Time Series & Clustering

Machine Learning Project

Presented by Arif Styawan



ABOUT

I am a mathematics student with a strong interest in pursuing a career in the field of data, focusing on **data analysis** and **machine learning** skills. I have proficiency in analytical skills and tools such as **Excel/Spreadsheets**, **Python**, **SQL**, and **Tableau**. I have a strong passion for learning and enjoy exploring new things.

EXPERIENCE

Bangkit Academy 2023

Machine Learning Student (Feb - July)

Pacmann Academy Scholarship

Analytics & Data Science Awardee (Feb - present)

TOOLS





linkedin.com/in/arifstyawan



arifstyawaan@gmail.com



github.com/arifstyawan

Overview

Saya berperan sebagai seorang **Data Scientist** di Kalbe Nutritionals dan sedang mendapatkan project baru dari tim inventory dan tim marketing.

Tujuan dari project ini yaitu:

- 1. Mengetahui **perkiraan** quantity product yang terjual sehingga tim inventory dapat membuat stock persediaan harian yang cukup.
- 2. Membuat cluster/segment customer berdasarkan beberapa kriteria.

Challenge

Objektif yang harus dicapai antara lain:

- Melakukan data ingestion dan exploratory data analysis ke dalam dbeaver, dengan query sebagai berikut.
 - a. Berapa rata-rata umur customer jika dilihat dari marital statusnya?
 - b. Berapa rata-rata umur customer jika dilihat dari gender nya?
 - c. Tentukan nama store dengan total quantity terbanyak!
 - d. Tentukan nama produk terlaris dengan total amount terbanyak!
- 2. Melakukan data ingestion dan membuat dashboard di Tableau
- 3. Membuat model machine learning time series ARIMA
- 4. Membuat model machine learning clustering dengan **K-Means**

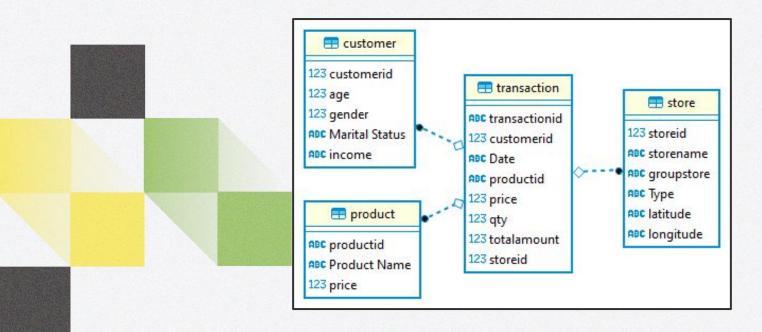
1

Melakukan **data ingestion** dan **exploratory data analysis** ke dalam dbeaver

Dataset yang digunakan berupa 4 file dengan format .csv yaitu:

Case Study - Customer.csv, Case Study - Product.csv, Case Study - Transaction.csv, dan Case Study - Store.csv.

Kolom untuk setiap tabel beserta hubungan antar tabel sebagai berikut.



• Berapa rata-rata umur customer jika dilihat dari marital statusnya?

	ABC Marital Status	123 average_age 🔻
1		31.3333333333
2	Married	43.0382352941
3	Single	29.3846153846

Rata-rata umur untuk customer dengan status **Married** adalah **43 tahun** dan rata-rata dengan status **Single** adalah **29,4 tahun**.

Terdapat NULL data dengan rata-rata 31,3 tahun

Berapa rata-rata umur customer jika dilihat dari gendernya?

	123 gender	*	123 average_age	
1		0	40.32644628	1
2		1	39.141463414	6

Rata-rata umur untuk customer **wanita** adalah **40 tahun** dan customer **pria** adalah **39 tahun**.

Tentukan nama store dengan total quantity terbanyak!

	ABC storename		123 total_qty	-
1	Lingga	Ī	2	2,777
2	Sinar Harapan		2	2,588

Nama store dengan total quantity terbanyak adalah **Lingga** dengan total quantity sebanyak **2.777**. Disusul oleh Store Sinar Harapan.

Tentukan nama produk terlaris dengan total amount terbanyak!

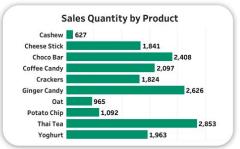
	ABC Product Name	•	123 total_amount
1	Cheese Stick		27,615,000
2	Choco Bar		21,190,400

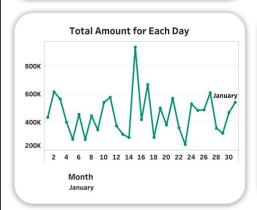
Nama produk terlaris dengan total amount terbanyak adalah **Cheese Stick** dengan total amount sebesar **27.615.000**. Disusul oleh produk Choco Bar.

Melakukan data ingestion dan membuat dashboard di Tableau



Jan Feb Mar Apr May Jun Jul Aug Sep Oct Nov Dec







Dashboard di samping terdiri dari 4 worksheet, sebagai berikut :

- 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

3

Membuat model machine learning time series ARIMA

IMPORT LIBRARY

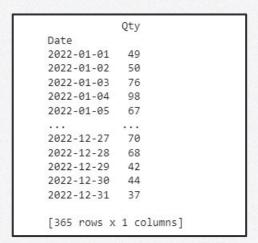
Library yang diimport berupa pandas, datetime, matplotlib, scikit-learn, dan library lainnya.

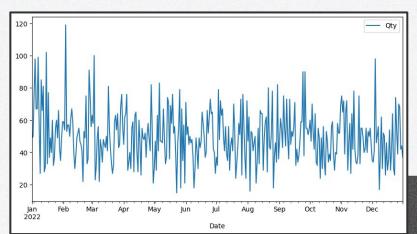
LOAD DATASET

Dataset yang digunakan yaitu 'Case Study - Transaction.csv'.

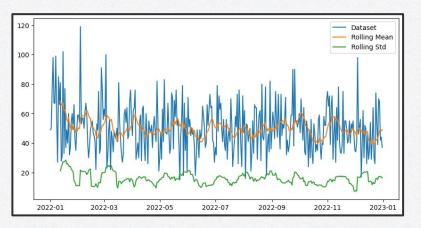
DATA CLEANSING & PREPROCESSING

- a. Mengecek **missing values** (tidak ada missing values)
- b. Mengecek tipe data kemudian mengubah kolom 'Date' menjadi tipe datetime
- c. Melakukan **grouping** berdasarkan 'Date' dan menjumlahkan 'Qty'.
- d. Mengubah kolom 'Date' menjadi **indeks** dan membuat **plot** pada 'Qty'.





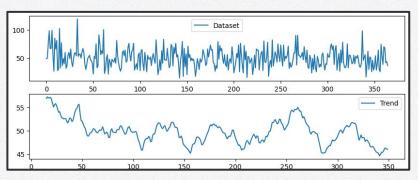
CHECK STATIONARITY

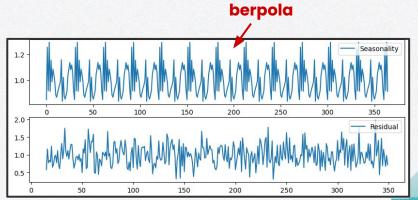


Mengecek stationarity menggunakan **rolling mean** dan **rolling standard deviation**.

Karena grafik rolling mean dan rolling std cenderung **konstan** dan tidak memiliki tren, maka data bersifat **stasioner**

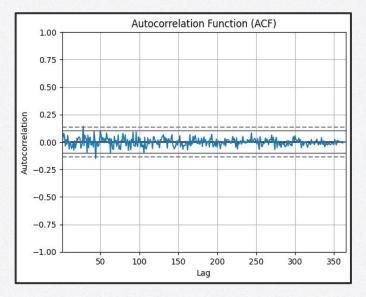
CHECK SEASONALITY

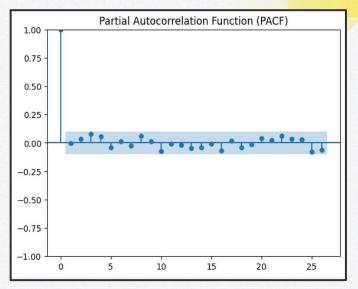




Terlihat bahwa data bersifat **seasonal (musiman)** karena memiliki pola setiap periode tertentu. Karena bersifat seasonal, maka nilai **d = 1**

• CHECK ACF(p) AND PACF(q)





Dari grafik **ACF** terlihat bahwa garis biru memotong 5 garis horizontal, sehingga **p = 5**Dari grafik **PACF** terlihat bahwa terdapat 1 titik terjauh/mendekati 1, sehingga **q = 1**

Jadi, dari check seasonality, ACF, dan PACF didapatkan d = 1, p = 5, q = 1

SPLIT DATASET Split dataset dengan pembagian train set = 335 data dan test set = 30 data.

BUILD MODEL

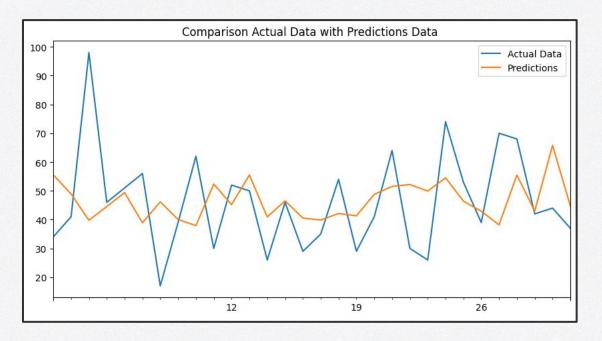
Model machine learning yang dibuat merupakan **Seasonal ARIMA (SARIMA)**. Training model menggunakan data training yang telah di split. Berikut model summary nya:

Dep. Variab	le:			Otv	No.	Observations:		335
Model:	SARI	IMAX(5, 1,	1)x(5, 1, 1	, 30)	Log	Likelihood		-1311.51
Date:			Wed, 26 Jul					2649.03
Time:			07:	27:30	BIC			2697.356
Sample:			01-01	-2022	HQIC			2668.369
SERVICE TO THE COURT OF THE			- 12-01	-2022				
Covariance				opg				
						[0.025		
ar.L1	-0.0091	0.065	-0.140	0.8	 88	-0.136	0.118	
						-0.147		
						-0.049		
						-0.023		
ar.L5	-0.0178	0.062	-0.287	0.7	74	-0.139	0.104	
ma.L1	-0.9867	0.033	-30.266	0.0	00	-1.051	-0.923	
ar.S.L30	-0.0276	0.167	-0.165	0.8	69	-0.354	0.299	
ar.S.L60	-0.0288	0.143	-0.202	0.8	40	-0.308	0.251	
ar.S.L90	-0.0805	0.127	-0.633	0.5	27	-0.330	0.169	
ar.S.L120	-0.0705	0.114	-0.616	0.5	38	-0.295	0.154	
ar.S.L150	-0.1235	0.111	-1.112	0.2	66	-0.341	0.094	
ma.S.L30	-0.9944	6.695	-0.149	0.8	82	-14.117	12.128	
sigma2	242.5074	1586.077	0.153	0.8	78	-2866.146	3351.161	
								===
Ljung-Box (L1) (Q):		0.04					.98
Prob(Q):			0.85):		6	1000
	sticity (H):		0.75					.09
Prob(H) (tw	o-sided):		0.15	Kurtosi	5:		9	.21

MODEL EVALUATION

Testing data menggunakan data testing. Selanjutnya dilakukan prediksi model untuk 30 data terakhir.

Grafik di samping adalah perbandingan antara actual data dengan prediksi model.



TES RMSE

Didapatkan RMSE sebesar **18,544** Angka ini menandakan bahwa model sudah berjalan cukup baik.

FORECASTING FOR NEXT MONTH

Berikut merupakan hasil prediksi jumlah quantity untuk **31 hari ke depan** dari hari ke hari. Berikut terdapat juga visualisasi training data, actual data, dan prediksi 31 hari kedepan.

2023-01-01 2023-01-02	56.098186 49.546683 45.055179	2023-01-16 2023-01-17 2023-01-18	35.546540 41.218269 38.742441	120 - Training Data Testing/Actual Data Forecasting
2023-01-03 2023-01-04 2023-01-05	43.692128 48.016123	2023-01-19 2023-01-20 2023-01-21	50.008640 50.975427 53.648984	80 -
2023-01-06 2023-01-07 2023-01-08	42.033181 44.946107 37.256635	2023-01-21 2023-01-22 2023-01-23 2023-01-24	45.616540 59.438613 48.081410	60 - 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
2023-01-09 2023-01-10 2023-01-11	40.647687 48.043476 44.107906	2023-01-25 2023-01-26 2023-01-27	41.508262 34.955224 55.409811	40 -
2023-01-12 2023-01-13	53.294111 41.294605	2023-01-28 2023-01-29	42.621646 65.269550	20 -
2023-01-14 2023-01-15	49.787157 39.550210	2023-01-30 2023-01-31	37.982738 54.615704	Jan Feb Mar Apr May Jun Jul Aug Sep Oct Nov Dec Jan 2022 2023 Date

4

Membuat model machine learning clustering K-Means

IMPORT LIBRARY

Library yang diimport berupa pandas, numpy, matplotlib, scikit-learn, dan seaborn.

LOAD DATASET

Dataset yang digunakan yaitu 'Case Study - Transaction.csv'.

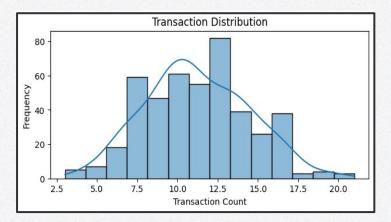
DATA CLEANSING & PREPROCESSING

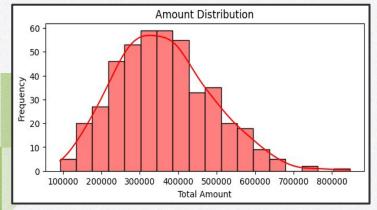
- a. Mengecek **missing values** (tidak ada missing values)
- b. Membuat data baru dengan **groupby** berdasarkan customerID dan di**aggregasi** dengan Transaction id (count), Qty (sum), Total amount (sum)
- c. Didapat data baru dengan 447 baris dan 4 kolom, diantaranya **CustomerID**, **TransactionCount**, **TotalQty**, dan **TotalAmount**
- d. Mengubah kolom CustomerID menjadi indeks kolom
- e. Mengecek konsistensi format/tipe data

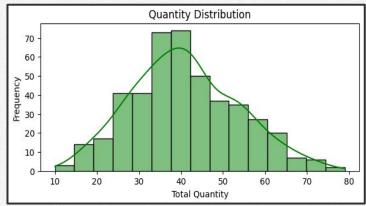
TransactionCount = Tipe integer
TotalQty = Tipe integer
TotalAmount = Tipe integer

CHECK DATA DISTRIBUTION

Berikut visualisasi persebaran dari tiap data, yaitu Transaction Count, Total Quantity, dan Total Amount.



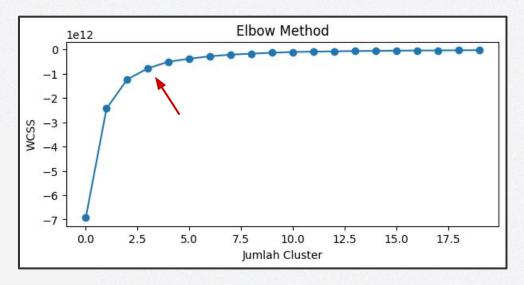




BUILD K-MEANS MODEL CLUSTERING

Langkah pembuatan K-Means model sebagai berikut

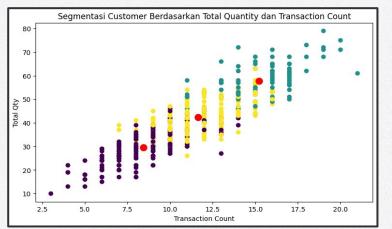
- 1. Mencari nilai jumlah **n cluster** yang paling optimal.
- 2. Menggunakan looping dengan n = 1 sampai n = 20
- 3. Visualisasikan hasilnya kemudian analisis menggunakan Elbow Method

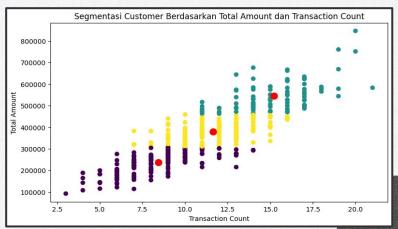


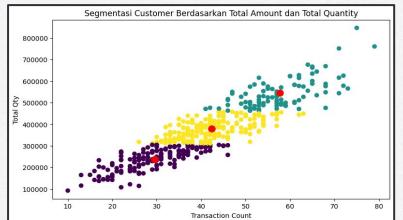
- 4. Dengan Elbow Method, didapatkan **n = 3** yang paling optimal
- 5. Build model K-Means menggunakan **n = 3**

```
model = KMeans(n_clusters=3)
model.fit(df)
```

CLUSTERING Berikut visualisasi persebaran data berdasarkan segmentasi/clustering







DASHBOARD VISUALIZATION:

https://public.tableau.com/views/SalesDashboardVisualization/Dashboard1?:lan guage=en-US&publish=yes&:display_count=n&:origin=viz_share_link

GITHUB REPOSITORY:

https://github.com/arifstyawan/TimeSeries-Clustering

VIDEO PRESENTATION:

https://youtu.be/3KAVt50YO_U

Thank You