

Virtual Internship Experience

Clustering Model

Disclaimer

“Dokumen ini memiliki hak cipta. Barang siapa yang menyebarluaskan atau menduplikasi tanpa izin dari instansi terkait dapat diproses sesuai dengan ketentuan hukum yang berlaku.”

Outline

- 1 ***Clustering Model***
- 2 ***K-Means Clustering***
- 3 **Aplikasi Clustering**
- 4 **Studi Kasus**

1

Clustering Model

❏ Pengertian

Clustering adalah algoritma machine learning unsupervised yang mengelompokkan titik data bersama berdasarkan kesamaan mereka.

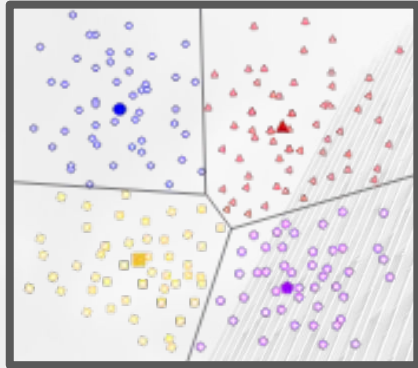
Tujuan dari clustering adalah untuk menemukan kelompok titik data yang mirip satu sama lain dan berbeda dari titik data dalam kelompok lain.



1

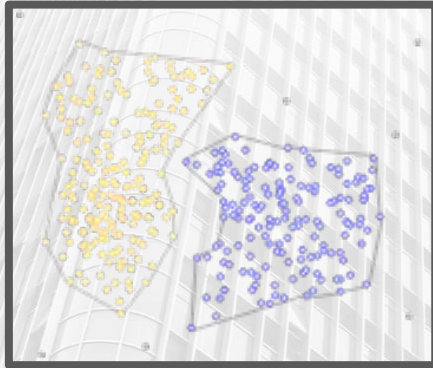
Clustering Model

❏ Jenis Clustering



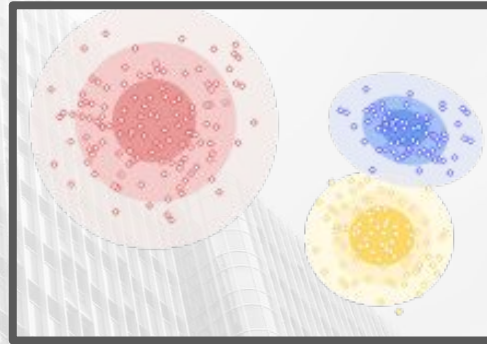
Gambar 1. Centroid-Based

Membagi data ke dalam sejumlah kluster tertentu, di mana setiap kluster memiliki sebuah sentroid yang mewakili rata-rata dari semua titik data dalam kluster tersebut.



Gambar 2. Density-Based

Algoritma ini menemukan kluster yang padat dan terpisah dengan baik satu sama lain.



Gambar 3. Distribution-Based

Algoritma ini mengasumsikan bahwa titik data dalam setiap kelompok mengikuti distribusi tertentu, seperti distribusi Gaussian.



Gambar 4. Hierarchical

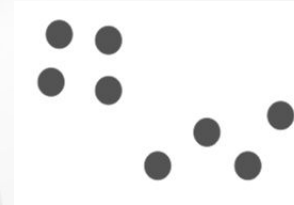
Algoritma ini membangun hierarki dari kelompok, dengan setiap kelompok menjadi anak dari kelompok yang lebih besar. Kelompok-kelompok tersebut digabungkan atau dibagi berdasarkan kesamaan mereka.

2

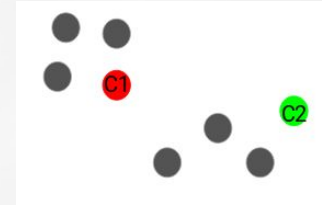
K-Means Clustering

Metode K-Means Clustering:

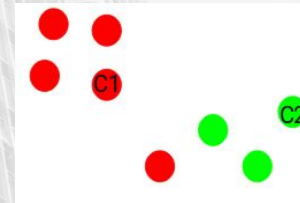
1. Pengguna menentukan jumlah kluster k .
2. Secara acak memilih k titik data sebagai pusat awal dari kluster.
3. mengelompokkan setiap titik data ke kluster dengan pusat terdekat.
4. Menghitung kembali nilai rata-rata centroid berdasarkan semua titik data yang telah dikelompokkan.
5. Mengulangi langkah 3 dan 4 sampai pusat tidak berubah lagi.



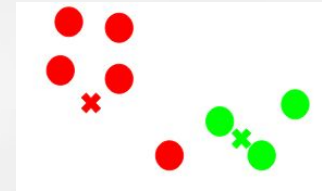
Gambar 1. Sebaran data



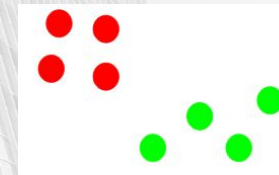
Gambar 2. pemilihan titik k secara random



Gambar 3. mengelompokkan data ke kluster terdekat



Gambar 4. Menghitung ulang pusat setiap kluster

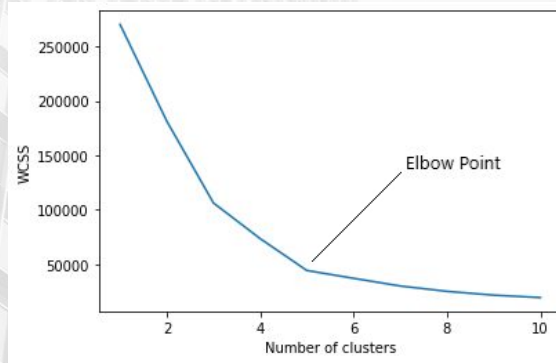


Gambar 5. mengelompokkan data ke pusat baru

2

K-Means Clustering

- ❏ Menentukan banyaknya K (kluster / kelompok)



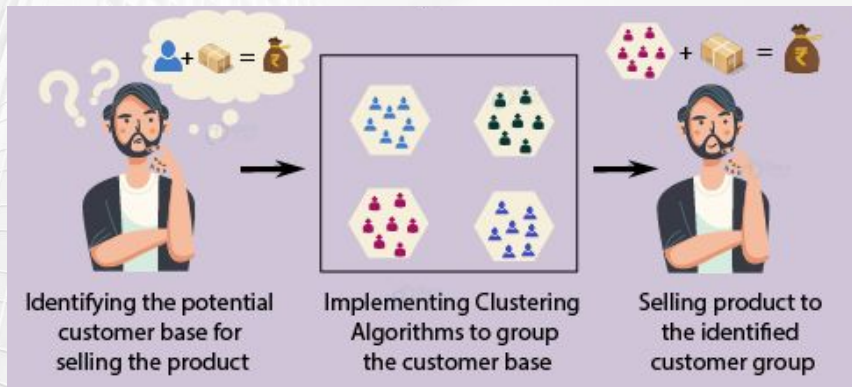
Gambar 6. Contoh Elbow Method

1. Mendefinisikan rentang nilai K untuk menjalankan *K-Means Clustering*
2. Mengevaluasi *Sum of Squared Errors* (SSE) untuk model menggunakan setiap jumlah klaster yang telah ditentukan.

3

Aplikasi Clustering

❏ Segmentasi Pelanggan



Gambar 1. Keuntungan Segmentasi pelanggan

❏ Implementasi K-Means Clustering



Gambar 2. Logo scikit-learn

4 Studi Kasus

Melihat penjualan produk yang menurun, tim marketing ingin membuat suatu program promo untuk meningkatkan penjualan produk. Sebelum memulai program, tim ingin membuat promo khusus kepada setiap pelanggan Kalbe Nutritionals. Sebagai seorang Data Scientist, anda memahami permasalahan tersebut dapat diselesaikan dengan clustering dan akan membuat model untuk permasalahan ini.

Petunjuk Pengerjaan:

1. Anda dapat mengerjakan menggunakan *jupyter notebook* atau [google colab](#)
2. Lakukan data cleaning dan eksplorasi data terdahulu.
3. Pastikan format data sudah sesuai, bila tidak anda terapkan *data transformation*.
4. selamat mencoba

[Link Dataset](#)



Solusi

- Memanggil Library yang diperlukan dan Load dataset
- Menampilkan informasi dataset
- Mengecek *missing values*
- Membuat histogram pada data pembelian dan data umur
- Membuat bar chart pada data jenis kelamin

Load Libraries dan Dataset

```
# library untuk pengolahan data
import pandas as pd
import numpy as np

# library untuk visualisasi data
import matplotlib.pyplot as plt
import seaborn as sns

# library untuk machine learning model
from sklearn.cluster import KMeans

df = pd.read_csv('Dataset GCV 6 - Clustering Model.csv')
df.head()

id_pelanggan  umur  jenis_kelamin  id_produk  nama_produk  kategori_produk  harga_pembelian
0             1    31             wanita      2    prenatal esensis  woman             54818
1             2    28             wanita      2    prenatal esensis  woman             55539
2             3    21             pria        3    hic1000 vitamin lemon  beverages          55226
3             4    38             pria        6    nutritive benecol    special needs      48295
4             5    42             pria        6    nutritive benecol    special needs      45963

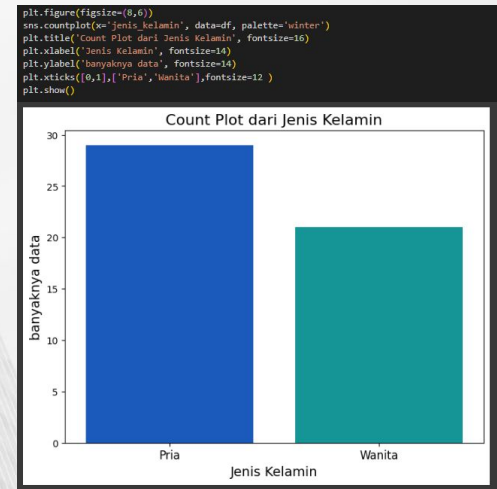
[3] df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50 entries, 0 to 49
Data columns (total 7 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   id_pelanggan  50 non-null     int64
 1   umur         50 non-null     int64
 2   jenis_kelamin  50 non-null     object
 3   id_produk     50 non-null     int64
 4   nama_produk   50 non-null     object
 5   kategori_produk  50 non-null     object
 6   harga_pembelian  50 non-null     int64
dtypes: int64(4), object(3)
memory usage: 2.9+ KB
```

Gambar 1. Import libraries dan dataset



Gambar 2. mengecek missing values dan visualisasi distribusi data



Gambar 3. membuat visualisasi jenis kelamin

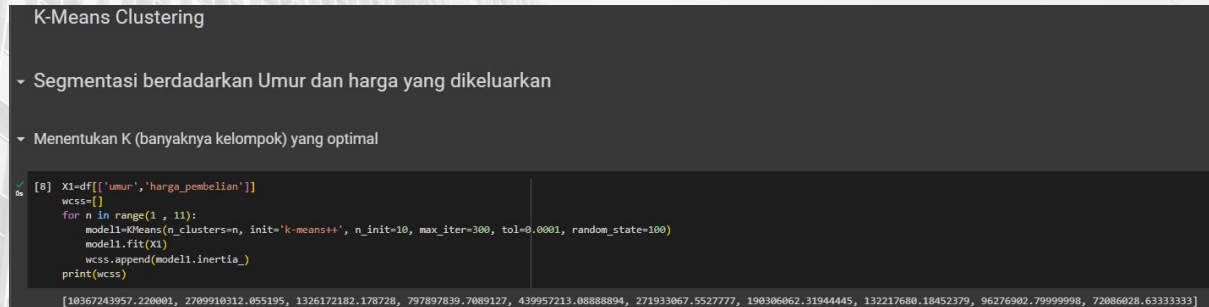


Solusi

- Membuat scatterplot dengan tiga variabel:
 - Umur
 - Harga pembelian
 - Jenis Kelamin
- Membuat model clustering dengan $K = 1$ sampai $K = 10$ untuk menentukan cluster yang optimal



Gambar 4. Membuat visualisasi sebaran data antara 3 variabel

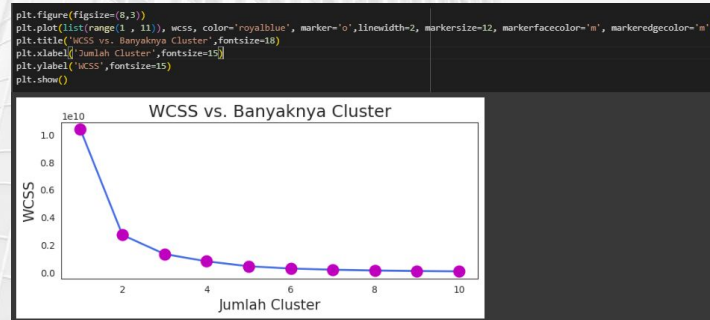


Gambar 5. Membangun model dengan $K=1$ sampai $K=10$



Solusi

- Membuat visualisasi antara jumlah cluster (K) dengan nilai WCSS (within cluster sum of squares) / elbow method
- Mendapati bahwa K optimal saat K= 3
- Membangun kembali model clustering dengan K = 3
- Melatih model dengan data harga pembelian dan umur pelanggan



Gambar 6. elbow method

Membuat model clustering dengan K yang optimal

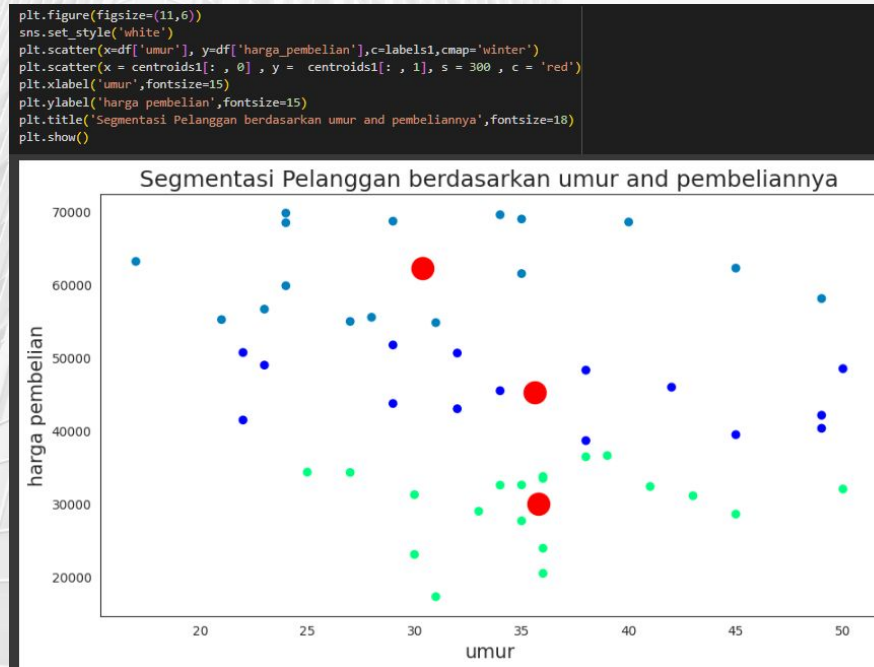
```
model1=KMeans(n_clusters=3, init='k-means++', n_init=10, max_iter=300, tol=0.0001, random_state=100)
model1.fit(X1)
labels1=model1.labels_
centroids1=model1.cluster_centers_
```

Gambar 7. Membuat model K-Means dengan K=3



Solusi

- Model telah membuat cluster pada sebaran data
- Membuat visualisasi persebaran data untuk menunjukkan segmentasi pelanggan dengan fitur umur dan harga pembelian produk



Gambar 8. Hasil akhir segmentasi pelanggan berdasar umur dan harga beli



Thank You!



KALBE
Nutritional