarak_baru, konsumsi_prediksi, color='green', label='Data Baru (Prediksi)')
4
5 plt.xlabel('jarak')
6 plt.ylabel('konsumsi bbm')
7 plt.title('Grafik Regresi Linier')
8 plt.legend()
9 plt.show()
```

#### 4. Evaluasi algoritma regresi linier



```
1 # Menghitung MSE dan R-squared  
2 mse = mean_squared_error(y_test, y_pred)  
3 r2 = r2_score(y_test, y_pred)
```



```
1 print(f'Mean Squared Error (MSE): {mse}')  
2 print(f'R-squared (R2): {r2}')
```

Nama Dosen : Teguh Iman Hermanto, M.Kom  
 Mata Kuliah : Machine Learning 1  
 Pembahasan : Simple Linear Regresion  
 Pokok Pemb : - Regresi Linier Berganda  
     - Membangun Model Regresi Linier Berganda  
     - Simulasi Algoritma Regresi Linier Berganda  
     - Evaluasi Algoritma Regresi Linier Berganda

## 1. Regresi Linier Berganda

### a. Konsep Regresi Linier Berganda

Regresi Linier Berganda adalah salah satu teknik statistik yang digunakan untuk memodelkan hubungan antara satu variabel dependen (output/label) dan lebih dari satu variabel independen (input/fitur). Regresi linier berganda memperluas konsep regresi linier sederhana, yang hanya melibatkan satu variabel independen, dengan memasukkan beberapa variabel independen dalam analisis.

Tujuan utama dari regresi linier berganda adalah untuk menemukan hubungan linear antara variabel-variabel independen dan variabel dependen, sehingga kita dapat memprediksi nilai variabel dependen berdasarkan nilai-nilai variabel independen.

### b. Aplikasi Regresi Linier Berganda

Regresi linier berganda sering digunakan dalam berbagai bidang seperti:

- **Ekonomi:** untuk memprediksi pendapatan berdasarkan berbagai faktor seperti tingkat pendidikan, pengalaman kerja, dan usia.
- **Pemasaran:** untuk menganalisis dampak harga, iklan, dan lokasi terhadap penjualan produk.
- **Kesehatan:** untuk memprediksi jumlah kalori yang dibakar berdasarkan usia, berat badan, tinggi badan, dan durasi olahraga.

### c. Rumus Regresi Linier Berganda

Secara matematis, persamaan regresi linier berganda dapat dinyatakan sebagai berikut:

$$Y = a + b_1 X_1 + b_2 X_2 + b_3 X_3 + \dots + b_n X_n$$

Di mana:

- $Y$  = adalah variabel dependen (yang ingin diprediksi).
- $X_1, X_2, X_3, \dots, X_n$  = adalah variabel independen (faktor yang memengaruhi variabel dependen).
- $a$  = adalah konstanta atau intercept, yaitu nilai  $Y$  saat semua variabel independen bernilai nol.
- $b_1, b_2, b_3, \dots, b_n$  adalah koefisien regresi, yang menunjukkan seberapa besar pengaruh setiap variabel independen terhadap variabel dependen.

## 2. Membangun Model Regresi Linier Berganda

```
● ● ●  
1 import pandas as pd  
2 import numpy as np  
3  
4 from sklearn.linear_model import LinearRegression  
5 from sklearn.model_selection import train_test_split  
6 from sklearn.metrics import mean_squared_error, r2_score  
7  
8 import matplotlib.pyplot as plt  
9 import seaborn as sns  
10 import plotly.express as px  
11 import plotly.graph_objects as go
```

```
● ● ●  
1 df = pd.read_csv('kalori_ganda.csv')
```

```
● ● ●  
1 # 1. Scatter plot untuk melihat hubungan antara Umur dan Kalori  
2 plt.figure(figsize=(7,5))  
3 plt.scatter(x='Umur', y='Kalori', data=df, color='blue')  
4 plt.title('Age vs Calorie Intake')  
5 plt.xlabel('Age (years)')  
6 plt.ylabel('Calorie Intake (kcal)')  
7 plt.grid(True)  
8 plt.show()
```



```
1 # 2. Berat Badan vs Konsumsi Kalori
2 plt.figure(figsize=(7,5))
3 plt.scatter(x='BB', y='Kalori', data=df, color='green')
4 plt.title('Weight vs Calorie Intake')
5 plt.xlabel('Weight (kg)')
6 plt.ylabel('Calorie Intake (kcal)')
7 plt.grid(True)
8 plt.show()
```



```
1 # Histogram untuk melihat distribusi BB
2 plt.figure(figsize=(8, 6))
3 sns.histplot(x='BB', data=df)
4 plt.title('Distribusi BB')
5 plt.show()
```



```
1 # menampilkan plot animasi 3d
2 fig = px.scatter_3d(df, x='Umur', y='BB', z='TB',
3                      color='Olahraga', size='Kalori',
4                      hover_name='Umur')
5 fig.show()
```

### 3. Simulasi Algoritma Regresi Linier Berganda

```
● ● ●  
1 # Memisahkan fitur (X) dan target (y)  
2 x = df[['Umur','BB','TB','Olahraga']]  
3 y = df['Kalori']
```

```
● ● ●  
1 # Membagi data menjadi data latih dan data uji  
2 X_train, X_test, y_train, y_test = train_test_split(  
3     x, y, test_size=0.2, random_state=42)
```

```
● ● ●  
1 # Membuat model regresi linier  
2 model = LinearRegression()  
3 model.fit(x, y)
```

```
● ● ●  
1 # Menampilkan intercept regresi dari model  
2 intercept = model.intercept_  
3 intercept
```



```
1 # Menampilkan koefisien regresi dari model
2 coefficients = model.coef_
3 coefficients
```



```
1 manual = intercept+(25*3.05138649)+(75*22.34269713)+(170*0.0466017)+(60*0.06622248)
2 manual
```



```
1 # Data input untuk prediksi
2 input_data = np.array([[25, 75, 170, 60]])
```



```
1 # Melakukan prediksi
2 predicted_calories = model.predict(input_data)
3 predicted_calories
```

#### 4. Evaluasi Regresi Linier Berganda



```
1 # Memprediksi nilai BB dengan data uji  
2 y_pred = model.predict(X_test)
```



```
1 # Menghitung MSE dan R2 score  
2 mse = mean_squared_error(y_test, y_pred)  
3 r2 = r2_score(y_test, y_pred)
```



```
1 # Menampilkan MSE dan R2 score  
2 print(f'Mean Squared Error (MSE): {mse}')  
3 print(f'R-squared (R2): {r2}')
```

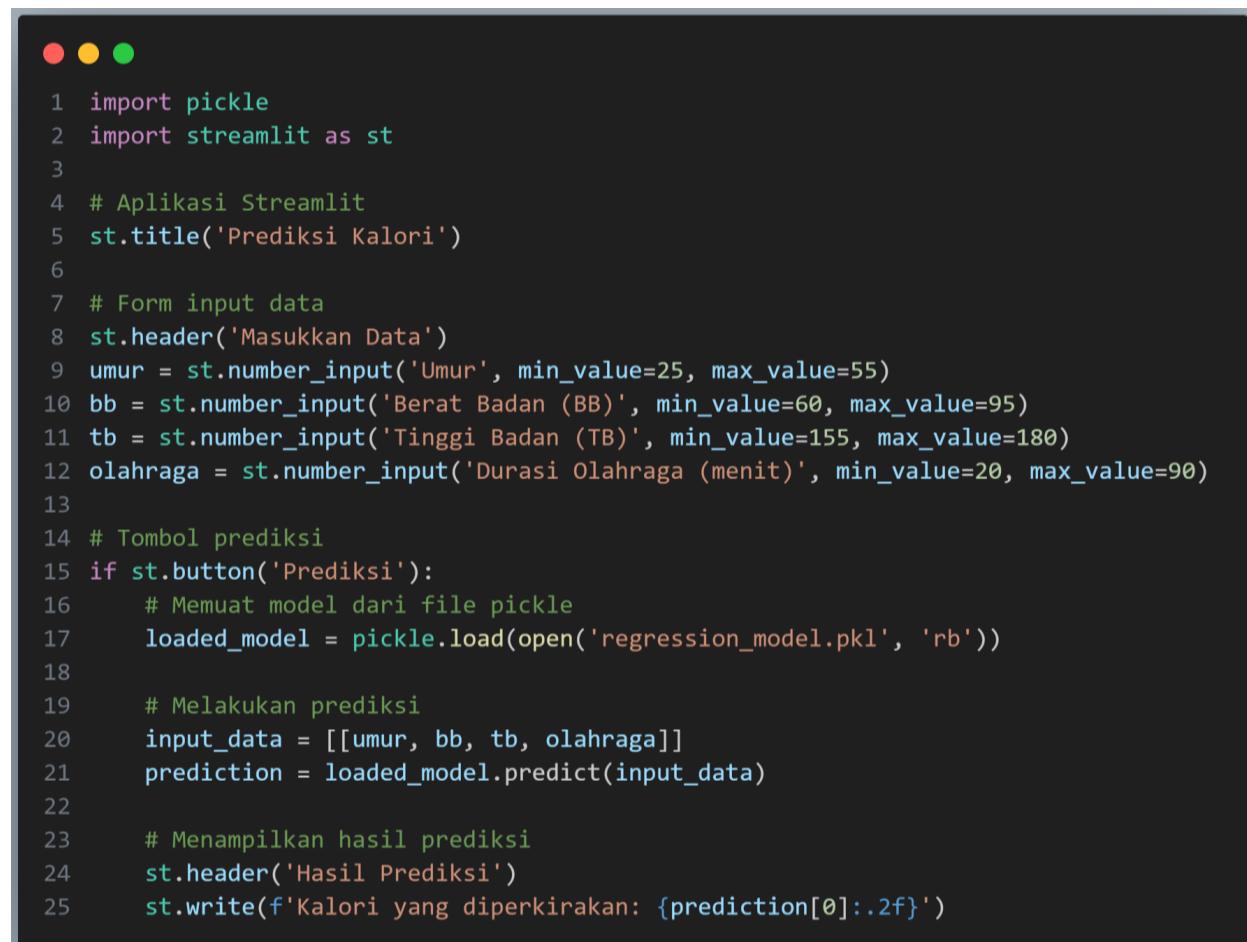
## 5. Simpan Model

```
● ● ●  
1 # Menyimpan model menggunakan pickle  
2 import pickle  
3  
4 filename = 'regression_model.pkl'  
5 pickle.dump(model, open(filename, 'wb'))
```

Dari perintah ini akan menghasilkan file baru yaitu **regression\_model.pkl**. file ini akan otomatis muncul di direktori project kalian. Apabila berhasil, file model ini bisa kita gunakan untuk membuat aplikasi machine learning berbasis web menggunakan streamlit.

## 6. Membuat Aplikasi

Buat file **regresi-ganda.py** pada direktori yang sama.

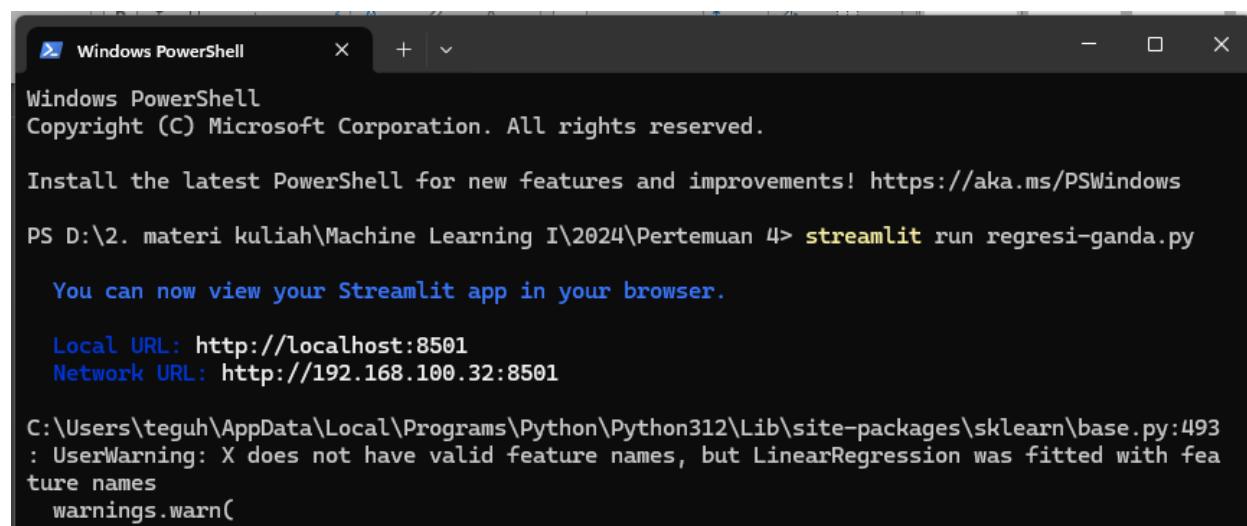


```

1 import pickle
2 import streamlit as st
3
4 # Aplikasi Streamlit
5 st.title('Prediksi Kalori')
6
7 # Form input data
8 st.header('Masukkan Data')
9 umur = st.number_input('Umur', min_value=25, max_value=55)
10 bb = st.number_input('Berat Badan (BB)', min_value=60, max_value=95)
11 tb = st.number_input('Tinggi Badan (TB)', min_value=155, max_value=180)
12 olahraga = st.number_input('Durasi Olahraga (menit)', min_value=20, max_value=90)
13
14 # Tombol prediksi
15 if st.button('Prediksi'):
16     # Memuat model dari file pickle
17     loaded_model = pickle.load(open('regression_model.pkl', 'rb'))
18
19     # Melakukan prediksi
20     input_data = [[umur, bb, tb, olahraga]]
21     prediction = loaded_model.predict(input_data)
22
23     # Menampilkan hasil prediksi
24     st.header('Hasil Prediksi')
25     st.write(f'Kalori yang diperkirakan: {prediction[0]:.2f}')

```

Untuk menjalankan aplikasinya buka command prompt, lalu arahkan ke direktori project dan masukan perintah **streamlit run nama\_aplikasi.py**



```

Windows PowerShell
Copyright (C) Microsoft Corporation. All rights reserved.

Install the latest PowerShell for new features and improvements! https://aka.ms/PSWindows

PS D:\2. materi kuliah\Machine Learning I\2024\Pertemuan 4> streamlit run regresi-ganda.py

You can now view your Streamlit app in your browser.

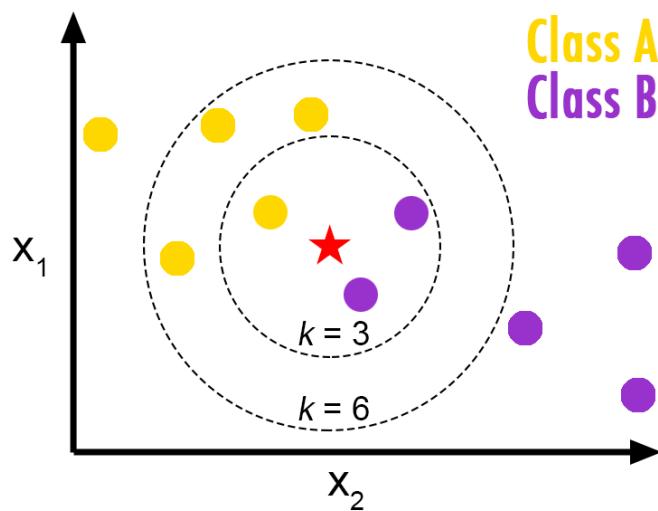
Local URL: http://localhost:8501
Network URL: http://192.168.100.32:8501

C:\Users\teguh\AppData\Local\Programs\Python\Python312\Lib\site-packages\sklearn\base.py:493
: UserWarning: X does not have valid feature names, but LinearRegression was fitted with fea
ture names
    warnings.warn(

```

Nama Dosen : Teguh Iman Hermanto, M.Kom  
 Mata Kuliah : Machine Learning 1  
 Pembahasan : K Nearest Neighbors  
 Pokok Pemb :  
 - Konsep Algoritma KNN  
 - Membangun Model KNN  
 - Simulasi Algoritma KNN  
 - Evaluasi Algoritma KNN  
 - Aplikasi menggunakan algoritma KNN

### 1. Konsep Algoritma KNN



- k-NN mengidentifikasi nilai “k” dalam data pelatihan yang merupakan “tetangga terdekat” (Lantz, 2013)
- Metode untuk melakukan klasifikasi terhadap objek berdasarkan pembelajaran yang jaraknya paling dekat dengan objek tersebut

Model algoritma KNN

- Euclidean Distance

$$Ed(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

- Manhattan Distance

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

Langkah-langkah Algoritma KNN :

1. Tentukan data testing dan nilai parameter k
2. Hitung jarak antara data testing dan data training
3. Urutkan data berdasarkan jarak terkecil
4. Kelompokkan data

## 2. Membangun Model KNN



```
1 import numpy as np
2 import pandas as pd
3
4 from sklearn.preprocessing import LabelEncoder
5 from sklearn.model_selection import train_test_split
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.metrics import accuracy_score, confusion_matrix
8
9 import matplotlib.pyplot as plt
10 import seaborn as sns
```



```
1 df = pd.read_csv('Pinjaman.csv')
```



```
1 df.head()
```

```
● ● ●  
1 # Distribution of Age  
2 plt.figure(figsize=(8, 6))  
3 sns.histplot(df['Usia'], bins=10, kde=True, color='blue')  
4 plt.title('Distribution of Age')  
5 plt.xlabel('Usia')  
6 plt.ylabel('Frequency')  
7 plt.show()
```

```
● ● ●  
1 # Boxplot of Income (Pendapatan) vs Lulus Kredit  
2 plt.figure(figsize=(8, 6))  
3 sns.boxplot(x='Lulus_Kredit', y='Pendapatan', data=df, palette='Set1')  
4 plt.title('Pendapatan vs Kelayakan Kredit')  
5 plt.xlabel('Lulus Kredit')  
6 plt.ylabel('Pendapatan (Juta)')  
7 plt.xticks(ticks=[0, 1], labels=['Tidak Layak', 'Layak'])  
8 plt.show()
```

```
● ● ●  
1 # Countplot of Marital Status  
2 plt.figure(figsize=(8, 6))  
3 sns.countplot(x='Status_Perkawinan', data=df, palette='Set2')  
4 plt.title('Count of Marital Status')  
5 plt.xticks(ticks=[0, 1], labels=['Belum Menikah', 'Menikah'])  
6 plt.ylabel('Frequency')  
7 plt.show()
```



```

1 # Pairplot for the numerical features
2 plt.figure(figsize=(8, 6))
3 sns.pairplot(df[['Usia', 'Pendapatan',
4                  'Jumlah_Pinjaman',
5                  'Durasi_Pinjaman',
6                  'Lulus_Kredit']],
7                  hue='Lulus_Kredit', palette='coolwarm')
8 plt.show()

```



```

1 # Violin plot of Loan Amount vs Credit Eligibility
2 plt.figure(figsize=(8, 6))
3 sns.violinplot(x='Lulus_Kredit', y='Jumlah_Pinjaman', data=df, palette='Set2')
4 plt.title('Distribusi Jumlah Pinjaman berdasarkan Kelayakan Kredit')
5 plt.xlabel('Lulus Kredit')
6 plt.ylabel('Jumlah Pinjaman (Juta)')
7 plt.xticks(ticks=[0, 1], labels=['Tidak Layak', 'Layak'])
8 plt.show()

```



```

1 label_encoder_status_perkawinan = LabelEncoder()
2 label_encoder_status_pekerjaan = LabelEncoder()
3 label_encoder_lulus_kredit = LabelEncoder()

```



```

1 df['Status_Perkawinan'] = label_encoder_status_perkawinan.fit_transform(df['Status_Perkawinan'])
2 df['Status_Pekerjaan'] = label_encoder_status_pekerjaan.fit_transform(df['Status_Pekerjaan'])
3 df['Lulus_Kredit'] = label_encoder_lulus_kredit.fit_transform(df['Lulus_Kredit'])

```

Status Perkawinan (0:Belum Menikah) | (1:Menikah)

Status Pekerjaan (0:Karyawan Kontrak) | (1:Karyawan Tetap) | (2:Pensiunan) | (3:Wirausaha)

Lulus Kredit (0:layak) | (1:tidak layak)



```
1 df.to_csv('Pinjaman_modif.csv', index=False)
```



```
1 # Splitting the dataset into features and labels
2 X = df.drop(columns=['Lulus_Kredit'])
3 y = df['Lulus_Kredit']
```



```
1 # Splitting data into training and testing sets (80% train, 20% test)
2 X_train, X_test, y_train, y_test = train_test_split(
3     X, y, test_size=0.2, random_state=42)
```



```
1 # Creating and training the KNN model
2 knn_model = KNeighborsClassifier(n_neighbors=5)
3 knn_model.fit(X_train, y_train)
```



```
1 # Predicting the test data
2 y_pred = knn_model.predict(X_test)
```

### 3. Simulasi Algoritma KNN

```
● ● ●  
1 # Example:  
2 # [Usia : 30,  
3 # Pendapatan : 80,  
4 # Status Perkawinan : 0,  
5 # Jumlah Pinjaman : 70,  
6 # Durasi Pinjaman : 10,  
7 # Status Pekerjaan : 1]  
8 new_data = np.array([[30, 80, 0, 70, 10, 1]])  
9 new_pred = knn_model.predict(new_data)
```

```
● ● ●  
1 print('Hasil prediksi: ', new_pred)
```

```
● ● ●  
1 label_encoder_lulus_kredit.inverse_transform(new_pred)
```

#### 4. Evaluasi Algoritma KNN

```
● ● ●  
1 # Calculating the accuracy  
2 accuracy = accuracy_score(y_test, y_pred)  
3 accuracy
```

```
● ● ●  
1 # Menyimpan model menggunakan pickle  
2 import pickle  
3  
4 filename = 'knn_pinjam_mod.pkl'  
5 pickle.dump(knn_model, open(filename, 'wb'))
```

## 5. Aplikasi Menggunakan Algoritma KNN

Buat file **knn\_pinjaman.py** pada direktori project

```

● ● ●
1 import streamlit as st
2 import numpy as np
3 import pickle
4 import plotly.express as px
5 import pandas as pd
6
7 # Load the trained KNN model
8 with open('knn_pinjam_mod.pkl', 'rb') as model_file:
9     knn = pickle.load(model_file)
10
11 # Create a title for the web app
12 st.title('K-Nearest Neighbor (KNN) Loan Eligibility Prediction')
13
14 # Input fields for new data
15 usia = st.sidebar.number_input('Usia', min_value=18, max_value=70, value=30)
16 pendapatan = st.sidebar.number_input('Pendapatan', min_value=10, max_value=200, value=50)
17 status_perkawinan = st.sidebar.selectbox('Status_Perkawinan',['Belum Menikah','Menikah'])
18 if status_perkawinan == 'Belum Menikah':
19     status_perkawinan = 0
20 else:
21     status_perkawinan = 1
22 jumlah_pinjaman = st.sidebar.number_input('Jumlah_Pinjaman', min_value=10, max_value=500, value=100)
23 durasi_pinjaman = st.sidebar.number_input('Durasi_Pinjaman', min_value=1, max_value=30, value=10)
24 #status_pekerjaan = st.sidebar.number_input('Status_Pekerjaan', min_value=0, max_value=3, value=0)
25 status_pekerjaan = st.sidebar.selectbox('Status_Pekerjaan',['Karyawan Kontrak','Karyawan Tetap','Pensiunan','Wirausaha'])
26 if status_pekerjaan == 'Karyawan Kontrak':
27     status_pekerjaan = 0
28 elif status_pekerjaan == 'Karyawan Tetap':
29     status_pekerjaan = 1
30 elif status_pekerjaan == 'Pensiunan':
31     status_pekerjaan = 2
32 elif status_pekerjaan == 'Wirausaha':
33     status_pekerjaan = 3
34
35 # New data for prediction
36 new_data = np.array([[usia, pendapatan, status_perkawinan, jumlah_pinjaman, durasi_pinjaman, status_pekerjaan]])
37
38 # Predict button
39 if st.sidebar.button('Predict'):
40     # Predict the result
41     prediction = knn.predict(new_data)
42     result = 'Layak' if prediction == 0 else 'Tidak Layak'
43
44     # Show the prediction result
45     st.write(f'Hasil prediksi: **{result}**')
46
47
48 # 3D Plot for KNN
49 df = pd.DataFrame({
50     'Usia': np.random.randint(18, 70, size=60),
51     'Pendapatan': np.random.randint(10, 200, size=60),
52     'Jumlah_Pinjaman': np.random.uniform(10, 500, size=60),
53     'Lulus_Kredit': np.random.choice(['0','1'], size=60)
54 })
55 fig = px.scatter_3d(df, x='Usia', y='Pendapatan', z='Jumlah_Pinjaman', color='Lulus_Kredit',
56                      title=f'KNN Prediction')
57 st.plotly_chart(fig)

```

Untuk menjalankannya buka command prompt, lalu arahkan ke direktori project dan masukan perintah `streamlit run nama_aplikasi.py`

Nama Dosen	: Teguh Iman Hermanto, M.Kom
Mata Kuliah	: Machine Learning 1
Pembahasan	: Naïve Bayes
Pokok Pemb	<ul style="list-style-type: none"> <li>- Konsep Algoritma Naïve Bayes</li> <li>- Membangun Model Naïve Bayes</li> <li>- Simulasi Algoritma Naïve Bayes</li> <li>- Evaluasi Algoritma Naïve Bayes</li> <li>- Aplikasi menggunakan algoritma Naïve Bayes</li> </ul>

### 1. Konsep Algoritma Naïve Bayes

Naïve Bayes Classifier merupakan sebuah metoda klasifikasi yang berakar pada teorema Bayes . Metode pengklasifikasian dg menggunakan metode probabilitas dan statistik yg dikemukakan oleh ilmuwan Inggris Thomas Bayes , yaitu memprediksi peluang di masa depan berdasarkan pengalaman di masa sebelumnya sehingga dikenal sebagai Teorema Bayes . Ciri utama dr Naïve Bayes Classifier ini adalah asumsi yg sangat kuat (naïf) akan independensi dari masing-masing kondisi / kejadian.

Model Algoritma Naïve Bayes

$$P(H|X) = \frac{P(X|H)P(H)}{P(X)} = P(X|H) \times P(H)/P(X)$$

x = Data dari class scoring yang belum diketahui

H = Hipotesis data X yang merupakan suatu class yang lebih spesifik

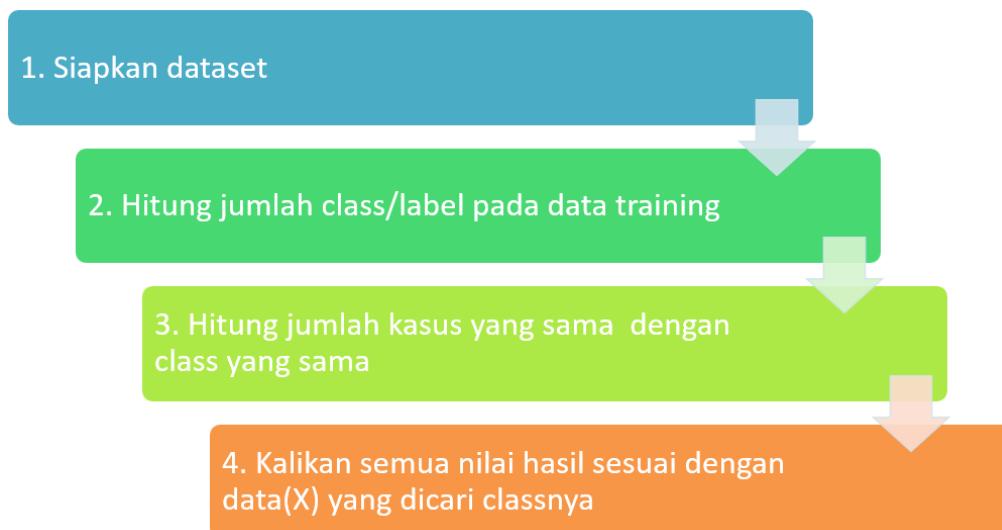
$P(H|X)$  = Probabilitas hipotesis H berdasarkan kondisi X

$P(H)$  = Probabilitas hipotesis H

$P(X|H)$  = Probabilitas hipotesis X berdasarkan kondisi H

$P(X)$  = Probabilitas X

Tahapan Algoritma Naïve Bayes



**a. Siapkan Dataset**

- Data training buys\_computer dataset

Age	Income	Student	Credit Rating	Buys_Computer
youth	high	no	fair	no
youth	high	no	excellent	no
middle_aged	high	no	fair	yes
senior	medium	no	fair	yes
senior	low	yes	fair	yes
senior	low	yes	excellent	no
middle_aged	low	yes	excellent	yes
youth	medium	no	fair	no
youth	low	yes	fair	yes
senior	medium	yes	fair	yes
youth	medium	yes	excellent	yes
middle_aged	medium	no	excellent	yes
middle_aged	high	yes	fair	yes
senior	medium	no	excellent	no

- Data testing

Age	Income	Student	Credit Rating	Buys_Computer
Senior	high	no	fair	?

Apakah calon pembeli akan membeli komputer?

**b. Hitung jumlah class/label pada data training**

Class 1 (C1) → Buys\_Computer yes → 9

Class 2 (C2) → Buys\_Computer no → 5

Total Record: 14

Maka:

$$P(C1) = 9/14 = 0.642857143$$

$$P(C2) = 5/14 = 0.357142857$$

**c. Hitung jumlah kasus yang sama dengan class yang sama**

Contoh penghitungan  $P(X|Ci)$  pada kasus Income dimana

i=1 dan i=2

Maka :

$$P(\text{Income}=\text{"high"} | \text{Buys_Computer}=\text{"yes"}) = 2/9 = 0.222222222$$

$$P(\text{Income}=\text{"high"} | \text{Buys_Computer}=\text{"no"}) = 2/5 = 0.4$$

d. Kalikan semua nilai hasil sesuai dengan data(X) yang dicari classnya

Jika semua data aribut pada data training dihitung, maka:

Atribut	Parameter	yes	no
Age	youth	0.222222222	0.6
Age	middle_aged	0.444444444	0
<b>Age</b>	<b>senior</b>	<b>0.333333333</b>	<b>0.4</b>
Income	low	0.333333333	0.2
Income	medium	0.444444444	0.4
<b>Income</b>	<b>high</b>	<b>0.222222222</b>	<b>0.4</b>
Student	yes	0.666666667	0.2
<b>Student</b>	<b>no</b>	<b>0.333333333</b>	<b>0.8</b>
<b>Credit Rating</b>	<b>fair</b>	<b>0.666666667</b>	<b>0.4</b>
Credit Rating	excellent	0.333333333	0.6

$$P(X| \text{Buys\_Computer}=\text{"yes"}) = 0.016460905$$

$$P(X| \text{Buys\_Computer}=\text{"no"}) = 0.0512$$

$$P(X| \text{Buys\_Computer}=\text{"yes"})*P(C1) = 0.010582011$$

$$P(X| \text{Buys\_Computer}=\text{"no"})*P(C2) = 0.018285714$$

Nilai "no" lebih besar daripada nilai "yes". Sehingga pembeli dengan atribut X tidak membeli komputer

## 2. Membangun Model Naïve Bayes



```

1 import pandas as pd
2 import numpy as np
3
4 from sklearn.model_selection import train_test_split
5 from sklearn.naive_bayes import CategoricalNB
6 from sklearn.metrics import accuracy_score, confusion_matrix
7
8 import matplotlib.pyplot as plt
9 import seaborn as sns

```



```
1 df = pd.read_csv('Diabetes1.csv')
```

	Age	BMI	Blood_Sugar_Level	Family_History	Diet	Diabetes
0	45	28.7	135	Yes	Poor	Yes
1	50	31.2	145	No	Moderate	Yes
2	30	22.0	95	No	Good	No
3	35	25.4	105	Yes	Poor	Yes
4	60	33.5	155	Yes	Poor	Yes



```
1 df.head()
```



```
1 df.info()
```



```
1 # 1. Bar Chart: Jumlah pasien diabetes berdasarkan jenis diet
2 plt.figure(figsize=(8, 5))
3 sns.countplot(x='Diet', hue='Diabetes', data=df)
4 plt.title("Jumlah pasien diabetes berdasarkan jenis diet")
5 plt.xlabel("Diet Type")
6 plt.ylabel("Count")
7 plt.legend(title="Diabetes")
8 plt.show()
```

```
● ● ●  
1 # 2. Pie Chart: Persentase pasien dengan riwayat keluarga  
2 family_history_counts = df['Family_History'].value_counts()  
3 plt.figure(figsize=(5, 5))  
4 family_history_counts.plot.pie(autopct='%1.1f%%', startangle=140, colors=['#ff9999', '#66b3ff'])  
5 plt.title("Persentase pasien dengan riwayat keluarga")  
6 plt.ylabel("")  
7 plt.show()
```

```
● ● ●  
1 # 3. Line Chart: Tren kadar gula darah berdasarkan usia  
2 df_sorted = df.sort_values(by='Age')  
3 plt.figure(figsize=(8, 5))  
4 plt.plot(df_sorted['Age'], df_sorted['Blood_Sugar_Level'], marker='o')  
5 plt.title("Tren kadar gula darah berdasarkan usia")  
6 plt.xlabel("Age")  
7 plt.ylabel("Blood Sugar Level (mg/dL)")  
8 plt.grid(True)  
9 plt.show()
```

```
● ● ●  
1 # 4. Scatter Plot: Hubungan antara BMI dan Kadar Gula Darah  
2 plt.figure(figsize=(8, 5))  
3 sns.scatterplot(x='BMI', y='Blood_Sugar_Level', hue='Diabetes', data=df)  
4 plt.title("Hubungan antara BMI dan Kadar Gula Darah")  
5 plt.xlabel("BMI")  
6 plt.ylabel("Blood Sugar Level (mg/dL)")  
7 plt.legend(title="Diabetes")  
8 plt.show()
```

```
● ● ●  
1 # 5. Histogram: Distribusi BMI  
2 plt.figure(figsize=(8, 5))  
3 plt.hist(df['BMI'], bins=5, color='skyblue', edgecolor='black')  
4 plt.title("Distribusi BMI")  
5 plt.xlabel("BMI")  
6 plt.ylabel("Frequency")  
7 plt.show()
```

```
● ● ●  
1 # 6. Box Plot: Tingkat Gula Darah berdasarkan Riwayat Keluarga  
2 plt.figure(figsize=(8, 5))  
3 sns.boxplot(x='Family_History', y='Blood_Sugar_Level', data=df)  
4 plt.title("Tingkat Gula Darah berdasarkan Riwayat Keluarga")  
5 plt.xlabel("Family History")  
6 plt.ylabel("Blood Sugar Level (mg/dL)")  
7 plt.show()
```

```
● ● ●  
1 # 7. Heatmap: Korelasi antar variabel numerik  
2 plt.figure(figsize=(8, 5))  
3 sns.heatmap(df[['Age', 'BMI', 'Blood_Sugar_Level']].corr(), annot=True, cmap='coolwarm')  
4 plt.title("Correlation Heatmap")  
5 plt.show()
```

```
● ● ●  
1 # 8. Area Chart: Kasus diabetes kumulatif berdasarkan usia (untuk tujuan ilustrasi)  
2 age_cum_cases = df[df['Diabetes'] == 'Yes'].groupby('Age').size().cumsum()  
3 plt.figure(figsize=(8, 5))  
4 plt.fill_between(age_cum_cases.index, age_cum_cases.values, color='lightgreen', alpha=0.6)  
5 plt.title("Kasus Diabetes Kumulatif Berdasarkan Usia")  
6 plt.xlabel("Age")  
7 plt.ylabel("Cumulative Cases")  
8 plt.grid(True)  
9 plt.show()
```

```
● ● ●  
1 # Membuat DataFrame  
2 df_encoded = pd.DataFrame(df)
```

```
● ● ●  
1 # Encode kolom categorical  
2 df_encoded['Family_History'] = df_encoded['Family_History'].map({'Yes': 1, 'No': 0})  
3 df_encoded['Diet'] = df_encoded['Diet'].map({'Poor': 0, 'Moderate': 1, 'Good': 2})  
4 df_encoded['Diabetes'] = df_encoded['Diabetes'].map({'Yes': 1, 'No': 0})
```



```
1 df_encoded.head()
```

	<b>Age</b>	<b>BMI</b>	<b>Blood_Sugar_Level</b>	<b>Family_History</b>	<b>Diet</b>	<b>Diabetes</b>
0	45	28.7		135	1	0
1	50	31.2		145	0	1
2	30	22.0		95	0	2
3	35	25.4		105	1	0
4	60	33.5		155	1	1

```
● ● ●
1 # mempersiapkan data untuk Categorical Naive Bayes
2 categorical_df = df_encoded.copy()
3 categorical_df['Age'] = pd.cut(categorical_df['Age'], bins=[0, 30, 40, 50, 60, 80], labels=[0, 1, 2, 3, 4])
4 categorical_df['BMI'] = pd.cut(categorical_df['BMI'], bins=[0, 18.5, 24.9, 29.9, 50], labels=[0, 1, 2, 3])
5 categorical_df['Blood_Sugar_Level'] = pd.cut(categorical_df['Blood_Sugar_Level'], bins=[0, 100, 125, 150, 200], labels=[0, 1, 2, 3])
```



```
1 categorical_df.head()
```

	<b>Age</b>	<b>BMI</b>	<b>Blood_Sugar_Level</b>	<b>Family_History</b>	<b>Diet</b>	<b>Diabetes</b>
0	2	2		2	1	0
1	2	3		2	0	1
2	0	1		0	0	2
3	1	2		1	1	0
4	3	3		3	1	1

```
● ● ●  
1 # tentukan features dan target untuk CategoricalNB  
2 X_categorical = categorical_df[['Age', 'BMI', 'Blood_Sugar_Level', 'Family_History', 'Diet']]  
3 y_categorical = categorical_df['Diabetes']
```

```
● ● ●  
1 # membagi (split) dataset  
2 X_train_cat, X_test_cat, y_train_cat, y_test_cat = train_test_split(X_categorical, y_categorical, test_size=0.3, random_state=42)
```

```
● ● ●  
1 # Creating and training the Categorical Naive Bayes model  
2 cat_nb_model = CategoricalNB()  
3 cat_nb_model.fit(X_train_cat, y_train_cat)
```

```
● ● ●  
1 # Predicting on the test set  
2 y_pred_cat = cat_nb_model.predict(X_test_cat)
```

### 3. Simulasi Algoritma Naïve Bayes

```
● ● ●  
1 # Simulating with new input data (categorical values)  
2 # Example: Age 45, BMI 27.5, Blood Sugar Level 125, Family History Yes, Diet Poor  
3 # Convert these values to categorical codes  
4 new_data_categorical = np.array([[3, 2, 1, 1, 0]])  
5 new_prediction_cat = cat_nb_model.predict(new_data_categorical)
```

```
● ● ●  
1 (new_prediction_cat)
```



```
1 if (new_prediction_cat[0]==0):
2     print ('Pasien tidak terkena diabetes')
3 else:
4     print ('Pasien terkena diabetes')
```

#### 4. Evaluasi Algoritma Naïve Bayes



```
1 # Calculating accuracy for Categorical Naive Bayes
2 cat_accuracy = accuracy_score(y_test_cat, y_pred_cat)
```



```
1 # Generating the confusion matrix
2 cat_conf_matrix = confusion_matrix(y_test_cat, y_pred_cat)
```



```
1 (cat_accuracy, cat_conf_matrix)
```



```
1 # Menyimpan model menggunakan pickle
2 import pickle
3
4 filename = 'cat_diabetes_mod.pkl'
5 pickle.dump(cat_nb_model, open(filename, 'wb'))
```

## 5. Aplikasi Menggunakan algoritma naïve bayes

```
● ● ●
1 import streamlit as st
2 import pickle
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 # Load the model from the .pkl file
7 with open('cat_diabetes_mod.pkl', 'rb') as model_file:
8     model = pickle.load(model_file)
9
10 # Streamlit app title
11 st.title("Categorical Naive Bayes - Diabetes Prediction")
12
13 # Sidebar inputs for user data
14 st.sidebar.header("Input Data Baru")
15 age = st.sidebar.selectbox("Age", ["<30", "30-40", "40-50", "50-60", "60+"])
16 bmi = st.sidebar.selectbox("BMI", ["<18.5", "18.5-24.9", "25-29.9", "30+"])
17 blood_sugar = st.sidebar.selectbox("Blood Sugar Level", ["<100", "100-125", "125-150", "150+"])
18 family_history = st.sidebar.selectbox("Family History of Diabetes", ["No", "Yes"])
19 diet = st.sidebar.selectbox("Diet", ["Good", "Moderate", "Poor"])
20
21 # Map inputs to numerical values (based on categories used for CategoricalNB)
22 age_map = {"<30": 0, "30-40": 1, "40-50": 2, "50-60": 3, "60+": 4}
23 bmi_map = {"<18.5": 0, "18.5-24.9": 1, "25-29.9": 2, "30+": 3}
24 blood_sugar_map = {"<100": 0, "100-125": 1, "125-150": 2, "150+": 3}
25 family_history_map = {"No": 0, "Yes": 1}
26 diet_map = {"Good": 2, "Moderate": 1, "Poor": 0}
27
28 # Convert user input to model input format
29 input_data = np.array([[age_map[age], bmi_map[bmi], blood_sugar_map[blood_sugar],
30                         family_history_map[family_history], diet_map[diet]]])
31
32 # Make prediction
33 prediction = model.predict(input_data)
34 prediction_prob = model.predict_proba(input_data)
35
36 # Display results
37 st.subheader("Prediction Result")
38 result = "Diabetes (Yes)" if prediction[0] == 1 else "No Diabetes"
39 st.write(f"The model predicts: **{result}**")
40
41 # Plot prediction probabilities
42 fig, ax = plt.subplots()
43 labels = ['No Diabetes', 'Diabetes']
44 ax.bar(labels, prediction_prob[0], color=['green', 'red'])
45 ax.set_ylabel('Probability')
46 ax.set_title('Prediction Probability')
47 st.pyplot(fig)
48
49 # Display input summary
50 st.subheader("Input Summary")
51 st.write(f"**Age**: {age}")
52 st.write(f"**BMI**: {bmi}")
53 st.write(f"**Blood Sugar Level**: {blood_sugar}")
54 st.write(f"**Family History**: {family_history}")
55 st.write(f"**Diet**: {diet}")
```

Nama Dosen	: Teguh Iman Hermanto, M.Kom
Mata Kuliah	: Machine Learning 1
Pembahasan	: Naïve Bayes
Pokok Pemb	<ul style="list-style-type: none"> <li>- Konsep Algoritma Naïve Bayes</li> <li>- Membangun Model Naïve Bayes</li> <li>- Simulasi Algoritma Naïve Bayes</li> <li>- Evaluasi Algoritma Naïve Bayes</li> <li>- Aplikasi menggunakan algoritma Naïve Bayes</li> </ul>

### 1. Konsep Algoritma Naïve Bayes

Naïve Bayes Classifier merupakan sebuah metoda klasifikasi yang berakar pada teorema Bayes . Metode pengklasifikasian dg menggunakan metode probabilitas dan statistik yg dikemukakan oleh ilmuwan Inggris Thomas Bayes , yaitu memprediksi peluang di masa depan berdasarkan pengalaman di masa sebelumnya sehingga dikenal sebagai Teorema Bayes . Ciri utama dr Naïve Bayes Classifier ini adalah asumsi yg sangat kuat (naïf) akan independensi dari masing-masing kondisi / kejadian.

Model Algoritma Naïve Bayes

$$P(H|X)=\frac{P(X|H)P(H)}{P(X)}=P(X|H)\times P(H)/P(X)$$

x = Data dari class scoring yang belum diketahui

H = Hipotesis data X yang merupakan suatu class yang lebih spesifik

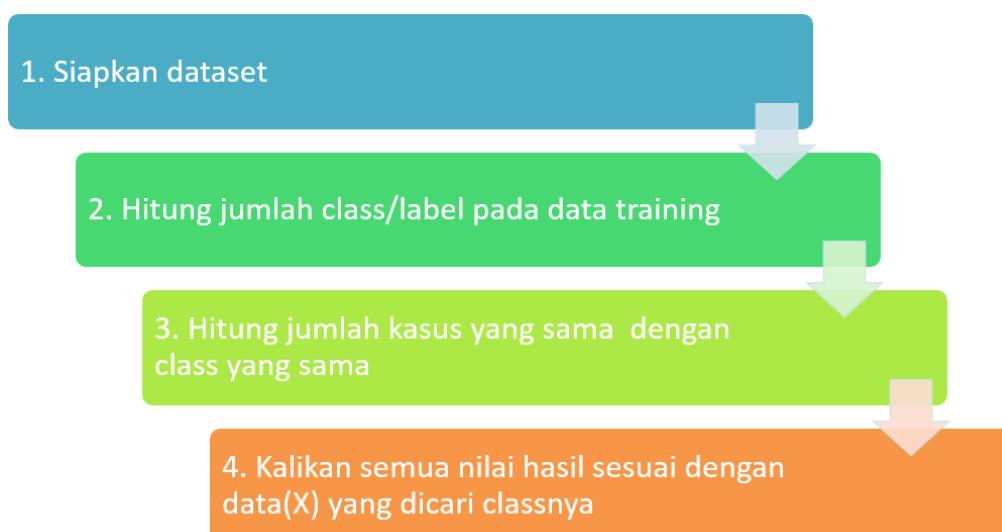
$P(H|X)$  = Probabilitas hipotesis H berdasarkan kondisi X

$P(H)$  = Probabilitas hipotesis H

$P(X|H)$  = Probabilitas hipotesis X berdasarkan kondisi H

$P(X)$  = Probabilitas X

Tahapan Algoritma Naïve Bayes



**a. Siapkan Dataset**

- Data training buys\_computer dataset

Age	Income	Student	Credit Rating	Buys_Computer
youth	high	no	fair	no
youth	high	no	excellent	no
middle_aged	high	no	fair	yes
senior	medium	no	fair	yes
senior	low	yes	fair	yes
senior	low	yes	excellent	no
middle_aged	low	yes	excellent	yes
youth	medium	no	fair	no
youth	low	yes	fair	yes
senior	medium	yes	fair	yes
youth	medium	yes	excellent	yes
middle_aged	medium	no	excellent	yes
middle_aged	high	yes	fair	yes
senior	medium	no	excellent	no

- Data testing

Age	Income	Student	Credit Rating	Buys_Computer
Senior	high	no	fair	?

Apakah calon pembeli akan membeli komputer?

**b. Hitung jumlah class/label pada data training**

Class 1 (C1) → Buys\_Computer yes → 9

Class 2 (C2) → Buys\_Computer no → 5

Total Record: 14

Maka:

$$P(C1) = 9/14 = 0.642857143$$

$$P(C2) = 5/14 = 0.357142857$$

**c. Hitung jumlah kasus yang sama dengan class yang sama**

Contoh penghitungan  $P(X|Ci)$  pada kasus Income dimana

i=1 dan i=2

Maka :

$$P(\text{Income}=\text{"high"} | \text{Buys_Computer}=\text{"yes"}) = 2/9 = 0.222222222$$

$$P(\text{Income}=\text{"high"} | \text{Buys_Computer}=\text{"no"}) = 2/5 = 0.4$$

d. Kalikan semua nilai hasil sesuai dengan data(X) yang dicari classnya

Jika semua data aribut pada data training dihitung, maka:

Atribut	Parameter	yes	no
Age	youth	0.222222222	0.6
Age	middle_aged	0.444444444	0
<b>Age</b>	<b>senior</b>	<b>0.333333333</b>	<b>0.4</b>
Income	low	0.333333333	0.2
Income	medium	0.444444444	0.4
<b>Income</b>	<b>high</b>	<b>0.222222222</b>	<b>0.4</b>
Student	yes	0.666666667	0.2
<b>Student</b>	<b>no</b>	<b>0.333333333</b>	<b>0.8</b>
<b>Credit Rating</b>	<b>fair</b>	<b>0.666666667</b>	<b>0.4</b>
Credit Rating	excellent	0.333333333	0.6

$$P(X| \text{Buys\_Computer}=\text{"yes"}) = 0.016460905$$

$$P(X| \text{Buys\_Computer}=\text{"no"}) = 0.0512$$

$$P(X| \text{Buys\_Computer}=\text{"yes"})*P(C1) = 0.010582011$$

$$P(X| \text{Buys\_Computer}=\text{"no"})*P(C2) = 0.018285714$$

Nilai "no" lebih besar daripada nilai "yes". Sehingga pembeli dengan atribut X tidak membeli komputer

## 2. Membangun Model Naïve Bayes



```

1 import pandas as pd
2 import numpy as np
3
4 from sklearn.model_selection import train_test_split
5 from sklearn.naive_bayes import CategoricalNB
6 from sklearn.metrics import accuracy_score, confusion_matrix
7
8 import matplotlib.pyplot as plt
9 import seaborn as sns

```



```
1 df = pd.read_csv('Diabetes1.csv')
```

	Age	BMI	Blood_Sugar_Level	Family_History	Diet	Diabetes
0	45	28.7	135	Yes	Poor	Yes
1	50	31.2	145	No	Moderate	Yes
2	30	22.0	95	No	Good	No
3	35	25.4	105	Yes	Poor	Yes
4	60	33.5	155	Yes	Poor	Yes



```
1 df.head()
```



```
1 df.info()
```



```
1 # 1. Bar Chart: Jumlah pasien diabetes berdasarkan jenis diet
2 plt.figure(figsize=(8, 5))
3 sns.countplot(x='Diet', hue='Diabetes', data=df)
4 plt.title("Jumlah pasien diabetes berdasarkan jenis diet")
5 plt.xlabel("Diet Type")
6 plt.ylabel("Count")
7 plt.legend(title="Diabetes")
8 plt.show()
```

```
● ● ●  
1 # 2. Pie Chart: Persentase pasien dengan riwayat keluarga  
2 family_history_counts = df['Family_History'].value_counts()  
3 plt.figure(figsize=(5, 5))  
4 family_history_counts.plot.pie(autopct='%1.1f%%', startangle=140, colors=['#ff9999', '#66b3ff'])  
5 plt.title("Persentase pasien dengan riwayat keluarga")  
6 plt.ylabel("")  
7 plt.show()
```

```
● ● ●  
1 # 3. Line Chart: Tren kadar gula darah berdasarkan usia  
2 df_sorted = df.sort_values(by='Age')  
3 plt.figure(figsize=(8, 5))  
4 plt.plot(df_sorted['Age'], df_sorted['Blood_Sugar_Level'], marker='o')  
5 plt.title("Tren kadar gula darah berdasarkan usia")  
6 plt.xlabel("Age")  
7 plt.ylabel("Blood Sugar Level (mg/dL)")  
8 plt.grid(True)  
9 plt.show()
```

```
● ● ●  
1 # 4. Scatter Plot: Hubungan antara BMI dan Kadar Gula Darah  
2 plt.figure(figsize=(8, 5))  
3 sns.scatterplot(x='BMI', y='Blood_Sugar_Level', hue='Diabetes', data=df)  
4 plt.title("Hubungan antara BMI dan Kadar Gula Darah")  
5 plt.xlabel("BMI")  
6 plt.ylabel("Blood Sugar Level (mg/dL)")  
7 plt.legend(title="Diabetes")  
8 plt.show()
```

```
● ● ●  
1 # 5. Histogram: Distribusi BMI  
2 plt.figure(figsize=(8, 5))  
3 plt.hist(df['BMI'], bins=5, color='skyblue', edgecolor='black')  
4 plt.title("Distribusi BMI")  
5 plt.xlabel("BMI")  
6 plt.ylabel("Frequency")  
7 plt.show()
```

```
● ● ●  
1 # 6. Box Plot: Tingkat Gula Darah berdasarkan Riwayat Keluarga  
2 plt.figure(figsize=(8, 5))  
3 sns.boxplot(x='Family_History', y='Blood_Sugar_Level', data=df)  
4 plt.title("Tingkat Gula Darah berdasarkan Riwayat Keluarga")  
5 plt.xlabel("Family History")  
6 plt.ylabel("Blood Sugar Level (mg/dL)")  
7 plt.show()
```

```
● ● ●  
1 # 7. Heatmap: Korelasi antar variabel numerik  
2 plt.figure(figsize=(8, 5))  
3 sns.heatmap(df[['Age', 'BMI', 'Blood_Sugar_Level']].corr(), annot=True, cmap='coolwarm')  
4 plt.title("Correlation Heatmap")  
5 plt.show()
```

```
● ● ●  
1 # 8. Area Chart: Kasus diabetes kumulatif berdasarkan usia (untuk tujuan ilustrasi)  
2 age_cum_cases = df[df['Diabetes'] == 'Yes'].groupby('Age').size().cumsum()  
3 plt.figure(figsize=(8, 5))  
4 plt.fill_between(age_cum_cases.index, age_cum_cases.values, color='lightgreen', alpha=0.6)  
5 plt.title("Kasus Diabetes Kumulatif Berdasarkan Usia")  
6 plt.xlabel("Age")  
7 plt.ylabel("Cumulative Cases")  
8 plt.grid(True)  
9 plt.show()
```

```
● ● ●  
1 # Membuat DataFrame  
2 df_encoded = pd.DataFrame(df)
```

```
● ● ●  
1 # Encode kolom categorical  
2 df_encoded['Family_History'] = df_encoded['Family_History'].map({'Yes': 1, 'No': 0})  
3 df_encoded['Diet'] = df_encoded['Diet'].map({'Poor': 0, 'Moderate': 1, 'Good': 2})  
4 df_encoded['Diabetes'] = df_encoded['Diabetes'].map({'Yes': 1, 'No': 0})
```



```
1 df_encoded.head()
```

	<b>Age</b>	<b>BMI</b>	<b>Blood_Sugar_Level</b>	<b>Family_History</b>	<b>Diet</b>	<b>Diabetes</b>
0	45	28.7	135	1	0	1
1	50	31.2	145	0	1	1
2	30	22.0	95	0	2	0
3	35	25.4	105	1	0	1
4	60	33.5	155	1	0	1

```
● ● ●
1 # mempersiapkan data untuk Categorical Naive Bayes
2 categorical_df = df_encoded.copy()
3 categorical_df['Age'] = pd.cut(categorical_df['Age'], bins=[0, 30, 40, 50, 60, 80], labels=[0, 1, 2, 3, 4])
4 categorical_df['BMI'] = pd.cut(categorical_df['BMI'], bins=[0, 18.5, 24.9, 29.9, 50], labels=[0, 1, 2, 3])
5 categorical_df['Blood_Sugar_Level'] = pd.cut(categorical_df['Blood_Sugar_Level'], bins=[0, 100, 125, 150, 200], labels=[0, 1, 2, 3])
```



```
1 categorical_df.head()
```

	<b>Age</b>	<b>BMI</b>	<b>Blood_Sugar_Level</b>	<b>Family_History</b>	<b>Diet</b>	<b>Diabetes</b>
0	2	2	2	1	0	1
1	2	3	2	0	1	1
2	0	1	0	0	2	0
3	1	2	1	1	0	1
4	3	3	3	1	0	1

```
● ● ●  
1 # tentukan features dan target untuk CategoricalNB  
2 X_categorical = categorical_df[['Age', 'BMI', 'Blood_Sugar_Level', 'Family_History', 'Diet']]  
3 y_categorical = categorical_df['Diabetes']
```

```
● ● ●  
1 # membagi (split) dataset  
2 X_train_cat, X_test_cat, y_train_cat, y_test_cat = train_test_split(X_categorical, y_categorical, test_size=0.3, random_state=42)
```

```
● ● ●  
1 # Creating and training the Categorical Naive Bayes model  
2 cat_nb_model = CategoricalNB()  
3 cat_nb_model.fit(X_train_cat, y_train_cat)
```

```
● ● ●  
1 # Predicting on the test set  
2 y_pred_cat = cat_nb_model.predict(X_test_cat)
```

### 3. Simulasi Algoritma Naïve Bayes

```
● ● ●  
1 # Simulating with new input data (categorical values)  
2 # Example: Age 45, BMI 27.5, Blood Sugar Level 125, Family History Yes, Diet Poor  
3 # Convert these values to categorical codes  
4 new_data_categorical = np.array([[3, 2, 1, 1, 0]])  
5 new_prediction_cat = cat_nb_model.predict(new_data_categorical)
```

```
● ● ●  
1 (new_prediction_cat)
```



```
1 if (new_prediction_cat[0]==0):
2     print ('Pasien tidak terkena diabetes')
3 else:
4     print ('Pasien terkena diabetes')
```

#### 4. Evaluasi Algoritma Naïve Bayes



```
1 # Calculating accuracy for Categorical Naive Bayes
2 cat_accuracy = accuracy_score(y_test_cat, y_pred_cat)
```



```
1 # Generating the confusion matrix
2 cat_conf_matrix = confusion_matrix(y_test_cat, y_pred_cat)
```



```
1 (cat_accuracy, cat_conf_matrix)
```



```
1 # Menyimpan model menggunakan pickle
2 import pickle
3
4 filename = 'cat_diabetes_mod.pkl'
5 pickle.dump(cat_nb_model, open(filename, 'wb'))
```

## 5. Aplikasi Menggunakan algoritma naïve bayes

```
● ● ●
1 import streamlit as st
2 import pickle
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 # Load the model from the .pkl file
7 with open('cat_diabetes_mod.pkl', 'rb') as model_file:
8     model = pickle.load(model_file)
9
10 # Streamlit app title
11 st.title("Categorical Naive Bayes - Diabetes Prediction")
12
13 # Sidebar inputs for user data
14 st.sidebar.header("Input Data Baru")
15 age = st.sidebar.selectbox("Age", ["<30", "30-40", "40-50", "50-60", "60+"])
16 bmi = st.sidebar.selectbox("BMI", ["<18.5", "18.5-24.9", "25-29.9", "30+"])
17 blood_sugar = st.sidebar.selectbox("Blood Sugar Level", ["<100", "100-125", "125-150", "150+"])
18 family_history = st.sidebar.selectbox("Family History of Diabetes", ["No", "Yes"])
19 diet = st.sidebar.selectbox("Diet", ["Good", "Moderate", "Poor"])
20
21 # Map inputs to numerical values (based on categories used for CategoricalNB)
22 age_map = {"<30": 0, "30-40": 1, "40-50": 2, "50-60": 3, "60+": 4}
23 bmi_map = {"<18.5": 0, "18.5-24.9": 1, "25-29.9": 2, "30+": 3}
24 blood_sugar_map = {"<100": 0, "100-125": 1, "125-150": 2, "150+": 3}
25 family_history_map = {"No": 0, "Yes": 1}
26 diet_map = {"Good": 2, "Moderate": 1, "Poor": 0}
27
28 # Convert user input to model input format
29 input_data = np.array([[age_map[age], bmi_map[bmi], blood_sugar_map[blood_sugar],
30                         family_history_map[family_history], diet_map[diet]]])
31
32 # Make prediction
33 prediction = model.predict(input_data)
34 prediction_prob = model.predict_proba(input_data)
35
36 # Display results
37 st.subheader("Prediction Result")
38 result = "Diabetes (Yes)" if prediction[0] == 1 else "No Diabetes"
39 st.write(f"The model predicts: **{result}**")
40
41 # Plot prediction probabilities
42 fig, ax = plt.subplots()
43 labels = ['No Diabetes', 'Diabetes']
44 ax.bar(labels, prediction_prob[0], color=['green', 'red'])
45 ax.set_ylabel('Probability')
46 ax.set_title('Prediction Probability')
47 st.pyplot(fig)
48
49 # Display input summary
50 st.subheader("Input Summary")
51 st.write(f"**Age**: {age}")
52 st.write(f"**BMI**: {bmi}")
53 st.write(f"**Blood Sugar Level**: {blood_sugar}")
54 st.write(f"**Family History**: {family_history}")
55 st.write(f"**Diet**: {diet}")
```

Nama Dosen : Teguh Iman Hermanto, M.Kom  
 Mata Kuliah : Machine Learning 1  
 Pembahasan : Decision Tree (C.45)  
 Pokok Pemb : - Membangun Model Decision Tree  
               - Simulasi Algoritma Decision Tree  
               - Evaluasi Algoritma Decision Tree  
               - Aplikasi menggunakan algoritma Decision Tree

## 1. Membangun Model Decision Tree

```
● ● ●
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score, confusion_matrix
7 from sklearn.preprocessing import LabelEncoder
8 from sklearn.tree import DecisionTreeClassifier
```

```
● ● ●
1 df = pd.read_csv('heart.csv')
```

```
● ● ●
1 df.head()
```

	<b>Age</b>	<b>Cholesterol Level</b>	<b>Blood Pressure</b>	<b>Smoking</b>	<b>Physical Activity</b>	<b>BMI</b>	<b>Heart Disease Risk</b>
0	55	High	High	Yes	Low	Overweight	High
1	43	Normal	Normal	No	Moderate	Normal	Low
2	60	High	High	Yes	Low	Obese	High
3	35	Normal	Low	No	High	Normal	Low
4	50	High	High	No	Low	Overweight	Medium



```
1 # 1. Histogram of Age
2 plt.figure(figsize=(10, 6))
3 plt.hist(df['Age'], bins=10, edgecolor='black')
4 plt.title("Distribution of Age")
5 plt.xlabel("Age")
6 plt.ylabel("Frequency")
7 plt.show()
```



```
1 # 2. Bar plot of Cholesterol Level Counts
2 plt.figure(figsize=(10, 6))
3 df['Cholesterol Level'].value_counts().plot(kind='bar')
4 plt.title("Count of Cholesterol Levels")
5 plt.xlabel("Cholesterol Level")
6 plt.ylabel("Count")
7 plt.show()
```



```
1 # 3. Bar plot of Blood Pressure Counts
2 plt.figure(figsize=(10, 6))
3 df['Blood Pressure'].value_counts().plot(kind='bar', color='orange')
4 plt.title("Count of Blood Pressure Levels")
5 plt.xlabel("Blood Pressure Level")
6 plt.ylabel("Count")
7 plt.show()
```



```
1 # 4. Pie chart of Smoking Status
2 plt.figure(figsize=(8, 8))
3 df['Smoking'].value_counts().plot(kind='pie', autopct='%1.1f%%', startangle=140)
4 plt.title("Proportion of Smoking Status")
5 plt.ylabel("") # Removes 'Smoking' label from pie
6 plt.show()
```



```
1 # 5. Count plot of BMI categories
2 plt.figure(figsize=(10, 6))
3 df['BMI'].value_counts().plot(kind='bar', color='purple')
4 plt.title("Count of BMI Categories")
5 plt.xlabel("BMI Category")
6 plt.ylabel("Count")
7 plt.show()
```



```
1 # 6. Violin Plot of Age by Cholesterol Level
2 plt.figure(figsize=(10, 6))
3 sns.violinplot(x='Cholesterol Level', y='Age', data=df)
4 plt.title("Violin Plot of Age by Cholesterol Level")
5 plt.xlabel("Cholesterol Level")
6 plt.ylabel("Age")
7 plt.show()
```



```
1 # 7. Swarm Plot of Age by BMI Category
2 plt.figure(figsize=(10, 6))
3 sns.swarmplot(x='BMI', y='Age', data=df)
4 plt.title("Swarm Plot of Age by BMI Category")
5 plt.xlabel("BMI Category")
6 plt.ylabel("Age")
7 plt.show()
```

```
● ● ●  
1 # 8. Box Plot of Age by Blood Pressure  
2 plt.figure(figsize=(10, 6))  
3 sns.boxplot(x='Blood Pressure', y='Age', data=df)  
4 plt.title("Box Plot of Age by Blood Pressure Level")  
5 plt.xlabel("Blood Pressure Level")  
6 plt.ylabel("Age")  
7 plt.show()
```

```
● ● ●  
1 # 9. Scatter Plot of Age vs Physical Activity with BMI Category as Hue  
2 plt.figure(figsize=(10, 6))  
3 sns.scatterplot(data=df, x='Age', y='Physical Activity', hue='BMI')  
4 plt.title("Scatter Plot of Age vs Physical Activity by BMI Category")  
5 plt.xlabel("Age")  
6 plt.ylabel("Physical Activity")  
7 plt.legend(title="BMI Category")  
8 plt.show()
```

```
● ● ●  
1 # membuat fungsi untuk mengkategorikan usia  
2 def categorize_age(age):  
3     if age < 35:  
4         return "Young"  
5     elif age < 55:  
6         return "Middle-aged"  
7     else:  
8         return "Senior"
```

```
● ● ●  
1 # terapkan fungsi pada kolom 'Age (Usia)' dan buat kolom baru 'Age Category'  
2 df['Age'] = df['Age'].apply(categorize_age)
```



```
1 df.head()
```

	Age	Cholesterol Level	Blood Pressure	Smoking	Physical Activity	BMI	Heart Disease Risk
0	Senior	High	High	Yes	Low	Overweight	High
1	Middle-aged	Normal	Normal	No	Moderate	Normal	Low
2	Senior	High	High	Yes	Low	Obese	High
3	Middle-aged	Normal	Low	No	High	Normal	Low
4	Middle-aged	High	High	No	Low	Overweight	Medium



```
1 # Encode semua variabel kategori
2 label_encoders = {}
3 for column in df.columns:
4     if df[column].dtype == 'object':
5         le = LabelEncoder()
6         df[column] = le.fit_transform(df[column])
7         label_encoders[column] = le
```



```
1 df.head()
```

	Age	Cholesterol Level	Blood Pressure	Smoking	Physical Activity	BMI	Heart Disease Risk
0	1	0	0	1		1	2
1	0	1	2	0		2	0
2	1	0	0	1		1	1
3	0	1	1	0		0	0
4	0	0	0	0		1	2



```
1 # Split the dataset
2 X = df.drop('Heart Disease Risk', axis=1)
3 y = df['Heart Disease Risk']
```



```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```



```
1 # Train the C4.5 model (Decision Tree)
2 model = DecisionTreeClassifier(criterion='entropy')
3 model.fit(X_train, y_train)
```



```
1 # Make predictions
2 y_pred = model.predict(X_test)
```

## 2. Simulasi Algoritma Decision Tree

```
● ● ●  
1 # Simulate new input data  
2 new_data = pd.DataFrame({  
3     'Age': [label_encoders['Age'].transform(['Middle-aged'])[0]],  
4     'Cholesterol Level': [label_encoders['Cholesterol Level'].transform(['High'])[0]],  
5     'Blood Pressure': [label_encoders['Blood Pressure'].transform(['High'])[0]],  
6     'Smoking': [label_encoders['Smoking'].transform(['Yes'])[0]],  
7     'Physical Activity': [label_encoders['Physical Activity'].transform(['Low'])[0]],  
8     'BMI': [label_encoders['BMI'].transform(['Overweight'])[0]]  
9 })
```

```
● ● ●  
1 new_prediction = model.predict(new_data)
```

```
● ● ●  
1 # Decode the new prediction  
2 new_prediction_decoded = label_encoders['Heart Disease Risk'].inverse_transform(new_prediction)
```

```
● ● ●  
1 new_prediction_decoded
```

## 3. Evaluasi Algoritma Decision Tree

```
● ● ●  
1 # Calculate accuracy and confusion matrix  
2 accuracy = accuracy_score(y_test, y_pred)  
3 conf_matrix = confusion_matrix(y_test, y_pred)
```

```
● ● ●  
1 accuracy, conf_matrix
```

#### 4. Plot Hasil Decision Tree

```
● ● ●
1 import matplotlib.pyplot as plt
2 from sklearn.tree import plot_tree
3
4 plt.figure(figsize=(20, 20))
5 plot_tree(model, feature_names=['Age',
6                               'Cholesterol Level',
7                               'Blood Pressure',
8                               'Smoking',
9                               'Physical Activity',
10                              'BMI'], class_names=['High Risk',
11                               'Low Risk',
12                               'Medium Risk'],
13                               filled=True, rounded=True)
14 plt.title("pohon keputusan")
15 plt.show()
```

## 5. Aplikasi menggunakan algoritma Decision Tree

```

● ● ●

1 import streamlit as st
2 import pickle
3 import numpy as np
4 import matplotlib.pyplot as plt
5 from sklearn.tree import plot_tree
6
7 # Load the trained model from the .pkl file
8 model_path = 'c45_pinjam_mod.pkl'
9 with open(model_path, 'rb') as model_file:
10     loaded_model = pickle.load(model_file)
11
12 st.title('C4.5 Decision Tree Model for Heart Disease Prediction')
13
14 # Input for new data
15 st.sidebar.header('Input Features')
16 age = st.sidebar.selectbox('Age', ('Senior', 'Middle-aged', 'Young'))
17 cholesterol = st.sidebar.selectbox('Cholesterol Level', ('High', 'Normal', 'Low'))
18 blood_pressure = st.sidebar.selectbox('Blood Pressure', ('High', 'Normal', 'Low'))
19 smoking = st.sidebar.selectbox('Smoking', ('Yes', 'No'))
20 physical_activity = st.sidebar.selectbox('Physical Activity', ('Low', 'Moderate', 'High'))
21 bmi = st.sidebar.selectbox('BMI', ('Overweight', 'Normal', 'Obese'))
22
23 # Create a mapping for the labels to numbers (as per how the model was trained)
24 label_encodings = {
25     'Age': {'Senior': 1, 'Middle-aged': 0, 'Young': 2},
26     'Cholesterol Level': {'High': 0, 'Normal': 1, 'Low': 2},
27     'Blood Pressure': {'High': 0, 'Normal': 1, 'Low': 2},
28     'Smoking': {'Yes': 1, 'No': 0},
29     'Physical Activity': {'Low': 0, 'Moderate': 1, 'High': 2},
30     'BMI': {'Overweight': 1, 'Normal': 0, 'Obese': 2}
31 }
32
33 # Encode the input data
34 new_data = np.array([[label_encodings['Age'][age],
35                      label_encodings['Cholesterol Level'][cholesterol],
36                      label_encodings['Blood Pressure'][blood_pressure],
37                      label_encodings['Smoking'][smoking],
38                      label_encodings['Physical Activity'][physical_activity],
39                      label_encodings['BMI'][bmi]]])
40
41 # Predict using the loaded model
42 if st.button('Predict'):
43     prediction = loaded_model.predict(new_data)
44     risk_mapping = {0: 'High Risk', 1: 'Low Risk', 2: 'Medium Risk'}
45     st.write(f'The predicted heart disease risk is: **{risk_mapping[prediction[0]]}**')
46
47 # Plot the decision tree
48 st.write("### Decision Tree Visualization")
49 fig = plt.figure(figsize=(20, 10))
50 plot_tree(loaded_model, feature_names=['Age',
51                                         'Cholesterol Level',
52                                         'Blood Pressure',
53                                         'Smoking',
54                                         'Physical Activity',
55                                         'BMI'],
56                                         class_names=['High Risk',
57                                         'Low Risk',
58                                         'Medium Risk'],
59                                         filled=True, rounded=True)
60 st.pyplot(fig)

```

Nama Dosen : Teguh Iman Hermanto, M.Kom  
Mata Kuliah : Machine Learning 1  
Pembahasan : K-Means Clustering  
Pokok Pemb : - Membangun Model K-Means Clustering  
- Simulasi Algoritma K-means

## 1. Membangun Model K-means

```
● ● ●  
1 import pandas as pd  
2  
3 from sklearn.cluster import KMeans  
4 from sklearn.preprocessing import StandardScaler  
5 import matplotlib.pyplot as plt  
6  
7 import seaborn as sns  
8 from sklearn.decomposition import PCA
```

```
● ● ●  
1 # Memuat dataset  
2 df = pd.read_csv("pelanggan.csv")
```

```
● ● ●  
1 # Pie chart distribusi gender pelanggan  
2 gender_counts = df['Gender'].value_counts()  
3 plt.pie(gender_counts,  
4          labels=gender_counts.index,  
5          autopct='%1.1f%%',  
6          colors=['skyblue', 'pink'],  
7          startangle=90)  
8 plt.title("Distribusi Gender Pelanggan")  
9 plt.show()
```

```
1 # Boxplot pengeluaran rata-rata per transaksi berdasarkan kategori produk
2 sns.boxplot(data=df,
3               x='Preferred_Product_Category',
4               y='Average_Spending_Per_Transaction',
5               palette='Set3')
6 plt.title("Pengeluaran Rata-Rata Per Transaksi Berdasarkan Kategori Produk")
7 plt.xlabel("Kategori Produk")
8 plt.ylabel("Rata-Rata Pengeluaran Per Transaksi (Rp)")
9 plt.xticks(rotation=45)
10 plt.show()
```

```
1 # Frekuensi pembelian berdasarkan kota
2 sns.barplot(data=df,
3               x='City',
4               y='Purchase_Frequency',
5               estimator=sum,
6               palette='pastel')
7 plt.title("Frekuensi Pembelian Berdasarkan Kota")
8 plt.xlabel("Kota")
9 plt.ylabel("Total Frekuensi Pembelian")
10 plt.show()
```

```
1 # Bar plot kategori produk yang disukai
2 sns.countplot(data=df,
3                x='Preferred_Product_Category',
4                order=df['Preferred_Product_Category'].value_counts().index,
5                palette='coolwarm')
6 plt.title("Kategori Produk yang Disukai Pelanggan")
7 plt.xlabel("Kategori Produk")
8 plt.ylabel("Jumlah Pelanggan")
9 plt.show()
```



```
1 # Distribusi frekuensi pembelian pelanggan
2 sns.histplot(df['Purchase_Frequency'],
3                 kde=True,
4                 color='orange',
5                 bins=15)
6 plt.title("Distribusi Frekuensi Pembelian Pelanggan")
7 plt.xlabel("Frekuensi Pembelian")
8 plt.ylabel("Jumlah Pelanggan")
9 plt.show()
```

## 2. Simulasi Algoritma K-Means



```
1 # Memilih fitur untuk clustering
2 features = ["Annual_Income",
3              "Purchase_Frequency",
4              "Average_Spending_Per_Transaction",
5              "Total_Spent"]
```



```
1 # Standarisasi data
2 scaler = StandardScaler()
3 scaled_features = scaler.fit_transform(df[features])
```



```
1 # Menentukan jumlah cluster menggunakan metode Elbow
2 inertia = []
3 k_values = range(1, 11)
4
5 for k in k_values:
6     kmeans = KMeans(n_clusters=k, random_state=42)
7     kmeans.fit(scaled_features)
8     inertia.append(kmeans.inertia_)
```



```
1 # Plot Elbow Method
2 plt.figure(figsize=(8, 5))
3 plt.plot(k_values, inertia, marker='o')
4 plt.title("Elbow Method for Optimal k")
5 plt.xlabel("Number of Clusters (k)")
6 plt.ylabel("Inertia")
7 plt.grid()
8 plt.show()
```



```
1 # Menggunakan jumlah cluster optimal
2 optimal_k = 4 # Berdasarkan elbow method (misalnya)
3 kmeans = KMeans(n_clusters=optimal_k, random_state=42)
4 df['Cluster'] = kmeans.fit_predict(scaled_features)
```



```
1 # Mengurangi dimensi data untuk visualisasi
2 pca = PCA(n_components=2)
3 reduced_features = pca.fit_transform(scaled_features)
```



```
1 # Warna untuk setiap cluster
2 colors = ['blue', 'green', 'red', 'purple']
```



```
1 # Plot hasil clustering
2 plt.figure(figsize=(10, 7))
3 for cluster in range(optimal_k):
4     cluster_points = reduced_features[df['Cluster'] == cluster]
5     plt.scatter(cluster_points[:, 0],
6                 cluster_points[:, 1],
7                 s=50,
8                 label=f'Cluster {cluster}',
9                 color=colors[cluster])
10
11 # Menambahkan centroid pada plot
12 centroids_reduced = pca.transform(kmeans.cluster_centers_)
13 plt.scatter(centroids_reduced[:, 0],
14             centroids_reduced[:, 1],
15             s=200,
16             c='black',
17             marker='X',
18             label='Centroids')
19
20 # Menambahkan detail plot
21 plt.title("Hasil Clustering K-Means (2D PCA Projection)")
22 plt.xlabel("Principal Component 1")
23 plt.ylabel("Principal Component 2")
24 plt.legend()
25 plt.grid()
26 plt.show()
```

Nama Dosen : Teguh Iman Hermanto, M.Kom  
 Mata Kuliah : Machine Learning 1  
 Pembahasan : Hierarchical Clustering  
 Pokok Pemb : - Membangun Model Hierarchical Clustering  
               - Simulasi Hierarchical Clustering  
               - Profiling Hasil CLustering

### 1. Load Library dan dataset pada file Notebook

```

● ● ●

1 import pandas as pd
2
3 from sklearn.preprocessing import StandardScaler
4 from scipy.cluster.hierarchy import linkage
5 from scipy.cluster.hierarchy import dendrogram
6 from scipy.cluster.hierarchy import cut_tree
7
8 import matplotlib.pyplot as plt
9 import seaborn as sns

```

```

● ● ●

1 # Membaca dataset
2 data = pd.read_csv('dataset_nasabah.csv')

```

```

data.head()
✓ 0.0s
Python

```

ID Nasabah	Usia	Jenis Kelamin	Status Merokok	BMI	Jumlah Klaim	Total Biaya Klaim
0	1	Pria	Ya	30.7	1	4744854
1	2	Pria	Ya	28.2	0	0
2	3	Pria	Ya	22.3	4	5780796
3	4	Pria	Tidak	22.1	4	22669060
4	5	Pria	Ya	30.0	2	11323814

2. Buatkan 5 variasi Exploratory data analysis dan berikan penjelasannya yang terdiri dari

- a. Pie Plot
- b. Box plot
- c. Violin plot
- d. Scatter plot
- e. Count plot

3. Buat model Hierarchical Clustering

```
● ● ●
1 # Pra-pemrosesan data
2 # Mengonversi kolom kategorikal menjadi numerik
3 data['Jenis Kelamin'] = data['Jenis Kelamin'].map({'Pria': 0, 'Wanita': 1})
4 data['Status Merokok'] = data['Status Merokok'].map({'Tidak': 0, 'Ya': 1})
```

`data.head()`

✓ 0.0s Python

	ID Nasabah	Usia	Jenis Kelamin	Status Merokok	BMI	Jumlah Klaim	Total Biaya Klaim
0	1	58	0	1	30.7	1	4744854
1	2	26	0	1	28.2	0	0
2	3	19	0	1	22.3	4	5780796
3	4	53	0	0	22.1	4	22669060
4	5	69	0	1	30.0	2	11323814

```
● ● ●
1 # Select features for clustering
2 features = ['Usia', 'Jenis Kelamin', 'Status Merokok', 'BMI', 'Jumlah Klaim', 'Total Biaya Klaim']
3 X = data[features]
```

```
● ● ●
1 # Standardize the data
2 scaler = StandardScaler()
3 X_scaled = scaler.fit_transform(X)
```

```
● ● ●
1 # Create linkage matrix for dendrogram
2 linked = linkage(X_scaled, method='ward')
```

```
● ● ●
1 # Plot dendrogram
2 plt.figure(figsize=(10, 7))
3 plt.title("Dendrogram for Hierarchical Clustering")
4 dendrogram(linked, truncate_mode='lastp', p=30, leaf_rotation=90, leaf_font_size=10)
5 plt.xlabel("Cluster Size")
6 plt.ylabel("Distance")
7 plt.show()
```

```
● ● ●
1 # membuat cluster label
2 cluster_labels = cut_tree(linked, n_clusters=5).reshape(-1, )
3 cluster_labels
```

```
● ● ●
1 # gabungkan cluster label dengan data
2 data['Cluster_Labels'] = cluster_labels
```

```
data.head()
```

✓ 0.0s

Python

ID Nasabah	Usia	Jenis Kelamin	Status Merokok	BMI	Jumlah Klaim	Total Biaya Klaim	Cluster_Labels
0	1	58	0	1	30.7	1	4744854
1	2	26	0	1	28.2	0	0
2	3	19	0	1	22.3	4	5780796
3	4	53	0	0	22.1	4	22669060
4	5	69	0	1	30.0	2	11323814

**4. Buatkan profiling masing-masing cluster**

- a. Buatkan profiling menggunakan scatterplot, lakukan perbandingan dengan merubah nilai x dan y lalu berikan kesimpulannya

```
● ● ●  
1 plt.figure()  
2 sns.scatterplot(data=data,  
3                   x='Usia',  
4                   y='Total Biaya Klaim',  
5                   hue='Cluster_Labels',  
6                   palette='viridis')  
7 plt.title("Hubungan Usia dan Total Biaya Klaim")  
8 plt.xlabel("Usia")  
9 plt.ylabel("Total Biaya Klaim")  
10 plt.show()
```

- b. Buatkan profiling menggunakan boxplot, lakukan perbandingan dengan merubah nilai y lalu berikan kesimpulannya

```
● ● ●  
1 sns.boxplot(data=data,  
2               x='Cluster_Labels',  
3               y='Usia',  
4               palette='Set3')  
5 plt.title("rata-rata usia berdasarkan cluster")  
6 plt.xlabel("Cluster")  
7 plt.ylabel("Usia")  
8 plt.show()
```

Nama Dosen : Teguh Iman Hermanto, M.Kom  
Mata Kuliah : Machine Learning 1  
Pembahasan : Density-Based Spatial Clustering of Applications with Noise (DBSCAN)  
Pokok Pemb : - Membangun Model DBSCAN  
- Simulasi Algoritma DBSCAN  
- Profiling Hasil Clustering DBSCAN

### 1. Load Library dan dataset pada file Notebook

```
1 import pandas as pd
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.cluster import DBSCAN
4 import matplotlib.pyplot as plt
5 import seaborn as sns
```

```
1 df = pd.read_csv('net_detection.csv')
```

df.head()

✓ 0.0s

	src_ip	dst_ip	protocol	src_port	dst_port	bytes_sent	bytes_received	duration	packet_count	attack_type
0	39.170.115.188	133.204.219.238	TCP	1762	62458	3422	5989	213	131	normal
1	80.35.125.105	246.113.106.207	TCP	32718	9699	3736	989	277	96	normal
2	49.134.137.30	151.26.62.67	TCP	1225	43970	2865	5943	305	89	ddos
3	157.51.229.193	175.153.3.55	TCP	20804	303	1852	9389	552	80	normal
4	121.123.112.174	72.234.63.118	UDP	15457	17942	8318	5160	533	46	normal

## 2. Exploratory data analysis

### a. Menghitung jumlah data yang ada pada kolom protocol

```
● ● ●  
1 # menghitung jumlah data yang ada pada kolom protocol  
2 protocol_counts = df['protocol'].value_counts()  
3 protocol_counts
```

```
● ● ●  
1 # membuat bar plot hasil perhitungan kolom protocol  
2 plt.figure(figsize=(10, 6))  
3 protocol_counts.plot(kind='bar')  
4 plt.title('Number of Data Points per Protocol')  
5 plt.xlabel('Protocol')  
6 plt.ylabel('Count')  
7 plt.show()
```

Lakukan perhitungan untuk mendapatkan jumlah pada kolom “attack\_type” dan buatkan bar plot nya

### b. Membuat box plot untuk menampilkan rata-rata nilai source plot berdasarkan jenis attack type

```
● ● ●  
1 plt.figure(figsize=(12, 6))  
2 sns.boxplot(x='attack_type', y='src_port', data=df)  
3 plt.title('Source Port Distribution by Attack Type')  
4 plt.xlabel('Attack Type')  
5 plt.ylabel('Source Port')  
6 plt.show()
```

Buatkan boxplot untuk menampilkan bytes\_sent dan berdasarkan senis protocol

### 3. Buat model DBSCAN

```
● ● ●
1 # pilih fitur untuk clustering
2 features = ['src_port','dst_port','bytes_sent','bytes_received']
3 X = df[features]
```

```
● ● ●
1 # Standardize the data
2 scaler = StandardScaler()
3 X_scaled = scaler.fit_transform(X)
```

```
● ● ●
1 # implementasi algoritma DBSCAN
2 # tentukan nilai epsilon dan min sample
3 dbSCAN = DBSCAN(eps=0.5, min_samples=5)
4 df['dbSCAN_cluster'] = dbSCAN.fit_predict(X_scaled)
```

```
● ● ●
1 # Plot hasil cluster
2 plt.figure(figsize=(10, 6))
3 sns.scatterplot(x='bytes_sent', y='bytes_received',
4                  hue='dbSCAN_cluster', data=df,
5                  palette='viridis')
6 plt.title('DBSCAN Clustering Results')
7 plt.show()
```

	df.head()										
	0.0s										
	src_ip	dst_ip	protocol	src_port	dst_port	bytes_sent	bytes_received	duration	packet_count	attack_type	dbSCAN_cluster
0	39.170.115.188	133.204.219.238	TCP	1762	62458	3422	5989	213	131	normal	-1
1	80.35.125.105	246.113.106.207	TCP	32718	9699	3736	989	277	96	normal	0
2	49.134.137.30	151.26.62.67	TCP	1225	43970	2865	5943	305	89	ddos	-1
3	157.51.229.193	175.153.3.55	TCP	20804	303	1852	9389	552	80	normal	-1
4	121.123.112.174	72.234.63.118	UDP	15457	17942	8318	5160	533	46	normal	0



```
1 # memisahkan data anomaly yang bernilai -1
2 df['anomaly'] = df['dbscan_cluster'] == -1
```



```
1 # Plot hasil deteksi anomaly
2 plt.figure(figsize=(10, 6))
3 sns.scatterplot(x='bytes_sent', y='bytes_received',
4                  hue='anomaly', data=df,
5                  palette=['blue', 'red'])
6 plt.title('Anomaly Detection')
7 plt.show()
```

df.head()

	src_ip	dst_ip	protocol	src_port	dst_port	bytes_sent	bytes_received	duration	packet_count	attack_type	dbscan_cluster	anomaly
0	39.170.115.188	133.204.219.238	TCP	1762	62458	3422	5989	213	131	normal	-1	True
1	80.35.125.105	246.113.106.207	TCP	32718	9699	3736	989	277	96	normal	0	False
2	49.134.137.30	151.26.62.67	TCP	1225	43970	2865	5943	305	89	ddos	-1	True
3	157.51.229.193	175.153.3.55	TCP	20804	303	1852	9389	552	80	normal	-1	True
4	121.123.112.174	72.234.63.118	UDP	15457	17942	8318	5160	533	46	normal	0	False



```
1 # simpan hasil cluster
2 df.to_csv('hasil_dbscan.csv', index=False)
```

**4. Buatkan profiling hasil cluster**

```
1 # membuat df untuk type serangan ddos  
2 ddos_df = df[df['attack_type'] == 'ddos']
```



```
1 # hitung jumlah anomalynya  
2 anomaly_counts = ddos_df['anomaly'].value_counts()  
3 anomaly_counts
```

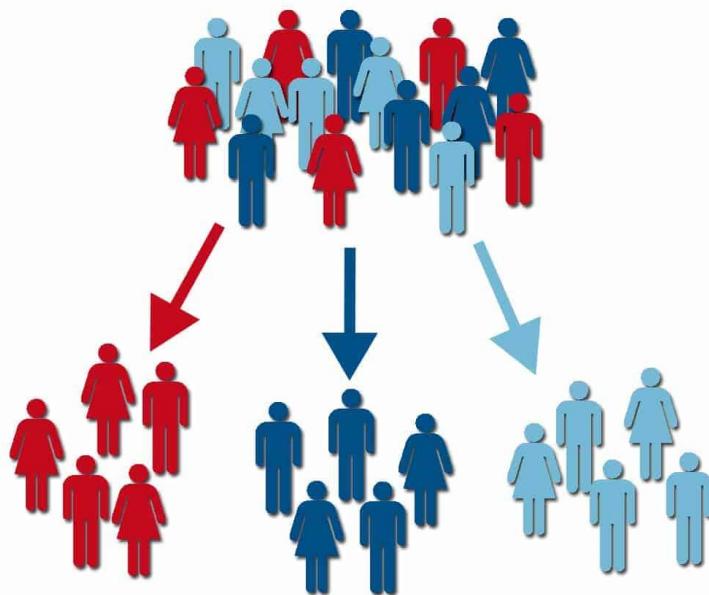


```
1 # bar plot untuk serangan DDOS  
2 plt.figure(figsize=(8, 6))  
3 anomaly_counts.plot(kind='bar')  
4 plt.title('Jumlah Anomaly untuk tipe DDoS')  
5 plt.xlabel('Anomaly')  
6 plt.ylabel('Count')  
7 plt.show()
```

Lakukan analisis untuk deteksi anomaly pada jenis serangan lainnya seperti brute force, icmp flood, dan port scan.

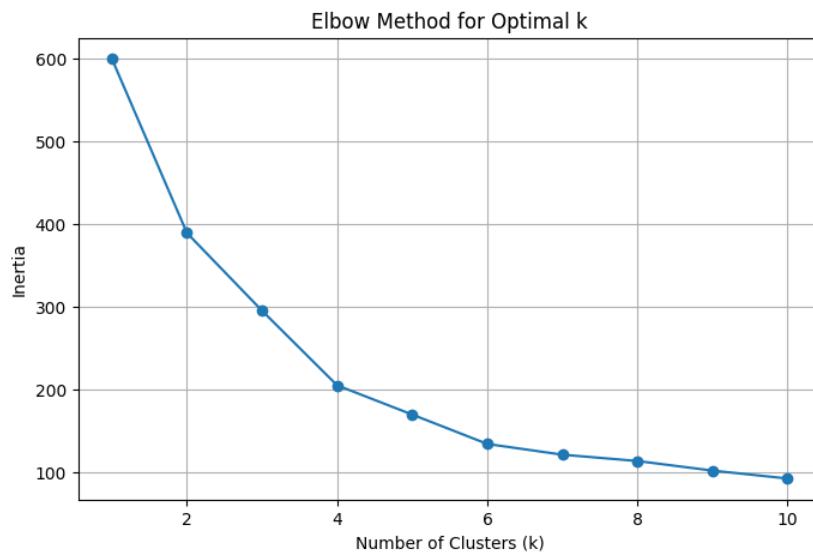
Nama Dosen : Teguh Iman Hermanto, M.Kom  
Mata Kuliah : Machine Learning 1  
Pembahasan : Komparasi Algoritma Clustering  
Pokok Pemb :  
- Membangun Model K-Means  
- Membangun Model Hierarchical Clustering  
- Membangun Model DBSCAN

### 1. Load dataset pada file Notebook

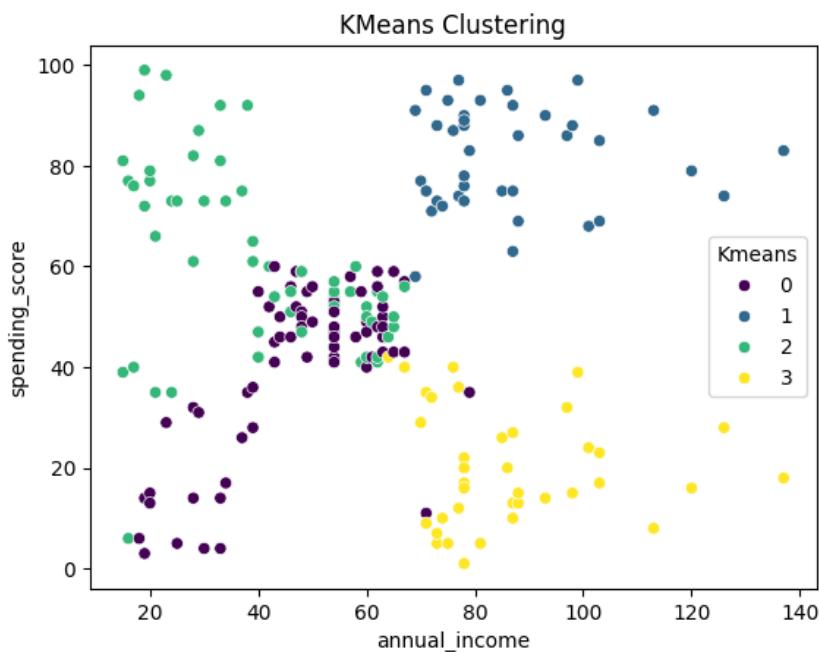


	customer_id	gender	age	annual_income	spending_score
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40

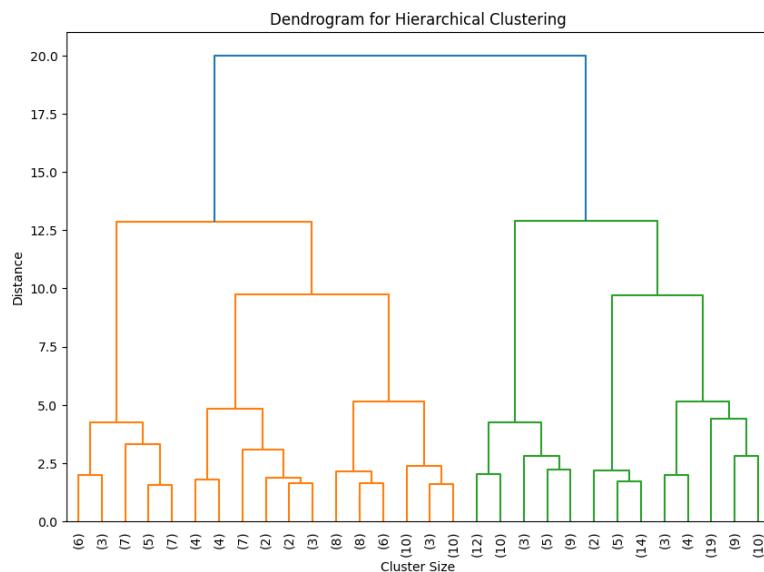
## 2. Membangun Model K-Means



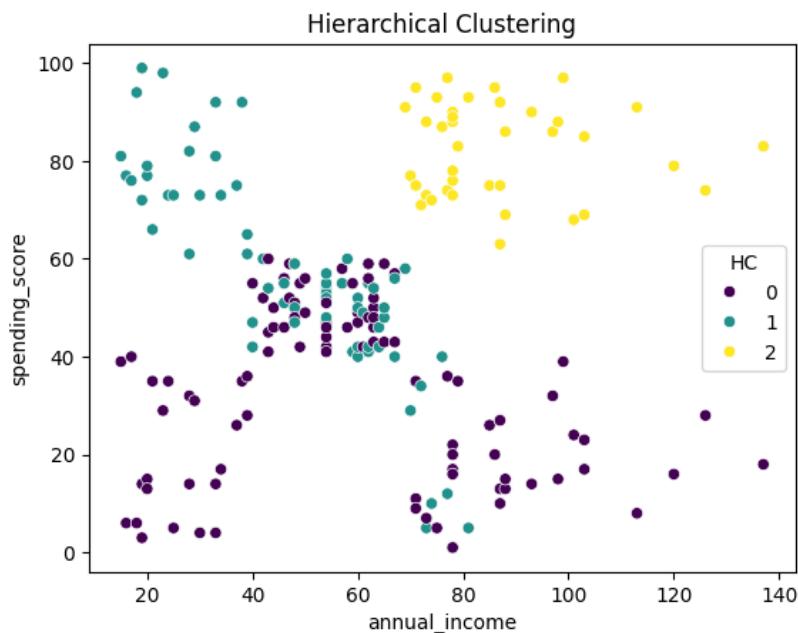
	customer_id	gender	age	annual_income	spending_score	Kmeans
0	1	Male	19		15	39
1	2	Male	21		15	81
2	3	Female	20		16	6
3	4	Female	23		16	77
4	5	Female	31		17	40



### 3. Membangun model Hierarchical Clustering

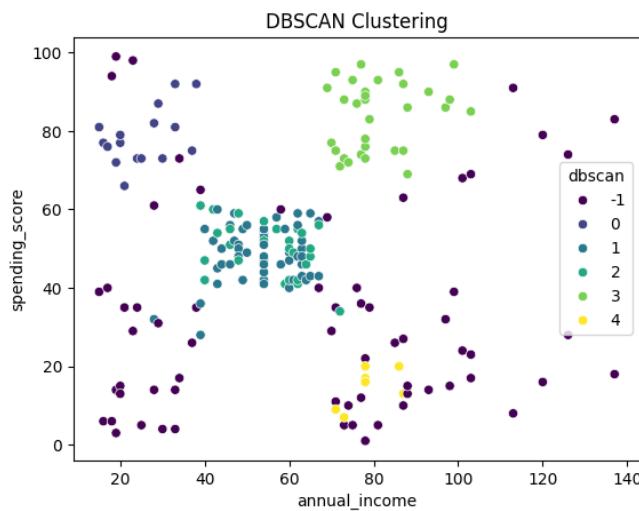


	customer_id	gender	age	annual_income	spending_score	Kmeans	HC
0	1	Male	19		15	39	2 0
1	2	Male	21		15	81	2 1
2	3	Female	20		16	6	2 0
3	4	Female	23		16	77	2 1
4	5	Female	31		17	40	2 0

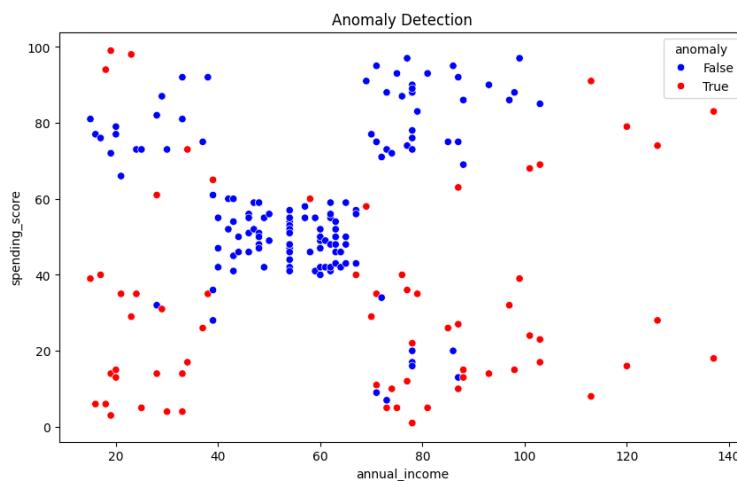


#### 4. Membangun Model DBSCAN

	customer_id	gender	age	annual_income	spending_score	Kmeans	HC	dbscan
0	1	Male	19		15	39	2	0
1	2	Male	21		15	81	2	1
2	3	Female	20		16	6	2	0
3	4	Female	23		16	77	2	1
4	5	Female	31		17	40	2	-1

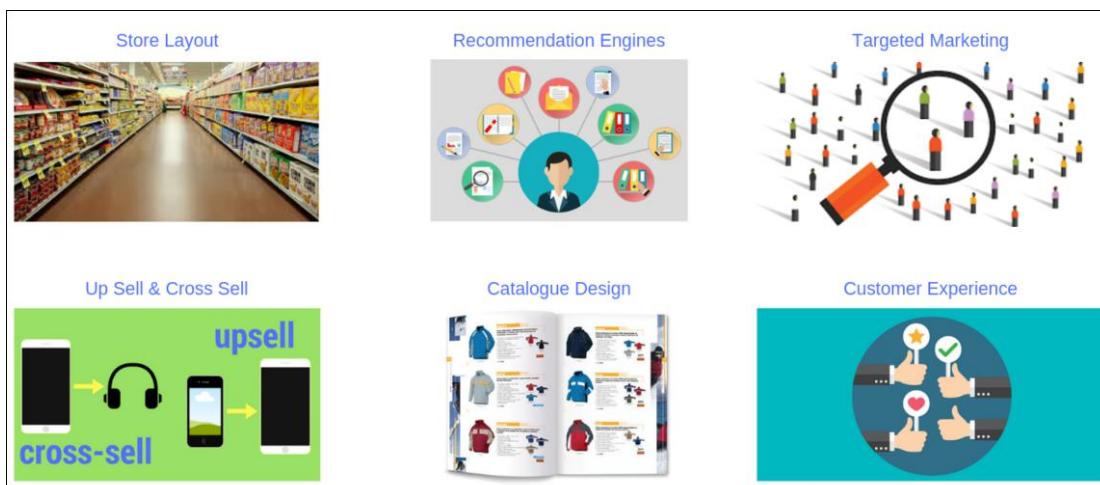
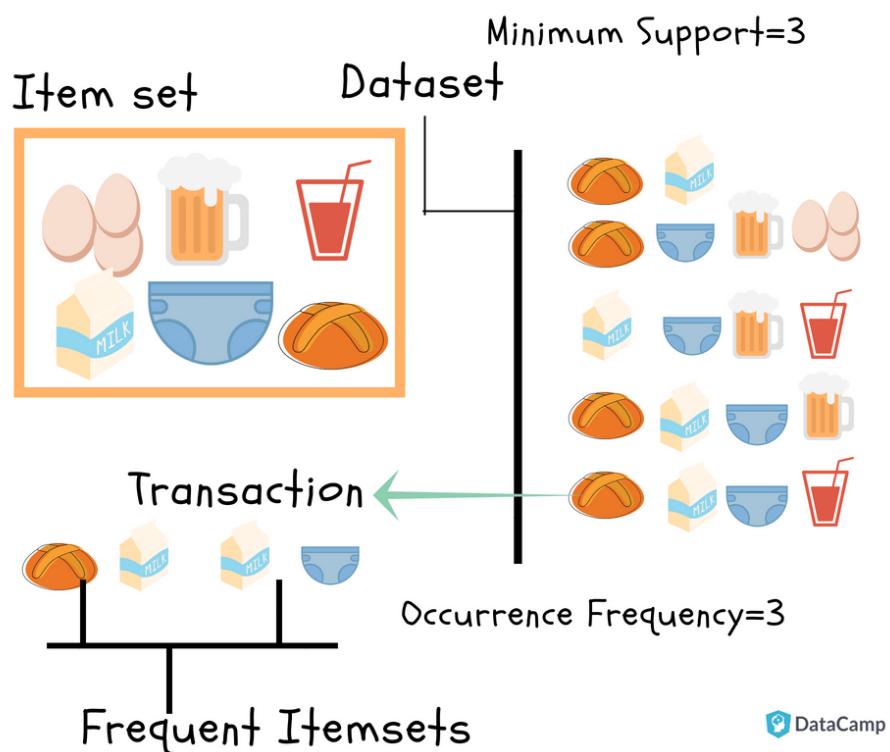


	customer_id	gender	age	annual_income	spending_score	Kmeans	HC	dbscan	anomaly
0	1	Male	19		15	39	2	0	-1
1	2	Male	21		15	81	2	1	0
2	3	Female	20		16	6	2	0	-1
3	4	Female	23		16	77	2	1	0
4	5	Female	31		17	40	2	0	-1



Nama Dosen : Teguh Iman Hermanto, M.Kom  
 Mata Kuliah : Machine Learning 1  
 Pembahasan : Algoritma Association Rules  
 Pokok Pemb :  
 - Mengenal Market Basket analysis  
 - Membangun Model Algoritma Apriori  
 - Membangun Model Algoritma FP-Growth

## 1. Market Basket Analysis



## 2. Load dataset pada file Notebook

```
● ● ●  
1 import pandas as pd  
2 import matplotlib.pyplot as plt  
3 from mlxtend.frequent_patterns import apriori, fpgrowth  
4 from mlxtend.frequent_patterns import association_rules  
5 from mlxtend.preprocessing import TransactionEncoder
```

```
● ● ●  
1 # Load data  
2 file_path = "Groceries_dataset.csv"  
3 data = pd.read_csv(file_path)
```

## 3. Analisa Data

```
● ● ●  
1 # Tampilkan informasi awal dataset  
2 print("Preview Dataset:")  
3 print(data.head())  
4 print("\n\nInfo Dataset:")  
5 print(data.info())  
6 print("\nStatistik Deskriptif Kolom Numerik:")  
7 print(data.describe())
```

```
● ● ●  
1 # 1. Analisis kolom 'Member_number'  
2 # Menghitung total member  
3 total_member = data['Member_number'].nunique()  
4 print(f"\n\nTotal member: {total_member}")
```



```
1 # Top 5 anggota dengan transaksi terbanyak
2 top_members = data['Member_number'].value_counts().head()
3 print("\nTop 5 anggota dengan transaksi terbanyak:")
4 print(top_members)
```



```
1 # 2. Analisis Kolom `Date`
2 data['Date'] = pd.to_datetime(data['Date'], format='%d-%m-%Y')
3
4 # Rentang waktu dataset
5 start_date = data['Date'].min()
6 end_date = data['Date'].max()
7 print(f"\nRentang waktu data: {start_date} hingga {end_date}")
```



```
1 # Hitung jumlah transaksi per bulan
2 monthly_transactions = data['Date'].dt.to_period('M').value_counts().sort_index()
3 print("\nJumlah transaksi per bulan:")
4 print(monthly_transactions)
```



```
1 # Visualisasi jumlah transaksi per bulan
2 plt.figure(figsize=(12, 6))
3 monthly_transactions.plot(kind='bar',
4                           title='Transaksi per Bulan',
5                           color='skyblue')
6 plt.xlabel('Bulan')
7 plt.ylabel('Jumlah Transaksi')
8 plt.xticks(rotation=45)
9 plt.show()
```

```
● ● ●  
1 # 3. Analisis Kolom `itemDescription`  
2 # Barang paling sering dibeli  
3 item_counts = data['itemDescription'].value_counts()  
4 print("\nTop 10 barang paling sering dibeli:")  
5 print(item_counts.head(10))
```

```
● ● ●  
1 # Visualisasi 10 barang teratas  
2 plt.figure(figsize=(10, 6))  
3 item_counts.head(10).plot(kind='bar',  
4                             title='Top 10 Barang Paling Sering Dibeli',  
5                             color='orange')  
6 plt.xlabel('Barang')  
7 plt.ylabel('Frekuensi')  
8 plt.xticks(rotation=45)  
9 plt.show()
```

```
● ● ●  
1 # 4. Analisis Gabungan  
2 # Barang yang paling sering dibeli oleh anggota tertentu  
3 member_id = 1808  
4 member_items = data[data['Member_number'] == member_id]['itemDescription'].value_counts()  
5 print(f"\nBarang paling sering dibeli oleh anggota {member_id}:")  
6 print(member_items)
```

```
● ● ●  
1 # Pola pembelian barang pada bulan tertentu  
2 month = '2015-07'  
3 monthly_items = data[data['Date'].dt.to_period('M') == month]['itemDescription'].value_counts()  
4 print(f"\nBarang paling sering dibeli pada {month}:")  
5 print(monthly_items)
```

```
● ● ●  
1 # Statistik tambahan: jumlah barang unik  
2 print("\nJumlah barang unik:")  
3 print(data['itemDescription'].nunique())
```

#### 4. Data Preprocessing

```
● ● ●
1 # Persiapan data: Mengelompokan item berdasarkan Member_number dan Date
2 transactions = data.groupby(['Member_number', 'Date'])['itemDescription'].apply(list)
```

```
● ● ●
1 te = TransactionEncoder()
2 te_ary = te.fit(transactions).transform(transactions)
3 df_encoded = pd.DataFrame(te_ary, columns=te.columns_)
```

#### 5. Modeling

```
● ● ●
1 # --- Algoritma Apriori ---
2 print("\n===== Apriori =====")
3 frequent_itemsets_apriori = apriori(df_encoded, min_support=0.005, use_colnames=True)
4 print("Frequent Itemsets (Apriori):")
5 print(frequent_itemsets_apriori)
```

```
● ● ●
1 # --- Algoritma FP-Growth ---
2 print("\n===== FP-Growth =====")
3 frequent_itemsets_fpgrowth = fpgrowth(df_encoded, min_support=0.005, use_colnames=True)
4 print("Frequent Itemsets (FP-Growth):")
5 print(frequent_itemsets_fpgrowth)
```

```
● ● ●
1 # Menampilkan hasil aturan asosiasi
2 print("\nAturan Asosiasi (Apriori):")
3 rules_apriori = association_rules(frequent_itemsets_apriori, metric="confidence", min_threshold=0.1)
4 print(rules_apriori[['antecedents', 'consequents', 'support', 'confidence', 'lift']])
```

```
● ● ●
1 print("\nAturan Asosiasi (FP-Growth):")
2 rules_fpgrowth = association_rules(frequent_itemsets_fpgrowth, metric="confidence", min_threshold=0.1)
3 print(rules_fpgrowth[['antecedents', 'consequents', 'support', 'confidence', 'lift']])
```