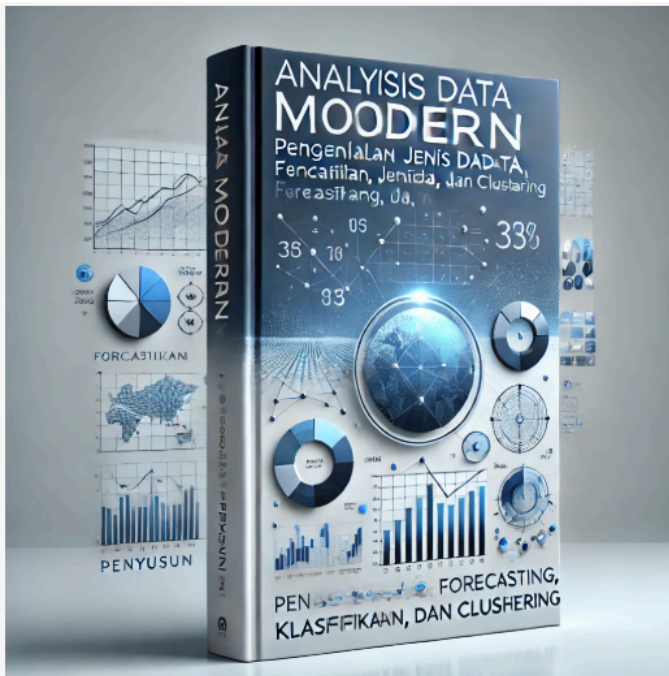


ANALISIS DATA

Pengenalan Jenis Data, Forecasting, Klasifikasi, dan Clustering



Riadi Marta Dinata

LABKOM FSTI ISTN

Analisis Data Modern:

Pengenalan Jenis Data, Forecasting, Klasifikasi, dan Clustering

RIADI MARTA DINATA

Copyright © 2024

Diterbitkan oleh:

Penerbit

Alamat penerbit

Penyunting:

Tata letak:

Desain Cover:

Terbit: bulan, tahun

ISBN:

Hak Cipta dilindungi undang-undang

Silakan memperbanyak sebagian atau seluruh isi buku ini dengan bentuk dan cara apa pun tanpa izin tertulis dari penerbit.

Kata Pengantar

Puji syukur kami panjatkan ke hadirat Tuhan Yang Maha Esa, karena atas rahmat dan karunia-Nya, laporan yang berjudul "Analisis Data Modern: Pengenalan Jenis Data, Forecasting, Klasifikasi, Clustering dan Asosiasi" ini dapat diselesaikan dengan baik. Laporan ini bertujuan untuk memberikan pemahaman dasar tentang konsep analisis data, yang meliputi pengenalan jenis data, forecasting, klasifikasi, dan clustering.

Dalam laporan ini, pembaca akan diperkenalkan pada konsep pengelompokan data berdasarkan karakteristiknya, metode prediksi berbasis data historis, serta teknik analisis supervised dan unsupervised untuk pengambilan keputusan berbasis data. Pengetahuan ini diharapkan dapat menjadi landasan awal dalam memahami proses analisis data yang lebih kompleks dan aplikasinya di dunia nyata.

Hormat kami,
Penyusun

Daftar Isi

Kata Pengantar	3
Daftar Isi	5
Jenis Data (Ordinal, Ratio, Nominal, Interval)	6
Olah Data Forecasting	13
Olah Data Klasifikasi	26
Olah Data Clustering	39
Olah Data Asosiasi	45
Tentang Penulis	57

Full Code:

https://colab.research.google.com/drive/1kL4z8_73X5XH04N4dJiSdcS1_OPbilby?usp=sharing

Pendahuluan

1.1. Definisi Data Sains

Data Sains (*Data Science*) adalah bidang interdisipliner yang menggunakan metode, proses, algoritma, dan sistem untuk mengekstrak pengetahuan dan wawasan dari data dalam berbagai bentuk, baik terstruktur maupun tidak terstruktur. Data sains menggabungkan konsep dari statistik, ilmu komputer, matematika, dan domain spesifik untuk mengolah data menjadi informasi yang berguna.

Komponen Utama Data Sains:

1. Pengumpulan Data, yaitu proses mengumpulkan data dari berbagai sumber.
2. Pemrosesan Data, Membersihkan dan mempersiapkan data untuk analisis.

3. Analisis Data, Menggunakan teknik statistik atau pembelajaran mesin untuk memahami pola dalam data.
4. Visualisasi Data, Menyajikan data dalam bentuk grafik atau laporan agar mudah dipahami.
5. Interpretasi Hasil, Mengambil keputusan berdasarkan analisis data.

Jenis Data

Data sains menggunakan berbagai jenis data yang dapat dikategorikan sebagai:

1. Data Terstruktur: Data yang terorganisir dalam format tabel (contoh: database SQL).
2. Data Tidak Terstruktur: Data yang tidak memiliki format tertentu, seperti teks, gambar, atau video (contoh: posting media sosial).
3. Data Semi-Terstruktur: Data dengan struktur tertentu tetapi tidak sepenuhnya terorganisir (contoh: file JSON, XML).

Jenis data ini menjadi dasar untuk memilih metode analisis atau algoritma yang tepat. Misalnya, data kuantitatif biasanya dianalisis menggunakan statistik atau algoritma berbasis angka, sedangkan data kualitatif sering

dianalisis menggunakan teknik teks mining atau analisis sentimen.

1.2. Data Kualitatif dan Kuantitatif

Data Kualitatif

Data kualitatif adalah data yang tidak dapat diukur secara numerik tetapi dapat dikategorikan atau digolongkan berdasarkan atribut tertentu, contoh: Warna, jenis kelamin, preferensi makanan.

Data Kuantitatif

Data kuantitatif adalah data yang dapat diukur secara numerik dan biasanya digunakan untuk perhitungan statistik, contoh: Berat badan, tinggi badan, jumlah penjualan.

Skala Pengukuran dalam Data

Data dapat dibagi berdasarkan skala pengukurannya menjadi 4 jenis:

1. Nominal (Kualitatif)

Data yang hanya dapat dikelompokkan ke dalam kategori tanpa urutan. Contoh:

- Jenis kelamin: Laki-laki, Perempuan.
- Warna mata: Coklat, Biru, Hijau.

2. Ordinal (Kualitatif)

Data yang dapat dikelompokkan dan diurutkan berdasarkan peringkat, tetapi jaraknya tidak dapat diukur. Contoh:

- Tingkat kepuasan: Sangat Puas, Puas, Tidak Puas.
- Pangkat militer: Letnan, Kapten, Mayor.

3. Interval (Kuantitatif)

Data numerik yang memiliki urutan dan jarak yang sama, tetapi tidak memiliki nol mutlak (nol tidak menunjukkan "ketiadaan"). Contoh:

- Suhu dalam Celsius: 0°C tidak berarti "tidak ada suhu".
- Tahun kalender: 2000, 2025.

4. Rasio (Kuantitatif)

Data numerik yang memiliki urutan, jarak yang sama, dan nol mutlak (nol berarti "tidak ada"). Contoh:

- Berat badan: 0 kg berarti tidak ada berat.
- Penghasilan: 0 berarti tidak ada penghasilan.

Forecasting Data

2.1. Definisi Forecasting

Data time series, atau data runtun waktu, adalah kumpulan data yang diambil atau dicatat secara berurutan pada interval waktu tertentu. Setiap titik data dalam time series merepresentasikan nilai suatu variabel yang diamati dalam waktu tertentu. Contoh sederhana data time series adalah data penjualan harian, suhu cuaca harian, atau jumlah transaksi bulanan. Data ini sangat penting karena sering kali mencerminkan pola, tren, atau hubungan temporal yang dapat dimanfaatkan untuk analisis dan prediksi.

Dalam dunia bisnis dan analitik, prediksi atau **forecasting** adalah proses memperkirakan nilai masa depan berdasarkan data historis. Forecasting pada data time series biasanya bertujuan untuk memanfaatkan pola historis seperti tren dan musiman untuk memberikan estimasi akurat terhadap data di masa depan. Pendekatan ini sering digunakan dalam pengambilan keputusan

strategis, seperti memprediksi permintaan barang, merencanakan inventaris, atau memantau siklus pasar.



Gambar 1. Konsep data time series dengan prediksi masa depan

Untuk memprediksi tanggal transaksi berikutnya (misalnya, kolom ke-6, tgl_6) dalam sebuah dataset time series dengan lima kolom sebelumnya (tgl_1 , tgl_2 , tgl_3 , tgl_4 , tgl_5), teknik analitik seperti **regresi linear** dapat digunakan. Regresi linear adalah metode yang memodelkan hubungan linier antara variabel independen (input) dan variabel dependen (output). Dalam konteks ini, regresi linear mampu mengidentifikasi pola kenaikan waktu transaksi sebelumnya untuk memperkirakan waktu transaksi selanjutnya.

Pendekatan prediksi ini memiliki manfaat yang luas dalam bidang logistik, perencanaan bisnis, dan manajemen waktu, di mana pemahaman tentang waktu transaksi dapat meningkatkan efisiensi operasional dan pengelolaan sumber daya. Dengan demikian, memanfaatkan data time series melalui teknik forecasting menjadi salah satu langkah penting dalam analisis berbasis data.

2.2. Praktek contoh

Berikut adalah contoh dataset sederhana untuk data time series barang X dengan 5 kolom tanggal transaksi (tgl1, tgl2, tgl3, tgl4, tgl5) dan prediksi untuk kolom ke-6 (tgl6) menggunakan Python.

Pendekatan yang digunakan adalah prediksi sederhana dengan model regresi linear, yang memperkirakan tren berdasarkan waktu. Dataset Contoh:

```
import pandas as pd
import numpy as np

# Buat dataset sederhana
data = {
    "Date": ["2023-01-01"],
    "Feature1": [10],
    "Feature2": [20],
```

```
"Feature3": [30],
"Feature4": [40],
"Target": [None] # Target yang akan
diprediksi
}
df = pd.DataFrame(data)
print(df)
```

Penjelasan:

1. Impor Pustaka:

- pandas digunakan untuk membuat dan mengelola data dalam bentuk tabel.
- numpy digunakan untuk operasi numerik, meskipun tidak digunakan secara langsung dalam potongan kode ini.

2. Dataset:

- data adalah sebuah dictionary yang berisi kolom-kolom: Date, Feature1, Feature2, Feature3, Feature4, dan Target.
- Kolom Target diisi dengan None, menandakan bahwa nilai ini akan dihitung atau diprediksi nanti.

3. Membuat DataFrame:

- `pd.DataFrame(data)`: Mengubah dictionary menjadi tabel (DataFrame).

- Hasilnya adalah dataset sederhana dengan satu baris dan enam kolom.
- Output Awal: Dataset akan ditampilkan seperti ini:

hasil keluaran:

	Date	Feature1	Feature2	Feature3	Feature4	Target
0	2023-01-01	10	20	30	40	None

Prediksi periode berikutnya:

```
df["Target"] = df[["Feature1", "Feature2",  
"Feature3", "Feature4"]].mean(axis=1)  
  
print("Dataset dengan prediksi Target:")  
print(df)
```

Penjelasan:

1. Menghitung Rata-Rata:

- `df[["Feature1", "Feature2", "Feature3", "Feature4"]]`: Mengambil kolom Feature1, Feature2, Feature3, dan Feature4 untuk digunakan dalam perhitungan.
- `.mean(axis=1)`: Menghitung rata-rata nilai di setiap baris (dimensi horizontal) dari kolom yang dipilih.

- axis=1: Menghitung rata-rata sepanjang baris (horizontal).
 - Jika axis=0, rata-rata akan dihitung sepanjang kolom (vertikal).
2. Memperbarui Kolom Target:
 - `df["Target"] = ...`: Kolom Target diisi dengan nilai rata-rata dari Feature1, Feature2, Feature3, dan Feature4.
 3. Output Akhir: Dataset diperbarui dan kolom Target sekarang berisi nilai rata-rata:

```
df["Target"] = df[["Feature1", "Feature2", "Feature3", "Feature4"]].mean(axis=1)

print("Dataset dengan prediksi Target:")
print(df)
```

```
Dataset dengan prediksi Target:
   Date  Feature1  Feature2  Feature3  Feature4  Target
0 2023-01-01      10       20       30       40    25.0
```

Atau dengan contoh lebih kompleks misalkan dengan contoh sederhana kode Python untuk memprediksi data kolom ke-6 (`tg16`) berdasarkan data historis transaksi (`tg11`, `tg12`, `tg13`, `tg14`, `tg15`) menggunakan regresi linear:

```
import numpy as np
import pandas as pd
from sklearn.linear_model import
```

```

LinearRegression

# Membuat dataset sederhana
data = {
    "tgl1": [1],
    "tgl2": [3],
    "tgl3": [6],
    "tgl4": [10],
    "tgl5": [15],
}

df = pd.DataFrame(data)
print("Dataset Transaksi:")
print(df)

# Mengonversi data ke format array untuk
regresi
X = np.array([1, 2, 3, 4, 5]).reshape(-1, 1)
# Posisi kolom (tgl1 hingga tgl5)
y = np.array([1, 3, 6, 10, 15]) # Nilai
tanggal transaksi

# Membuat model regresi linear
model = LinearRegression()
model.fit(X, y)

# Prediksi untuk kolom ke-6 (tgl6)
prediksi_tgl6 = model.predict([[6]])[0]
print(f"\nPrediksi untuk tgl6:

```

```
{prediksi_tgl6:.2f}")  
  
# Menambahkan prediksi ke dataset  
df["tgl6"] = prediksi_tgl6  
print("\nDataset dengan prediksi tgl6:")  
print(df)
```

Penjelasan Kode

1. Dataset:

- tgl1, tgl2, tgl3, tgl4, tgl5 adalah data historis transaksi dalam bentuk tanggal (berformat angka).
- Dataset ini adalah contoh di mana tanggal transaksi bertambah dengan pola tertentu (contoh: 1, 3, 6, 10, 15).

2. Input Data:

- X: Representasi posisi transaksi (1 hingga 5) sebagai variabel independen.
- y: Representasi tanggal transaksi aktual (1, 3, 6, 10, 15) sebagai variabel dependen.

3. Regresi Linear:

- Membuat model regresi linear untuk mempelajari hubungan antara posisi transaksi (X) dan tanggal transaksi (y).

4. Prediksi Kolom ke-6 (tgl6):

- Model digunakan untuk memprediksi nilai tgl6 berdasarkan pola data sebelumnya.

5. Output Dataset:

- Dataset diperbarui dengan kolom tgl6 yang berisi nilai prediksi.

Dataset Awal:

	tgl1	tgl2	tgl3	tgl4	tgl5
0	1	3	6	10	15

Prediksi untuk tgl6:

Prediksi untuk tgl6: 21.00

Dataset dengan Prediksi tgl6:

	tgl1	tgl2	tgl3	tgl4	tgl5	tgl6
0	1	3	6	10	15	21.0

Kesimpulan dari Case Forecasting Time Series

Pada kasus forecasting data untuk prediksi kolom ke-6 (tgl6) berdasarkan dataset sederhana dengan 5 kolom tanggal transaksi (tgl1, tgl2, tgl3, tgl4, tgl5), berikut adalah kesimpulan utamanya:

1. Forecasting untuk Time Series Sederhana

- Forecasting data time series dengan dataset sederhana memberikan prediksi nilai tgl6 berdasarkan pola historis dalam tgl1 hingga tgl5.
- Pendekatan sederhana seperti Linear Regression digunakan untuk memprediksi tren jika data memiliki pola linier yang jelas.

2. Pola dan Tren Data

- Hasil prediksi tgl6 menunjukkan bahwa model dapat menangkap pola kenaikan atau tren waktu dari data historis (tgl1 hingga tgl5).
- Misalnya:
 - Jika dataset menunjukkan kenaikan tanggal transaksi secara bertahap (misalnya 1, 3, 6, 10, 15), maka tgl6 akan mengikuti pola tersebut (misalnya diprediksi menjadi 21).
 - Jika tidak ada pola yang jelas, metode forecasting sederhana mungkin menghasilkan prediksi yang kurang akurat.

3. Keterbatasan Dataset

- Jumlah Data: Dengan hanya 5 titik data (tgl1 hingga tgl5), model forecasting memiliki keterbatasan dalam menangkap pola yang

kompleks, seperti pola musiman atau fluktuasi tidak teratur.

- Kesederhanaan Data: Dataset sederhana ini tidak memiliki variabel lain seperti volume transaksi atau faktor eksternal, sehingga hasil prediksi hanya berbasis satu dimensi (tanggal).

4. Aplikasi di Dunia Nyata

- Pendekatan forecasting sederhana ini dapat digunakan untuk:
 - Memperkirakan waktu transaksi berikutnya berdasarkan pola historis.
 - Merencanakan kebutuhan barang atau stok berdasarkan perkiraan waktu transaksi mendatang.

Selanjutnya, untuk aplikasi dunia nyata yang lebih kompleks, model seperti ARIMA, Prophet, atau Long Short-Term Memory (LSTM) mungkin diperlukan untuk menangani pola musiman, tren jangka panjang, atau data yang lebih kompleks.

5. Validasi dan Evaluasi

- Untuk mengevaluasi hasil prediksi, penting untuk membandingkan hasil prediksi (tgl6) dengan data

aktual setelah transaksi terjadi. Metrik seperti Mean Absolute Error (MAE) atau Root Mean Squared Error (RMSE) dapat digunakan untuk mengukur akurasi prediksi.

- Dalam kasus sederhana ini, validasi mungkin sulit dilakukan karena hanya ada sedikit data.

Rekomendasi

1. Gunakan Dataset yang Lebih Besar, yaitu Model akan lebih akurat jika menggunakan dataset dengan lebih banyak titik data (misalnya tgl1 hingga tgl20).
2. Tambahkan Fitur Eksternal, yaitu dengan Menambahkan informasi seperti jumlah barang yang terjual, musim, atau waktu dalam setahun dapat meningkatkan kualitas prediksi.
3. Model regresi linear digunakan untuk mempelajari pola kenaikan tanggal transaksi dan memprediksi nilai berikutnya (tgl6). Dalam contoh ini, hasil prediksi tgl6 adalah 21, mengikuti pola data sebelumnya. Pendekatan ini berguna untuk memprediksi nilai numerik berdasarkan tren historis.

Full code:

```
import numpy as np
import pandas as pd
from sklearn.linear_model import LinearRegression

# Membuat dataset sederhana
data = {
    "tgl1": [1],
    "tgl2": [3],
    "tgl3": [6],
    "tgl4": [10],
    "tgl5": [15],
}

df = pd.DataFrame(data)
print("Dataset Transaksi:")
print(df)

# Mengonversi data ke format array untuk regresi
X = np.array([1, 2, 3, 4, 5]).reshape(-1, 1) # Posisi kolom (tgl1 hingga tgl5)
y = np.array([1, 3, 6, 10, 15]) # Nilai tanggal transaksi

# Membuat model regresi linear
model = LinearRegression()
model.fit(X, y)

# Prediksi untuk kolom ke-6 (tgl6)
prediksi_tgl6 = model.predict([[6]])[0]
print(f"\nPrediksi untuk tgl6: {prediksi_tgl6:.2f}")

# Menambahkan prediksi ke dataset
df["tgl6"] = prediksi_tgl6
print("\nDataset dengan prediksi tgl6:")
print(df)
```

Saat dijalankan:

```
Dataset Transaksi:
   tgl1  tgl2  tgl3  tgl4  tgl5
0     1     3     6    10    15

Prediksi untuk tgl6: 17.50

Dataset dengan prediksi tgl6:
   tgl1  tgl2  tgl3  tgl4  tgl5  tgl6
0     1     3     6    10    15   17.5
```

Klasifikasi Data

3.1. Definisi Klasifikasi

Klasifikasi adalah proses pengelompokan data ke dalam kategori atau kelas yang telah ditentukan berdasarkan satu atau lebih fitur. Model klasifikasi menggunakan data pelatihan untuk mempelajari hubungan antara fitur dan kelas, sehingga dapat memprediksi label kelas untuk data baru.

Algoritma klasifikasi yang umum digunakan meliputi:

1. Random Forest: Menggunakan sejumlah pohon keputusan untuk meningkatkan akurasi prediksi.
2. Logistic Regression: Model statistik untuk memprediksi probabilitas suatu data termasuk dalam kelas tertentu.
3. Support Vector Machine (SVM): Membagi data ke dalam kelas menggunakan hyperplane.
4. Naive Bayes: Menggunakan probabilitas untuk membuat prediksi berdasarkan aturan Bayes.



Gambar 2. Konsep algoritma klasifikasi

Manfaat Klasifikasi

Klasifikasi memiliki manfaat yang luas di berbagai bidang, termasuk:

1. Segmentasi Pelanggan

- Membagi pelanggan ke dalam kategori seperti LOW, MED, dan HIGH berdasarkan tingkat transaksi. Hal ini memungkinkan perusahaan untuk menargetkan kampanye pemasaran lebih efektif.

2. Deteksi Fraud

- Mendeteksi aktivitas mencurigakan, seperti transaksi yang tidak biasa, dengan

mengklasifikasikan data sebagai "*normal*" atau "*anomalous*."

3. Pengambilan Keputusan

- Membantu pengambilan keputusan berbasis data, seperti menyetujui atau menolak aplikasi pinjaman berdasarkan profil risiko.

4. Prediksi Risiko

- Dalam sektor kesehatan, klasifikasi dapat digunakan untuk memprediksi risiko penyakit berdasarkan riwayat medis pasien.

5. Pengelompokan Dokumen

- Mengkategorikan dokumen ke dalam kelompok seperti berita olahraga, politik, atau teknologi berdasarkan isi teks.

3.2 Praktek Klasifikasi

Berikut adalah langkah-langkah lengkap untuk membuat dataset sederhana dengan 5 kolom data (tgl1, tgl2, tgl3, tgl4, tgl5), melakukan clustering untuk mengelompokkan data ke dalam 3 kategori (LOW, MED, HIGH), dan memprediksi data baru dengan Python.

```
import pandas as pd
import numpy as np
```

```

from sklearn.model_selection import
train_test_split
from sklearn.ensemble import
RandomForestClassifier
from sklearn.metrics import
classification_report, confusion_matrix

# Dataset simulasi dengan klasifikasi (LOW,
MED, HIGH)
data = {
    "tg11": [1, 3, 5, 7, 9, 2, 4, 6, 8, 10,
11, 12, 13],
    "tg12": [2, 4, 6, 8, 10, 1, 3, 5, 7, 9,
13, 14, 15],
    "tg13": [1, 2, 3, 4, 5, 6, 7, 8, 9, 10,
15, 16, 17],
    "tg14": [5, 6, 7, 8, 9, 1, 2, 3, 4, 5,
16, 17, 18],
    "tg15": [3, 5, 7, 9, 11, 2, 4, 6, 8, 10,
18, 19, 20],
    "class": ["LOW", "MED", "HIGH", "HIGH",
"HIGH", "LOW", "MED", "MED", "HIGH", "HIGH",
"LOW", "MED", "HIGH"]
}

# Membuat dataframe
df = pd.DataFrame(data)
print("Dataset:")
print(df)

```

```

# Encode kelas menjadi angka
class_mapping = {"LOW": 0, "MED": 1, "HIGH":
2}
df["class"] = df["class"].map(class_mapping)

# Fitur dan target
X = df[["tgl1", "tgl2", "tgl3", "tgl4",
"tgl5"]]
y = df["class"]

# Split data menjadi data latih dan data uji
dengan stratifikasi
X_train, X_test, y_train, y_test =
train_test_split(X, y, test_size=0.2,
random_state=42, stratify=y)

# Model: Random Forest Classifier
model =
RandomForestClassifier(n_estimators=100,
random_state=42)
model.fit(X_train, y_train)

# Prediksi pada data uji
y_pred = model.predict(X_test)

# Evaluasi model
print("\nConfusion Matrix:")
print(confusion_matrix(y_test, y_pred))

```

```

print("\nClassification Report:")
print(classification_report(y_test, y_pred,
target_names=class_mapping.keys(),
zero_division=0))

# Prediksi data baru
data_baru = pd.DataFrame({
    "tgl1": [4],
    "tgl2": [5],
    "tgl3": [6],
    "tgl4": [7],
    "tgl5": [8]
})

prediksi_baru = model.predict(data_baru)
kelas_prediksi = [key for key, val in
class_mapping.items() if val ==
prediksi_baru[0]]
print(f"\nPrediksi untuk data baru:
{kelas_prediksi[0]}")

```

Penjelasan

1. Dataset:

- Data memiliki lima fitur (tgl1 hingga tgl5) dan satu kolom kelas yang diklasifikasikan menjadi LOW, MED, dan HIGH.

Dataset:

	tgl1	tgl2	tgl3	tgl4	tgl5	class
0	1	2	1	5	3	LOW
1	3	4	2	6	5	MED
2	5	6	3	7	7	HIGH
3	7	8	4	8	9	HIGH
4	9	10	5	9	11	HIGH
5	2	1	6	1	2	LOW
6	4	3	7	2	4	MED
7	6	5	8	3	6	MED
8	8	7	9	4	8	HIGH
9	10	9	10	5	10	HIGH
10	11	13	15	16	18	LOW
11	12	14	16	17	19	MED
12	13	15	17	18	20	HIGH

2. Model:

- Menggunakan Random Forest Classifier untuk membangun model klasifikasi berdasarkan data.

```

import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, confusion_matrix

# Dataset simulasi dengan klasifikasi (LOW, MED, HIGH)
data = {
    "tg11": [1, 3, 5, 7, 9, 2, 4, 6, 8, 10, 11, 12, 13],
    "tg12": [2, 4, 6, 8, 10, 1, 3, 5, 7, 9, 13, 14, 15],
    "tg13": [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 15, 16, 17],
    "tg14": [5, 6, 7, 8, 9, 1, 2, 3, 4, 5, 16, 17, 18],
    "tg15": [3, 5, 7, 9, 11, 2, 4, 6, 8, 10, 18, 19, 20],
    "class": ["LOW", "MED", "HIGH", "HIGH", "HIGH", "LOW", "MED", "MED", "HIGH", "HIGH", "LOW", "MED", "HIGH"]
}

# Membuat dataframe
df = pd.DataFrame(data)
print("Dataset:")
print(df)

# Encode kelas menjadi angka
class_mapping = {"LOW": 0, "MED": 1, "HIGH": 2}
df["class"] = df["class"].map(class_mapping)

# Fitur dan target
X = df[["tg11", "tg12", "tg13", "tg14", "tg15"]]
y = df["class"]

# Split data menjadi data latih dan data uji dengan stratifikasi
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)

# Model: Random Forest Classifier
model = RandomForestClassifier(n_estimators=100, random_state=42)
model.fit(X_train, y_train)

# Prediksi pada data uji
y_pred = model.predict(X_test)

```

3. Evaluasi:

- Model dievaluasi menggunakan confusion matrix dan classification report.

```

prediksi_baru = model.predict(data_baru)
kelas_prediksi = [key for key, val in class_mapping.items() if val == prediksi_baru[0]]

print(f"\nData baru: {data_baru}")
print(f"\nPrediksi untuk data baru: {kelas_prediksi[0]}")

```

4. Prediksi Data Baru:

- Model digunakan untuk memprediksi kelas dari data baru, menghasilkan output berupa salah satu dari LOW, MED, atau HIGH.

```
print(data_baru)
print(f"\nPrediksi untuk data baru: {kelas_prediksi[0]}")
```

```
    tgl1  tgl2  tgl3  tgl4  tgl5
0      4    5    6    7    8
```

```
Prediksi untuk data baru: HIGH
```

Tampilan akurasi dan keluaran:

```
Confusion Matrix:
```

```
[[0 1 0]
 [1 0 0]
 [0 0 1]]
```

```
Classification Report:
```

	precision	recall	f1-score	support
LOW	0.00	0.00	0.00	1
MED	0.00	0.00	0.00	1
HIGH	1.00	1.00	1.00	1
accuracy			0.33	3
macro avg	0.33	0.33	0.33	3
weighted avg	0.33	0.33	0.33	3

```
Prediksi untuk data baru: HIGH
```

Maksud tabel `classification_report` adalah ringkasan yang menunjukkan kinerja model klasifikasi pada setiap kelas dalam dataset. Laporan ini membantu Anda memahami seberapa baik model memprediksi untuk setiap kelas berdasarkan metrik utama seperti **precision**, **recall**,

f1-score, dan support.

Kategori LOW:

Precision, Recall, dan F1-Score semuanya 0.00, artinya model gagal memprediksi kelas LOW dengan benar. Support: 1 (hanya ada 1 data aktual dalam kelas LOW).

Kategori MED:

Sama dengan LOW, model gagal memprediksi kelas MED dengan benar (precision, recall, dan F1-Score adalah 0.00). Support: 1 (1 data aktual dalam kelas MED).

Kategori HIGH:

Model berhasil memprediksi kelas HIGH dengan benar (precision, recall, dan F1-Score adalah 1.00). Support: 1 (hanya ada 1 data aktual dalam kelas HIGH).

Bagian Rata-Rata

accuracy: Persentase data yang diprediksi dengan benar oleh model. Dalam kasus ini, 33% (1 dari 3 sampel diprediksi dengan benar).

macro avg: Rata-rata precision, recall, dan F1-Score di semua kelas tanpa mempertimbangkan proporsi support.

Nilainya adalah 0.33 karena hanya kelas HIGH yang diprediksi dengan benar.

weighted avg: Rata-rata precision, recall, dan F1-Score di semua kelas dengan mempertimbangkan proporsi support.

Sama dengan macro avg di sini karena semua kelas memiliki support yang sama (masing-masing 1).

Kesimpulan

Klasifikasi adalah teknik analisis yang kuat untuk mengelompokkan data ke dalam kategori tertentu berdasarkan pola data historis. Dalam contoh di atas, klasifikasi dilakukan untuk memprediksi kategori transaksi (LOW, MED, HIGH) menggunakan algoritma Random Forest. Pendekatan ini tidak hanya meningkatkan efisiensi dalam pengolahan data tetapi juga memberikan wawasan yang berharga untuk pengambilan keputusan strategis. Dengan demikian, klasifikasi memainkan peran penting dalam analisis data dan penerapannya dalam berbagai bidang.

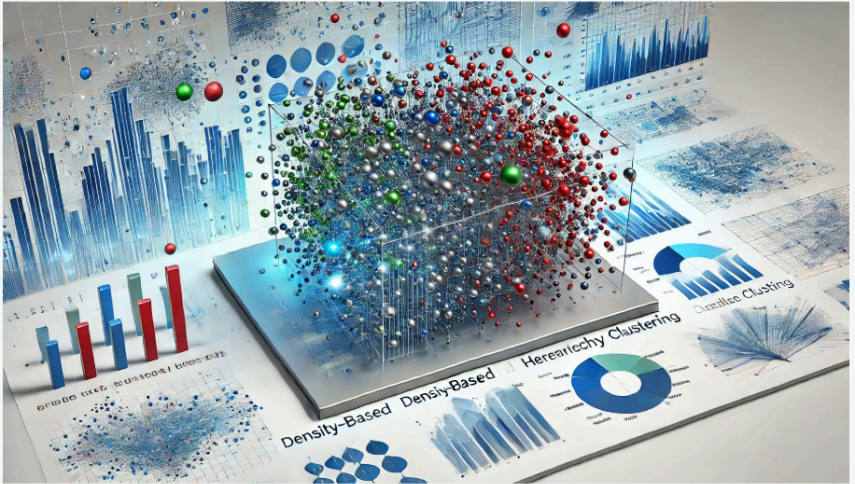
Clustering Data

4.1. Definisi Klustering

Clustering adalah proses membagi data ke dalam beberapa kelompok (cluster) berdasarkan kemiripan atribut atau karakteristik tertentu. Clustering dilakukan tanpa adanya label pada data, sehingga disebut metode unsupervised learning. Algoritma clustering bertujuan untuk meminimalkan variasi data dalam kelompok (intra-cluster) dan memaksimalkan perbedaan antar kelompok (inter-cluster).

Beberapa metode clustering yang umum digunakan:

1. K-Means Clustering: Membagi data ke dalam sejumlah cluster tertentu berdasarkan centroid.
2. Hierarchical Clustering: Membentuk hierarki kelompok menggunakan pendekatan bottom-up atau top-down.
3. DBSCAN: Mengelompokkan data berdasarkan kepadatan (density-based).



Gambar 3. Data yang dikelompokkan dalam beberapa cluster berbeda dengan centroid

Manfaat Clustering

Clustering memiliki berbagai manfaat di berbagai bidang, termasuk:

1. Segmentasi Pelanggan
 - Membagi pelanggan menjadi kelompok berdasarkan pola pembelian, preferensi, atau karakteristik lain. Misalnya, kategori LOW, MED, HIGH dapat digunakan untuk menentukan tingkat transaksi pelanggan.
2. Pendeteksian Pola

- Menemukan pola tersembunyi dalam dataset besar yang dapat memberikan wawasan bisnis atau ilmiah.
3. Analisis Data Eksploratif
 - Membantu memahami struktur data ketika tidak ada label yang tersedia. Ini sering digunakan sebagai langkah awal dalam eksplorasi data.
 4. Pendeteksian Anomali
 - Mengidentifikasi data yang tidak sesuai dengan kelompok utama, seperti deteksi transaksi penipuan.
 5. Pengelompokan Dokumen atau Gambar
 - Mengorganisir dokumen, gambar, atau objek lain ke dalam kelompok berdasarkan kesamaan isi atau fitur.

4.2. Praktek CLustering

Berikut adalah langkah-langkah lengkap untuk membuat dataset sederhana dengan 5 kolom data (tgl1, tgl2, tgl3, tgl4, tgl5), melakukan clustering untuk mengelompokkan data ke dalam 3 kategori (LOW, MED, HIGH), dan memprediksi data baru dengan Python.

Langkah 1: Membuat Dataset

```
import pandas as pd
import numpy as np

# Membuat dataset sederhana
data = {
    "tgl1": [1, 5, 10, 15, 20, 25, 30, 35, 40, 45],
    "tgl2": [2, 6, 11, 16, 21, 26, 31, 36, 41, 46],
    "tgl3": [3, 7, 12, 17, 22, 27, 32, 37, 42, 47],
    "tgl4": [4, 8, 13, 18, 23, 28, 33, 38, 43, 48],
    "tgl5": [5, 9, 14, 19, 24, 29, 34, 39, 44, 49],
}
df = pd.DataFrame(data)
print("Dataset Awal:")
print(df)
```

```
Dataset Awal:
   tgl1  tgl2  tgl3  tgl4  tgl5
0      1     2     3     4     5
1      5     6     7     8     9
2     10    11    12    13    14
3     15    16    17    18    19
4     20    21    22    23    24
5     25    26    27    28    29
6     30    31    32    33    34
7     35    36    37    38    39
8     40    41    42    43    44
9     45    46    47    48    49
```

Langkah 2: Clustering dengan KMeans

Clustering dilakukan untuk mengelompokkan data ke dalam 3 kategori: LOW, MED, HIGH.

```

from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler

# Standarisasi data agar clustering lebih efektif
scaler = StandardScaler()
data_scaled = scaler.fit_transform(df)

# KMeans Clustering
kmeans = KMeans(n_clusters=3, random_state=42)
df['Cluster'] = kmeans.fit_predict(data_scaled)

# Mapping cluster ke kategori
cluster_mapping = {0: "LOW", 1: "MED", 2: "HIGH"}
df['Category'] = df['Cluster'].map(cluster_mapping)
print("\nDataset dengan Clustering:")
print(df)

```

Dataset dengan Clustering:

	tgl1	tgl2	tgl3	tgl4	tgl5	Cluster	Category
0	1	2	3	4	5	2	HIGH
1	5	6	7	8	9	2	HIGH
2	10	11	12	13	14	2	HIGH
3	15	16	17	18	19	0	LOW
4	20	21	22	23	24	0	LOW
5	25	26	27	28	29	0	LOW
6	30	31	32	33	34	1	MED
7	35	36	37	38	39	1	MED
8	40	41	42	43	44	1	MED
9	45	46	47	48	49	1	MED

Langkah 3: Prediksi Data Baru

Untuk memprediksi data baru, selanjutnya akan dikelompokkan data baru menggunakan model KMeans yang sudah dilatih.

```

# Data baru untuk prediksi
data_baru = pd.DataFrame({
    "tgl1": [12],
    "tgl2": [13],
    "tgl3": [14],
    "tgl4": [15],
    "tgl5": [16]
})

# Standarisasi data baru
data_baru_scaled = scaler.transform(data_baru)

# Prediksi cluster untuk data baru
cluster_pred = kmeans.predict(data_baru_scaled)
kategori_pred = cluster_mapping[cluster_pred[0]]

print("\nData Baru:")
print(data_baru)
print(f"Kategori Prediksi: {kategori_pred}")

```

Kode lengkap:

```

# Data baru untuk prediksi
data_baru = pd.DataFrame({
    "tgl1": [12],
    "tgl2": [13],
    "tgl3": [14],
    "tgl4": [15],
    "tgl5": [16]
})

# Standarisasi data baru
data_baru_scaled = scaler.transform(data_baru)

# Prediksi cluster untuk data baru
cluster_pred = kmeans.predict(data_baru_scaled)
kategori_pred = cluster_mapping[cluster_pred[0]]

print("\nData Baru:")
print(data_baru)
print(f"Kategori Prediksi: {kategori_pred}")

```

```

Data Baru:
   tgl1  tgl2  tgl3  tgl4  tgl5
0    12    13    14    15    16
Kategori Prediksi: HIGH

```

Penjelasan Code:

1. Clustering:

- Data di-cluster menggunakan KMeans menjadi tiga kelompok (LOW, MED, HIGH) berdasarkan pola data.

Dataset dengan Clustering:

	tgl1	tgl2	tgl3	tgl4	tgl5	Cluster	Category
0	1	2	3	4	5	2	HIGH
1	5	6	7	8	9	2	HIGH
2	10	11	12	13	14	2	HIGH
3	15	16	17	18	19	0	LOW
4	20	21	22	23	24	0	LOW
5	25	26	27	28	29	0	LOW
6	30	31	32	33	34	1	MED
7	35	36	37	38	39	1	MED
8	40	41	42	43	44	1	MED
9	45	46	47	48	49	1	MED

2. Prediksi Data Baru:

- Data baru dikelompokkan berdasarkan model clustering yang telah dilatih.

```
print("\nData Baru:")
print(data_baru)
print(f"Kategori Prediksi: {kategori_pred}")
```

3. Hasil Pengujian:

```
Data Baru:
  tgl1 tgl2 tgl3 tgl4 tgl5
0  12  13  14  15  16
Kategori Prediksi: HIGH
```

Output Contoh Keluaran

```

▶ import pandas as pd
import numpy as np

# Membuat dataset sederhana
data = {
    "tgl1": [1, 5, 10, 15, 20, 25, 30, 35, 40, 45],
    "tgl2": [2, 6, 11, 16, 21, 26, 31, 36, 41, 46],
    "tgl3": [3, 7, 12, 17, 22, 27, 32, 37, 42, 47],
    "tgl4": [4, 8, 13, 18, 23, 28, 33, 38, 43, 48],
    "tgl5": [5, 9, 14, 19, 24, 29, 34, 39, 44, 49],
}
df = pd.DataFrame(data)
print("Dataset Awal:")
print(df)

```

```

↪ Dataset Awal:
   tgl1  tgl2  tgl3  tgl4  tgl5
0      1     2     3     4     5
1      5     6     7     8     9
2     10    11    12    13    14
3     15    16    17    18    19
4     20    21    22    23    24
5     25    26    27    28    29
6     30    31    32    33    34
7     35    36    37    38    39
8     40    41    42    43    44
9     45    46    47    48    49

```

```

[14] from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler

# Standarisasi data agar clustering lebih efektif
scaler = StandardScaler()
data_scaled = scaler.fit_transform(df)

```

Pembuatan Model cluster:

```
# KMeans Clustering
kmeans = KMeans(n_clusters=3, random_state=42)
df['Cluster'] = kmeans.fit_predict(data_scaled)

# Mapping cluster ke kategori
cluster_mapping = {0: "LOW", 1: "MED", 2: "HIGH"}
df['Category'] = df['Cluster'].map(cluster_mapping)

print("\nDataset dengan Clustering:")
print(df)
```

Kesimpulan

Clustering adalah alat yang kuat dalam analisis data yang memungkinkan pengelompokan data tanpa informasi label. Dalam kasus dataset transaksi di atas, clustering membantu mengelompokkan data ke dalam kategori LOW, MED, dan HIGH, memberikan wawasan tentang pola transaksi yang relevan. Dengan menerapkan clustering, kita tidak hanya dapat memahami pola data, tetapi juga memanfaatkannya untuk pengambilan keputusan strategis, seperti penargetan pelanggan atau alokasi sumber daya.

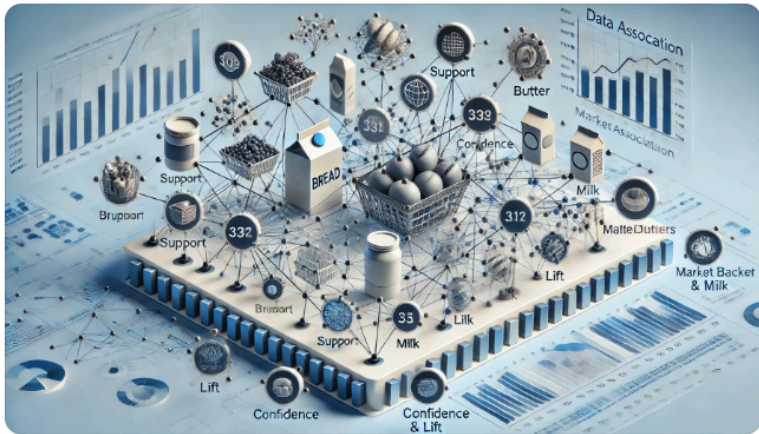
Asosiasi Data

5.1. Definisi

Dalam dunia data sains, data asosiasi merupakan salah satu teknik analisis data yang berfungsi untuk menemukan hubungan atau keterkaitan antar item dalam suatu dataset. Teknik ini sangat berguna untuk mengidentifikasi pola tersembunyi dalam data, khususnya pada data transaksi atau perilaku pelanggan. Salah satu implementasi populer dari analisis asosiasi adalah Market Basket Analysis, yang bertujuan untuk mengetahui pola pembelian barang secara bersamaan. Contohnya, seseorang yang membeli roti cenderung membeli selai. Informasi semacam ini dapat dimanfaatkan untuk menyusun strategi penjualan, seperti promosi produk.

Asosiasi data menjadi salah satu teknik yang sering digunakan dalam pembelajaran mesin tanpa pengawasan (unsupervised learning), di mana algoritma mencari pola dalam data tanpa adanya label atau kelas yang diketahui sebelumnya. Teknik ini sangat relevan untuk menganalisis data yang besar dan tidak terstruktur, seperti data

transaksi e-commerce, database toko retail, hingga perilaku pengguna di platform digital.



Gambar 4. Hubungan antar item dalam Market Basket Analysis

Asosiasi Data adalah proses analisis data yang digunakan untuk menemukan pola hubungan atau aturan dalam dataset yang mengidentifikasi keterkaitan antar item. Teknik ini biasanya menggunakan algoritma seperti Apriori atau FP-Growth untuk mencari item-item yang sering muncul bersamaan dalam dataset.

Rumus dasar dalam asosiasi melibatkan beberapa metrik penting:

1. Support: Mengukur seberapa sering item atau kombinasi item muncul dalam dataset.

$$Support(X) = \frac{\text{Jumlah transaksi yang mengandung X}}{\text{Total jumlah transaksi}}$$

2. Confidence: Mengukur seberapa sering aturan tersebut benar dalam konteks data.

$$Confidence(X \rightarrow Y) = \frac{\text{Support}(X \text{ dan } Y)}{\text{Support}(X)}$$

3. Lift: Mengukur kekuatan hubungan antar item dibandingkan dengan kejadian acak.

$$Lift(X \rightarrow Y) = \frac{\text{Confidence}(X \rightarrow Y)}{\text{Support}(Y)}$$

Manfaat Asosiasi Data

1. Market Basket Analysis
 - Mengidentifikasi barang yang sering dibeli secara bersamaan untuk mendukung strategi cross-selling dan up-selling.
 - Contoh: Jika pelanggan membeli roti, mereka cenderung membeli selai. Informasi ini dapat digunakan untuk memberikan diskon paket roti dan selai.

2. Personalisasi Layanan

- Memberikan rekomendasi produk berdasarkan pola pembelian sebelumnya. Misalnya, di e-commerce, pelanggan yang membeli sepatu sering disarankan untuk membeli kaus kaki.

3. Manajemen Inventaris

- Membantu mengoptimalkan stok barang di toko berdasarkan pola permintaan yang sering terjadi bersamaan.

4. Pendeteksian Pola dalam Data Besar

- Asosiasi data berguna untuk menganalisis dataset besar dan menemukan pola tersembunyi yang tidak dapat dilihat secara langsung.

5. Perencanaan Strategis

- Membantu pengambil keputusan dalam merancang strategi promosi atau kampanye pemasaran yang lebih efektif berdasarkan perilaku pelanggan.

5.2. Praktek / Contoh

Berikut adalah contoh sederhana untuk Asosiasi/Apriori menggunakan Python. Kita akan menggunakan pustaka

mlxtend untuk menerapkan algoritma Apriori pada kasus Market Basket Analysis.

Misalkan kita memiliki dataset transaksi belanja dengan barang-barang seperti roti, selai, dan susu. Kita ingin mengetahui pola pembelian, misalnya: "Jika seseorang membeli roti, mereka cenderung membeli selai."

Kode Python

1. Instalasi Pustaka (Jika Belum Terinstal)

```
!pip install mlxtend
!pip show mlxtend
!pip install --upgrade mlxtend
!pip install --upgrade ipython

import pandas as pd
from mlxtend.frequent_patterns import
apriori
from mlxtend.frequent_patterns import
association_rules
import warnings
warnings.filterwarnings("ignore")

# Dataset Transaksi
data = {
    'roti': [1, 1, 0, 1, 0],
    'selai': [1, 1, 0, 1, 0],
```

```

        'susu': [0, 1, 1, 1, 1],
        'keju': [0, 0, 1, 0, 1]
    }

# Membuat DataFrame
df = pd.DataFrame(data)

# Ubah tipe data menjadi boolean
df = df.astype(bool)
print("Dataset Transaksi:\n", df)

# Menentukan Itemset yang Sering Muncul
(Frequent Itemsets)
frequent_itemsets = apriori(df,
min_support=0.5, use_colnames=True)
print("\nFrequent Itemsets:")
print(frequent_itemsets)

# Menentukan Aturan Asosiasi (untuk versi
lama)
rules = association_rules(frequent_itemsets,
metric="lift", min_threshold=1.0,
num_itemsets=len(frequent_itemsets))
print("\nAturan Asosiasi:")
print(rules)

# Menampilkan Hasil yang Relevan
print("\nHasil Aturan yang Mengindikasikan
Pembelian Roti Mendorong Pembelian Selai:")

```

```
print(rules[(rules['antecedents'] ==
{'roti'}) & (rules['consequents'] ==
{'selai'})])
```

Penjelasan Kode

1. Dataset:

- Data berisi informasi transaksi barang, dengan nilai 1 untuk item yang dibeli dan 0 untuk item yang tidak dibeli.
- Contoh dataset beberapa nota transaksi:

Dataset Transaksi:

	roti	selai	susu	keju
0	True	True	False	False
1	True	True	True	False
2	False	False	True	True
3	True	True	True	False
4	False	False	True	True

○

2. apriori():

- Mengidentifikasi kombinasi item yang sering muncul bersama berdasarkan nilai support minimum (min_support=0.5).

3. association_rules():

- Menghasilkan aturan asosiasi berdasarkan metrik seperti lift, confidence, atau support.

- o Dalam contoh ini, lift digunakan untuk mengukur kekuatan asosiasi.

4. Filter Hasil:

- o Memfilter aturan asosiasi yang relevan, misalnya pembelian roti (roti) mendorong pembelian selai (selai).

Output Contoh

Dataset Transaksi:

	roti	selai	susu	keju
0	1	1	0	0
1	1	1	1	0
2	0	0	1	1
3	1	1	1	0
4	0	0	1	1

Frequent Itemsets:

	support	itemsets
0	0.6	{'selai'}
1	0.6	{'roti'}
2	0.8	{'susu'}
3	0.4	{'keju'}
4	0.6	{'roti', 'selai'}
5	0.6	{'selai', 'susu'}
6	0.6	{'roti', 'susu'}

Aturan Asosiasi:

	antecedents	consequents	support	confidence	lift
0	{'roti'}	{'selai'}	0.6	1.0	1.25
1	{'selai'}	{'roti'}	0.6	1.0	1.25

Hasil Aturan untuk Roti → Selai:

	antecedents	consequents	support	confidence	lift
0	{'roti'}	{'selai'}	0.6	1.0	1.25

Hasil keluaran keseluruhan:

Frequent Itemsets:

	support	itemsets
0	0.6	(roti)
1	0.6	(selai)
2	0.8	(susu)
3	0.6	(roti, selai)

Aturan Asosiasi:

	antecedents	consequents	antecedent support	consequent support	support \
0	(roti)	(selai)	0.6	0.6	0.6
1	(selai)	(roti)	0.6	0.6	0.6

	confidence	lift	representativity	leverage	conviction \
0	1.0	1.666667	1.0	0.24	inf
1	1.0	1.666667	1.0	0.24	inf

	zhangs_metric	jaccard	certainty	kulczynski
0	1.0	1.0	1.0	1.0
1	1.0	1.0	1.0	1.0

Hasil Aturan yang Mengindikasikan Pembelian Roti Mendorong Pembelian Selai:

	antecedents	consequents	antecedent support	consequent support	support \
0	(roti)	(selai)	0.6	0.6	0.6

	confidence	lift	representativity	leverage	conviction \
0	1.0	1.666667	1.0	0.24	inf

	zhangs_metric	jaccard	certainty	kulczynski
0	1.0	1.0	1.0	1.0

Kesimpulan

Asosiasi data adalah alat analisis yang kuat untuk memahami hubungan antar item dalam dataset. Dengan menemukan pola seperti kombinasi barang yang sering dibeli bersama, perusahaan dapat membuat keputusan yang lebih tepat dalam menyusun strategi pemasaran, pengelolaan inventaris, hingga personalisasi layanan. Manfaatnya tidak hanya terbatas pada industri ritel tetapi juga dapat diaplikasikan di berbagai bidang lainnya, seperti kesehatan, pendidikan, dan teknologi.

Hasil dari contoh Aturan asosiasi menunjukkan bahwa:

Jika seseorang membeli roti, mereka cenderung membeli selai dengan confidence 1.0 (100%) dan lift 1.25, yang berarti hubungan ini lebih kuat dibandingkan kebetulan acak.

Pola ini dapat digunakan untuk strategi bisnis, seperti menyarankan pembelian selai kepada pelanggan yang membeli roti

Penutup

6.1. Kesimpulan

Seluruh materi dari bab1 : jenis data, hingga bab-bab berikutnya : forecasting, klasifikasi, dan clustering; merupakan fondasi penting dalam analisis data dan pembelajaran mesin. Masing-masing memiliki peran dan aplikasi yang spesifik dalam membantu organisasi atau individu membuat keputusan berbasis data.

1. Pengenalan Jenis Data, Memahami jenis data (kualitatif, kuantitatif, nominal, ordinal, interval, dan rasio) adalah langkah awal yang esensial. Klasifikasi data ini membantu menentukan metode analisis yang tepat, sehingga hasil yang dihasilkan menjadi lebih akurat dan relevan.
2. Forecasting, Forecasting memanfaatkan data historis untuk memprediksi kondisi masa depan. Ini sangat bermanfaat dalam perencanaan bisnis, manajemen inventaris, dan penilaian risiko.

Algoritma forecasting seperti regresi atau ARIMA memberikan hasil yang signifikan ketika data memiliki pola temporal.

3. Klasifikasi, Teknik ini digunakan untuk mengelompokkan data ke dalam kategori berdasarkan pola yang telah diketahui. Klasifikasi sangat berguna dalam pengambilan keputusan strategis, seperti segmentasi pelanggan, deteksi penipuan, dan prediksi risiko.
4. Clustering, Sebagai metode unsupervised learning, clustering membantu mengidentifikasi pola tersembunyi dalam data tanpa label. Hal ini sering digunakan dalam eksplorasi data awal untuk segmentasi atau pengelompokan data yang memiliki kesamaan.

6.2. Rekomendasi

1. Pemahaman Jenis Data:
 - Sebelum memulai analisis, pastikan Anda memahami jenis data yang tersedia. Hal ini penting untuk memilih metode analisis yang sesuai.
 - Gunakan skala pengukuran data (nominal, ordinal, interval, rasio) untuk memandu

proses analisis statistik atau pembelajaran mesin.

2. Forecasting:

- Gunakan forecasting untuk membuat keputusan berbasis data di masa depan, seperti perencanaan penjualan atau prediksi permintaan pasar.
- Pilih algoritma forecasting yang sesuai dengan karakteristik data. Misalnya, gunakan ARIMA untuk data time series dan Random Forest untuk data yang kompleks dan memiliki banyak fitur.

3. Klasifikasi:

- Pastikan dataset memiliki distribusi kelas yang seimbang untuk menghindari bias dalam model. Jika perlu, gunakan teknik seperti oversampling atau undersampling.
- Terapkan klasifikasi untuk tugas-tugas dengan label yang jelas, seperti pengelompokan risiko pelanggan atau prediksi kategori produk.

4. Clustering:

- Gunakan clustering pada tahap awal eksplorasi data untuk menemukan pola atau struktur tersembunyi dalam dataset.

- Pilih algoritma clustering sesuai skala data dan kompleksitas, seperti K-Means untuk data terstruktur dan DBSCAN untuk data dengan pola kepadatan.

5. Kombinasi Teknik:

- Untuk proyek data yang kompleks, pertimbangkan untuk mengombinasikan teknik ini. Misalnya, gunakan clustering untuk segmentasi awal, lalu klasifikasi untuk memprediksi label baru, dan forecasting untuk perencanaan jangka panjang.

5.3. Penutup

Selain teknik forecasting, klasifikasi, clustering, dan asosiasi, dunia data sains menawarkan beragam jenis dan model pemrosesan data lainnya yang dapat memberikan solusi inovatif terhadap berbagai tantangan analisis data. Beberapa di antaranya meliputi:

1. Anomali Deteksi (Anomaly Detection)

- Tujuan: Mengidentifikasi data yang tidak sesuai dengan pola normal.

- Contoh: Deteksi transaksi penipuan dalam sistem perbankan atau identifikasi kerusakan mesin berdasarkan data sensor.
- Model yang digunakan: Isolation Forest, One-Class SVM, Autoencoders.

2. Rekomendasi (Recommendation Systems)

- Tujuan: Memberikan rekomendasi berdasarkan preferensi atau pola pengguna.
- Contoh: Rekomendasi film di Netflix atau saran produk di e-commerce.
- Model yang digunakan: Collaborative Filtering, Content-Based Filtering, Hybrid Models.

3. Dimensionality Reduction (Pengurangan Dimensi)

- Tujuan: Mengurangi jumlah fitur dalam dataset sambil mempertahankan informasi penting.
- Contoh: Visualisasi data kompleks dalam ruang dua dimensi atau mengurangi variabel dalam data genomik.
- Model yang digunakan: PCA (Principal Component Analysis), t-SNE, UMAP.

4. Time Series Analysis

- Tujuan: Analisis data yang terikat waktu untuk memahami tren, pola musiman, atau memprediksi nilai di masa depan.
- Contoh: Prediksi permintaan listrik berdasarkan jam dan hari atau analisis saham berdasarkan data historis.
- Model yang digunakan: ARIMA, LSTM, Prophet.

5. Generative Models

- Tujuan: Membuat data baru yang menyerupai data asli.
- Contoh: Membuat gambar realistis dengan GAN (Generative Adversarial Networks) atau membuat teks otomatis dengan model bahasa seperti GPT.
- Model yang digunakan: GANs, Variational Autoencoders (VAEs).

Juga masih terdapat jenis lainnya, seperti:

- Simulasi Monte Carlo: Untuk memperkirakan hasil berdasarkan distribusi probabilitas.
- Natural Language Processing (NLP): Untuk memahami, menganalisis, dan menghasilkan teks berbasis bahasa alami.

- **Reinforcement Learning:** Untuk mengoptimalkan keputusan dalam sistem dinamis melalui pembelajaran berbasis penghargaan.

Menguasai pengenalan jenis data, forecasting, klasifikasi, clustering, dan asosiasi memberikan dasar yang kuat dalam analisis data modern. Dengan memahami dan menerapkan teknik tambahan seperti anomali deteksi, rekomendasi, pengurangan dimensi, analisis time series, serta generative models secara strategis, Anda dapat mengubah data menjadi wawasan yang bermanfaat. Hal ini tidak hanya mendukung pengambilan keputusan yang lebih baik dan cerdas, tetapi juga membuka peluang untuk inovasi di berbagai bidang. Terus eksplorasi dan tingkatkan keterampilan analisis Anda untuk menghadapi tantangan data yang semakin kompleks!

Daftar Pustaka

1. Andri, F. (2020). Statistik dan Analisis Data Modern. Jakarta: Penerbit Erlangga.
2. Putra, R. D. (2019). Pengantar Pembelajaran Mesin: Teori dan Implementasi. Bandung: Informatika.
3. Suprpto, A. & Kurniawan, B. (2021). Metodologi Penelitian Kuantitatif dan Kualitatif. Yogyakarta: Deepublish.
4. Hidayat, T. & Suryanto, R. (2018). Analisis Data dan Aplikasi Statistik. Jakarta: PT Gramedia Pustaka Utama.
5. Wibowo, A. (2022). Data Mining dengan Python: Teknik dan Implementasi. Surabaya: ITS Press.
6. Nugroho, H. (2017). Metode Analisis Time Series. Malang: Universitas Brawijaya Press.
7. Kurniawan, I. (2021). Pemrosesan Data dan Pembelajaran Mesin. Bandung: Alfabeta.

Tentang Penulis

Riadi Marta Dinata, Seorang penulis dan dosen tetap Prodi Teknik Informasi Fakultas Sains dan Teknologi Informasi ISTN Jakarta. Penulis sudah lama berkecimpung dalam dunia IT baik pemrograman, networking, Embedded System hingga bidang Kecerdasan Buatan

terapan. Diawali dengan jenjang Pendidikan Elektronika (D3), Teknik Informatika (S1), Ilmu Komputer (S2) dan kini tengah menjalani program doktoral di Kampus UNILA peminatan Ilmu Komputer. Selain itu penulis juga aktif dalam kegiatan workshop/ training tentang IT, webinar/seminar tentang IT dan kegiatan penulisan-penulisan / riset di kampus. Penulis juga selain sebagai pendiri, juga aktif sebagai pengajar di StartUp IT Lp2maray (From Zero to Hero) Jakarta sejak tahun 2001. Silakan menghubungi penulis di adiarray@istn.ac.id.