

Machine Learning Project

Data Scientist Project Based Internship Program

Presented by Karina Khairunnisa Putri



Karina Khairunnisa Putri

About Me

I am an enthusiastic, highly motivated and analytical individual. As a 5th semester student majoring in Informatics who's passionate about Data Science. Also, I took some certificates are Introduction to Data Science with R and maximize data science skills. With my educational background in computer science, artificial intelligence and statistics, I have sufficient knowledge in data processing, data analysis, and data modeling. Skilled in operating python, SQL, C++, and R program languages. Looking for an opportunity to develop my skills for an internship in Data Science.



My Experience



Participant of Scientific / Interest Competitions in data science with R language

Experience 2

Project IMPACT (Coding competition event for knowledge in Informatics major) in 2022

2. Case Study



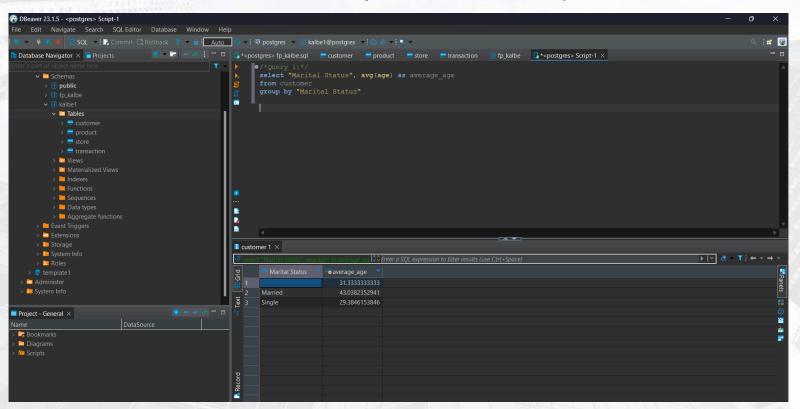
Peserta dapat melakukan exploratory data analysis di dbeaver

Peserta dapat melakukan exploratory data analysis di dbeaver

- 1. query 1 : Berapa rata-rata umur customer jika dilihat dari marital statusnya?
 - 2. query 2 : Berapa rata-rata umur customer jika dilihat dari gender nya ?
 - 3. query 3 : Tentukan nama store dengan total quantity terbanyak!
 - 4. query 4 : Tentukan nama produk terlaris dengan total amount terbanyak!

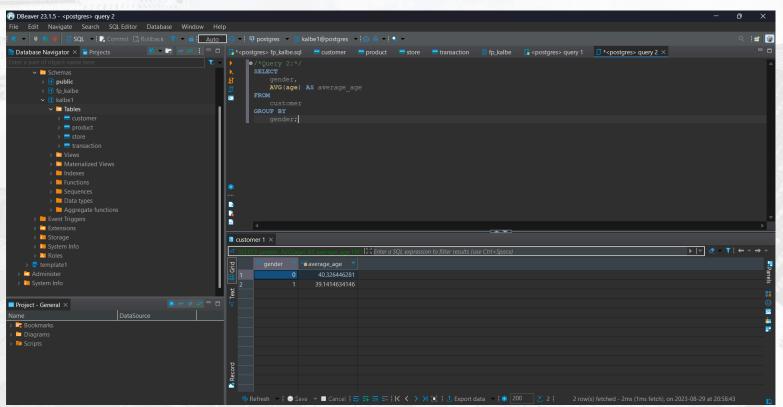


Rata-rata umur customer jika dilihat dari marital statusnya



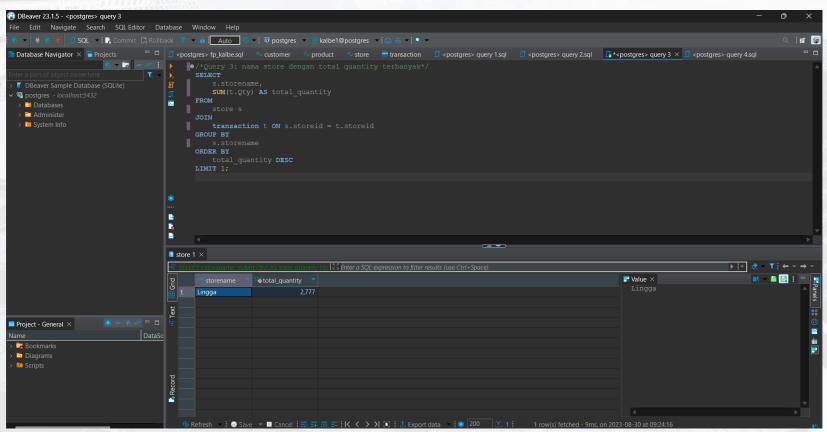


Rata-rata umur customer jika dilihat dari gender



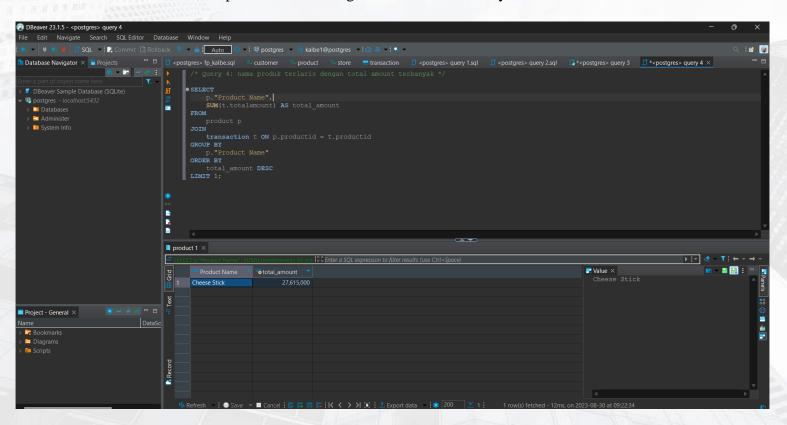


Nama store dengan total quantity terbanyak adalah store Lingga





Nama produk terlaris dengan total amount terbanyak adalah cheese stick



1. Case Study



Machine Learning Regression (Time Series)

Tujuan dari pembuatan model machine learning ini adalah untuk dapat memprediksi total quantity harian dari product yang terjual.

- Data cleansing terlebih dahulu, merubah tipe data supaya sesuai
 - Data merge untuk menggabungkan semua data
- Membuat data baru untuk regression, yaitu groupby by date lalu yang di aggregasi adalah qty di sum
 - Akan ada sekitar 365 rows
 - Menggunakan metode time series ARIMA



```
#Library import statements
import pandas as pd
import numpy as np

import seaborn as sns
import matplotlib.pyplot as plt

from sklearn.metrics import mean_absolute_error, mean_squared_error
from sklearn import preprocessing
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_score

from statsmodels.tsa.seasonal import seasonal_decompose
from statsmodels.tsa.holtwinters import SimpleExpSmoothing, Holt
from statsmodels.tsa.arima.model import ARIMA
from pandas.plotting import autocorrelation_plot
import warnings
```

```
[56] customer_data = pd.read_csv('Case Study - Customer.csv')

product_data = pd.read_csv('Case Study - Product.csv')

store_data = pd.read_csv('Case Study - Store.csv')

transaction_data = pd.read_csv('Case Study - Transaction.csv')
```



#menampilkan bentuk (shape) dari empat dataframe yang berbeda customer_data.shape, product_data.shape, store_data.shape, transaction_data.shape

((447, 5), (10, 3), (14, 6), (5020, 8))

os [7]	cus	tomer_data	ı.h	ead())			
		CustomerI	D	Age	Gender	Marital Status	Income	
	0		1	55	1	Married	5	ıl.
	1	:	2	60	1	Married	6	
	2	;	3	32	1	Married	9	
	3		4	31	1	Married	4	
	4	:	5	58	1	Married	3	





os [9] t	ran	saction_data.h	nead()							
		TransactionID	CustomerID	Date	ProductID	Price	Qty	TotalAmount	StoreID	
	0	TR11369	328	1/1/2022	P3	7500	4	30000	12	11.
	1	TR16356	165	1/1/2022	P9	10000	7	70000		
	2	TR1984	183	1/1/2022	P1	8800	4	35200	4	
	3	TR35256	160	1/1/2022	P1	8800	7	61600	4	
	4	TR41231	386	1/1/2022	P9	10000		10000	4	

vos [10]	sto	re_data.h	nead()					
		StoreID	StoreName	GroupStore	Туре	Latitude	Longitude	
	0		Prima Tendean	Prima	Modern Trade	-6	NaN	ili
	1	2	Prima Kelapa Dua	Prima	Modern Trade	-6	NaN	
	2	3	Prima Kota	Prima	Modern Trade	-7	NaN	
	3	4	Gita Ginara	Gita	General Trade	-6	NaN	
	4	5	Bonafid	Gita	General Trade	-7	NaN	





[15] trans	action_data [tra	nsaction_dat	t a[' Transact	tionID'] ==	'TR713	13']			
	TransactionID	CustomerID	Date	ProductID	Price	Qty	TotalAmount	StoreID	
1982	TR71313	117	2022-05-21	P1	8800	10	88000	8	
3336	TR71313	401	2022-08-30	P3	7500	6	45000	11	
3722	TR71313	370	2022-09-26	P3	7500	2	15000	3	

```
Gabung Data

[16] merge_data = pd.merge(transaction_data, customer_data, on = ['CustomerID'])

merge_data = pd.merge(merge_data, product_data.drop(columns = ['Price']), on = ['ProductID'])

merge_data = pd.merge(merge_data, store_data, on = ['StoreID'])
```

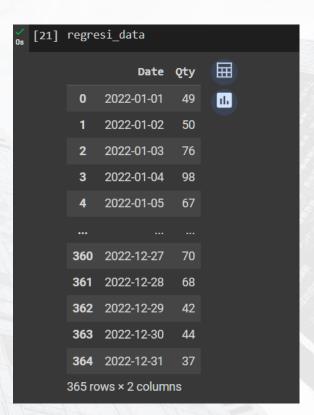


[17]	merge	e_data.head()																		
	т	ransactionID	CustomerID	Date	ProductID	Price	Qty	TotalAmount	StoreID	Age	Gender	Marital Status	Income	Longitude_x	Product Name	StoreName	GroupStore	Туре	Latitude	Longitude_y
	0	TR11369	328	2022- 01-01	P3	7500	4	30000	12	36		Married	10.0	NaN	Crackers	Prestasi Utama	Prestasi	General Trade	-2.0	NaN
	1	TR89318	183	2022- 07-17	P3	7500		7500	12	27		Single	0.0	NaN	Crackers	Prestasi Utama	Prestasi	General Trade	-2.0	NaN
	2	TR9106	123	2022- 09-26	P3	7500	4	30000	12	34		Married	4.0	NaN	Crackers	Prestasi Utama	Prestasi	General Trade	-2.0	NaN
	3	TR4331	335	2022- 08-01	P3	7500	3	22500	12	29		Single	4.0	NaN	Crackers	Prestasi Utama	Prestasi	General Trade	-2.0	NaN
	4	TR6445	181	2022- 10-01	P3	7500	4	30000	12	33		Married	9.0	NaN	Crackers	Prestasi Utama	Prestasi	General Trade	-2.0	NaN
	4																			

```
[18] #menyimpan data cleansed
merge_data.to_csv('cleaned_data.csv', index=False)

[19] cleaned_data = pd.read_csv('cleaned_data.csv')
```







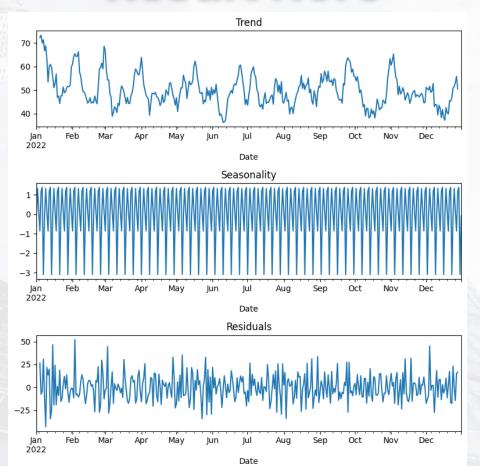
```
#digunakan untuk menghasilkan tiga subplot dalam satu gambar, masing-masing mewakili komponen utama dari analisis dekomposisi musiman pada data time series decomposed = seasonal_decompose(regresi_data.set_index('Date'))

plt.figure(figsize = (8,8))

plt.subplot(311)
    decomposed.trend.plot(ax = plt.gca())
    plt.title('Trend')
    plt.subplot(312)
    decomposed.seasonal.plot(ax = plt.gca())
    plt.title('Seasonality')
    plt.subplot(313)
    decomposed.resid.plot(ax = plt.gca())
    plt.title('Residuals')

plt.tight_layout()
```

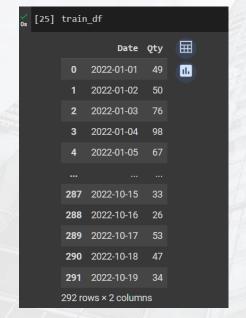


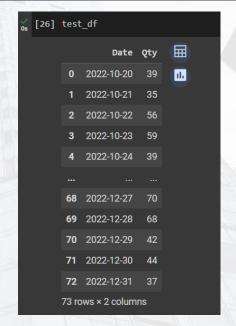




```
[24] #untuk membagi data time series menjadi dua set: set pelatihan (train) dan set pengujian (test)
cut_off = round(regresi_data.shape[0]*0.8)
train_df = regresi_data[:cut_off]
test_df = regresi_data[cut_off:].reset_index(drop=True)
train_df.shape, test_df.shape

((292, 2), (73, 2))
```







```
#membuat plot garis menggunakan library seaborn untuk memvisualisasikan data dalam set pelatihan (train) dan set pengujian (test) dari time series.
plt.figure(figsize=(20,5))
sns.lineplot(data=train_df, x=train_df['Date'], y=train_df['Qty']);
sns.lineplot(data=test_df, x=test_df['Date'], y=test_df['Qty']);
   120
   100
    80
 oth
    60
    40
    20
            2022-01
                                      2022-03
                                                                 2022-05
                                                                                            2022-07
                                                                                                                        2022-09
                                                                                                                                                   2022-11
                                                                                                                                                                              2023-01
```



#mengilustrasikan seberapa berkorelasinya suatu variabel dengan dirinya sendiri pada pergeseran waktu yang berbeda (lag). autocorrelation_plot(regresi_data['Qty']); 1.00 0.75 0.50 Autocorrelation 0.25 -0.25-0.50-0.75-1.0050 100 150 200 250 300 350 Lag



[28] #untuk mengambil baris-baris dari dataframe customer_data di mana kolom 'Marital Status' memiliki nilai yang kosong customer_data[customer_data['Marital Status'].isnull()]

	CustomerID	Age	Gender	Marital Status	Income	Longitude
9	10	34	1	NaN	4.0	NaN
415	416	27	1	NaN	3.0	NaN
442	443	33	1	NaN	9.0	NaN

	mbil baris-bari _data[merge_data				li mana		m 'Marital St	atus' mem	iliki	. nilai y	ang kosoi								
	TransactionID	CustomerID	Date	ProductID	Price	Qty	TotalAmount	StoreID	Age	Gender	Marital Status	Income	Longitude_x	Product Name	StoreName	GroupStore	Туре	Latitude	Lon
143	TR65931	443	2022- 05-19	P1	8800		44000	12	33		NaN	9.0	NaN	Choco Bar	Prestasi Utama	Prestasi	General Trade	-2.0	
175	TR92340	416	2022- 03-14	P7	9400	6	56400	12	27		NaN	3.0	NaN	Coffee Candy	Prestasi Utama	Prestasi	General Trade	-2.0	
295	TR28303	10	2022- 12-30	P10	15000	4	60000	12	34		NaN	4.0	NaN	Cheese Stick	Prestasi Utama	Prestasi	General Trade	-2.0	
387	TR84178	10	2022- 04-25	P3	7500		7500	13	34		NaN	4.0	NaN	Crackers	Buana	Buana	General Trade	-1.0	
458	TR78858	443	2022-	P5	4200	5	21000	13	33		NaN	9.0	NaN	Thai Tea	Buana	Buana	General	-1.0	



```
[30] #menghitung dan mengevaluasi kinerja model prediksi menggunakan metrik MAE dan RMSE dengan menggunakan contoh data aktual dan prediksi.
    def rmse(y actual, y pred):
        mse = mean squared error(y actual, y pred)
        rmse value = mse ** 0.5
        return rmse value
    def eval(y_actual, y_pred):
        mae value = mean absolute error(y actual, y pred)
        rmse value = rmse(y actual, y pred)
        print(f'Nilai MAE: {mae value}')
        print(f'Nilai RMSE: {rmse value}')
    y_actual = [3, 5, 7, 9]
    y \text{ pred} = [2.8, 5.2, 6.6, 9.3]
    eval(y actual, y pred)
    Nilai MAE: 0.275000000000000036
    Nilai RMSE: 0.2872281323269018
```

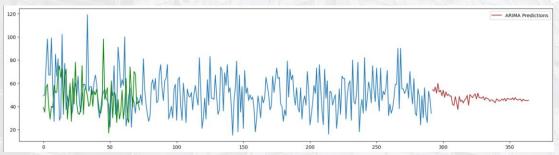


```
    Metode time series ARIMA

  [32] #mencetak informasi tentang kolom-kolom dan beberapa baris pertama dari dataframe train df dan test df
       print(train df.columns)
       print(test df.columns)
       print(train df.head())
       print(test df.head())
       Index(['Date', 'Qty'], dtype='object')
       Index(['Date', 'Qty'], dtype='object')
               Date Oty
       0 2022-01-01
       1 2022-01-02
       2 2022-01-03
       3 2022-01-04
       4 2022-01-05
               Date Oty
       0 2022-10-20
       1 2022-10-21
       2 2022-10-22 56
       3 2022-10-23
       4 2022-10-24
```



```
[33] #membangun model ARIMA, membuat prediksi, mengevaluasi performa prediksi, dan membuat plot visualisasi dari hasil prediksi
        y = train_df['Qty']
        ARIMAmodel = ARIMA(y, order=(40, 2, 1))
        ARIMAmodel fit = ARIMAmodel.fit()
        y pred = ARIMAmodel fit.get forecast(steps=len(test df))
        y pred df = y pred.conf int()
        y pred df['predictions'] = ARIMAmodel fit.predict(start=y pred df.index[0], end=y pred df.index[-1])
        y pred out = y pred df['predictions']
        eval(test df['Oty'], y pred out)
        plt.figure(figsize=(20, 5))
        plt.plot(train df['Qty'])
        plt.plot(test df['Qty'], color='green')
        plt.plot(y pred out, color='brown', label='ARIMA Predictions')
        plt.legend()
        plt.show()
        /usr/local/lib/python3.10/dist-packages/statsmodels/tsa/statespace/sarimax.py;966: UserWarning: Non-stationary starting autoregressive parame
          warn('Non-stationary starting autoregressive parameters'
        /usr/local/lib/python3.10/dist-packages/statsmodels/tsa/statespace/sarimax.py:978: UserWarning: Non-invertible starting MA parameters found.
          warn('Non-invertible starting MA parameters found.'
        /usr/local/lib/python3.10/dist-packages/statsmodels/base/model.py:607: ConvergenceWarning: Maximum Likelihood optimization failed to converge
         warnings.warn("Maximum Likelihood optimization failed to "
        Nilai MAE: 12.357079825645519
        Nilai RMSE: 15.472386623448042
```



2. Case Study



Machine Learning Clustering

- Tujuan dari pembuatan model machine learning ini adalah untuk dapat membuat cluster customer-customer yang mirip
 - Data cleansing terlebih dahulu, merubah tipe data supaya sesuai
 - Data merge untuk menggabungkan semua data
 - Membuat data baru untuk clustering, yaitu groupby by customerID lalu yang di aggregasi adalah :
 - Transaction id count
 - Qty sum
 - Total amount sum
 - Menggunakan metode clustering KMeans
 - Untuk proses queri data, dapat memilih alternatif pemrograman yang sesuai preferensi kamu.



- -- Model Machine Learning Clustering --

	TransactionID	CustomerID	Date	ProductID	Price	Qty	TotalAmount	StoreID	Age	Gender	Marital Status	Income	Longitude_x	Product Name	StoreName	GroupStore	Туре	Latitude	Longitude_y
0	TR11369	328	2022- 01-01	P3	7500	4	30000	12	36		Married	10.0	NaN	Crackers	Prestasi Utama	Prestasi	General Trade	-2.0	NaN
1	TR89318	183	2022- 07-17	P3	7500		7500	12	27		Single	0.0	NaN	Crackers	Prestasi Utama	Prestasi	General Trade	-2.0	NaN
2	TR9106	123	2022- 09-26	P3	7500	4	30000	12	34		Married	4.0	NaN	Crackers	Prestasi Utama	Prestasi	General Trade	-2.0	NaN
3	TR4331	335	2022- 08-01	P3	7500	3	22500	12	29		Single	4.0	NaN	Crackers	Prestasi Utama	Prestasi	General Trade	-2.0	NaN
4	TR6445	181	2022- 10-01	P3	7500	4	30000	12	33		Married	9.0	NaN	Crackers	Prestasi Utama	Prestasi	General Trade	-2.0	NaN



[35] #menghitung dan menampilkan matriks korelasi antar kolom-kolom dalam dataframe merge_data.corr()

<ipython-input-35-f0b9e936349b>:2: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will
merge_data.corr()

CustomerID 1.000000 -0.016423 -0.009755 -0.024915 0.004129 -0.025952 -0.009947 0.007537 NaN 0.001060 NaN Price -0.016423 1.000000 -0.353640 0.440632 -0.032863 0.014693 0.010705 0.001439 NaN -0.028502 NaN Qty -0.009755 -0.353640 1.000000 0.621129 0.014365 -0.027768 -0.010542 -0.029113 NaN -0.002765 NaN
Qty -0.009755 -0.353640 1.000000 0.621129 0.014365 -0.027768 -0.010542 -0.029113 NaN -0.002765 NaN
4 ,
TotalAmount -0.024915 0.440632 0.621129 1.000000 -0.010722 -0.016900 -0.008774 -0.025680 NaN -0.028421 NaN
StoreID 0.004129 -0.032863 0.014365 -0.010722 1.000000 -0.003872 -0.000189 0.001189 NaN 0.488507 NaN
Age -0.025952 0.014693 -0.027768 -0.016900 -0.003872 1.000000 -0.033183 0.485788 NaN 0.009051 NaN
Gender -0.009947 0.010705 -0.010542 -0.008774 -0.000189 -0.033183 1.000000 -0.072268 NaN -0.005635 NaN
Income 0.007537 0.001439 -0.029113 -0.025680 0.001189 0.485788 -0.072268 1.000000 NaN 0.014446 NaN
Longitude_x NaN NaN NaN NaN NaN NaN NaN NaN NaN Na
Latitude 0.001060 -0.028502 -0.002765 -0.028421 0.488507 0.009051 -0.005635 0.014446 NaN 1.000000 NaN
Longitude_y NaN NaN NaN NaN NaN NaN NaN NaN NaN Na



[38]	clu	ster_df.head	I()			
		CustomerID	TransactionID	Qty	TotalAmount	
	0	1	17	60	623300	ıl.
	1	2	13	57	392300	
	2	3	15	56	446200	
	3	4	10	46	302500	
	4	5	7	27	268600	

nID Q	ty Tot	alAmount	
17	60	623300	11.
13	57	392300	
15	56	446200	
10	46	302500	
7	27	268600	
16	59	485100	
18	62	577700	
18	68	587200	
11	42	423300	
13	42	439300	
	13	13 42	13 42 439300

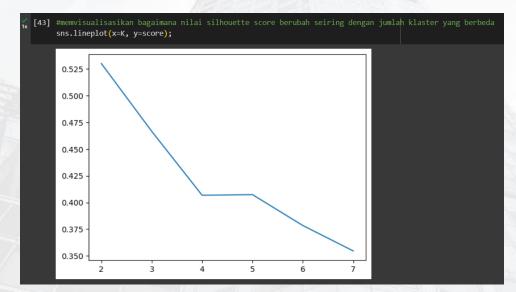


```
[40] #pra-pemrosesan data sebelum dilakukan analisis klasterisasi (clustering)
cluster_data = cluster_df.drop(columns = ['CustomerID'])
cluster_data_normalize = preprocessing.normalize(cluster_data)
```



```
[42] #algoritma K-Means untuk melakukan analisis klasterisasi pada data yang telah di normalisasi.
    K = range(2,8)
    fits = []
    score = []

for k in K:
    model = KMeans(n_clusters = k, random_state = 0, n_init = 'auto').fit(cluster_data_normalize)
    fits.append(model)
    score.append(silhouette_score(cluster_data_normalize, model.labels_, metric = 'euclidean'))
```





```
[44] fits[1]

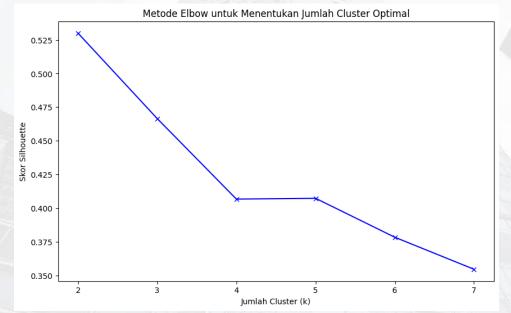
* KMeans

KMeans(n_clusters=3, n_init='auto', random_state=0)
```

[45] #menambahkan kolom 'cluster_label' ke dalam dataframe cluster_df yang berisi label klaster hasil dari analisis K-Means cluster_df['cluster_label'] = fits[2].labels_



```
[48] # Memvisualisasikan hasil siluet untuk menentukan jumlah cluster yang optimal plt.figure(figsize=(10, 6)) plt.plot(K, score, 'bx-') plt.xlabel('Jumlah Cluster (k)') plt.ylabel('Skor Silhouette') plt.title('Metode Elbow untuk Menentukan Jumlah Cluster Optimal') plt.show()
```





```
[49] # Pilih jumlah cluster berdasarkan hasil visualisasi siluet
selected_k = 3
selected_model = fits[selected_k - 2] # Indeks array dimulai dari 0
```

```
[50] # Tambahkan informasi klaster ke dalam DataFrame cluster_df
cluster_df['cluster_label'] = selected_model.labels_
```

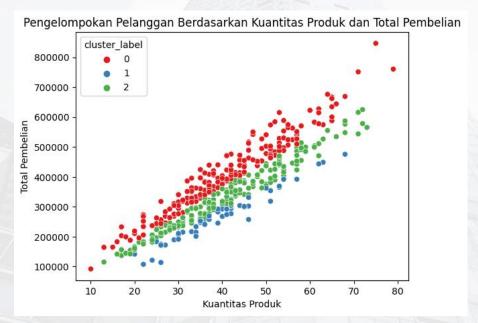
```
[52] #analisis klasterisasi
    cluster_analysis = cluster_df.groupby('cluster_label').agg({
        'CustomerID': 'count',
        'TransactionID': 'mean',
        'Qty': 'mean',
        'TotalAmount': 'mean'
})
```



```
[53] # Memberikan rekomendasi berdasarkan analisis klaster
     for cluster label, data in cluster analysis.iterrows():
         print(f"Cluster {cluster label}:")
         print(f"Jumlah Pelanggan: {data['CustomerID']}")
         print(f"Rata-rata Transaksi per Pelanggan: {data['TransactionID']:.2f}")
         print(f"Rata-rata Kuantitas Produk per Pelanggan: {data['Oty']:.2f}")
         print(f"Rata-rata Total Pembelian per Pelanggan: {data['TotalAmount']:.2f}")
         print("\n")
     Cluster 0:
     Jumlah Pelanggan: 205.0
     Rata-rata Transaksi per Pelanggan: 11.65
     Rata-rata Kuantitas Produk per Pelanggan: 40.88
     Rata-rata Total Pembelian per Pelanggan: 409866.83
     Cluster 1:
     Jumlah Pelanggan: 61.0
     Rata-rata Transaksi per Pelanggan: 10.34
     Rata-rata Kuantitas Produk per Pelanggan: 39.23
     Rata-rata Total Pembelian per Pelanggan: 268652.46
     Cluster 2:
     Jumlah Pelanggan: 181.0
     Rata-rata Transaksi per Pelanggan: 11.05
     Rata-rata Kuantitas Produk per Pelanggan: 41.56
     Rata-rata Total Pembelian per Pelanggan: 340511.05
```



```
[54] # Visualisasi klaster dengan scatter plot
sns.scatterplot(data=cluster_df, x='Qty', y='TotalAmount', hue='cluster_label', palette='Set1')
plt.xlabel('Kuantitas Produk')
plt.ylabel('Total Pembelian')
plt.title('Pengelompokan Pelanggan Berdasarkan Kuantitas Produk dan Total Pembelian')
plt.show()
```



Case Study



Peserta dapat membuat dashboard di tableau

Sebelum membuat dashboard terlebih dahulu membuat worksheet sebanyak 4.

- Worksheet 1 Jumlah qty dari bulan ke bulan
- Worksheet 2 Jumlah total amount dari hari ke hari
- Worksheet 3 Jumlah penjualan (qty) by product
- Worksheet 4 Jumlah penjualan (total amount) by store name
- Setelah itu bisa membuat dashboard dengan menggabungkan 4 worksheet.



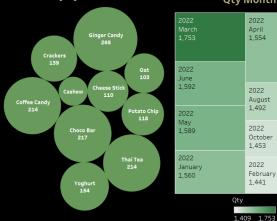


TRANSACTIONS DASHBOARD

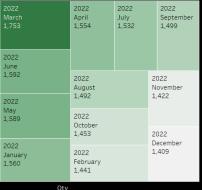
Total Amount Day by Day



Qty by Product



Qty Month by Month



Amount by Store Name



Total Amount

162,043,000



Link Tableau

https://public.tableau.com/views/TRANSACTIONSDASHBOARD_KARIN/Dashboard1?:language=e n-US&:display_count=n&:origin=viz_share_link



Link Github

https://github.com/KarinaKhPutri/FinalTask_Kalbe_DS_KarinaKhairunnisaPutri



Link Video Presentation

https://drive.google.com/drive/folders/1h0W44zwPoUjrexL-Vx-9_RghuaebhXm0?usp=sharing

Thank You





