# Kalbe Nutritionals Data Scientist
## Project Based Internship Program
## Machine Learning Project Using ARIMA and K-Means

Presented By Radisha Fanni Sianti

# About Me

" I am a fresh graduate of Master of Statistics and Bachelor of Mathematics from Institut Teknologi Sepuluh Nopember. I have a passion for continuing to learn and develop myself by constantly trying to broaden my knowledge in statistics and other related disciplines. I'm enthusiastic about deriving valuable insights from data and leveraging them to facilitate well-informed decision-making. "

# Table of contents

# 01
# Overview

# Overview

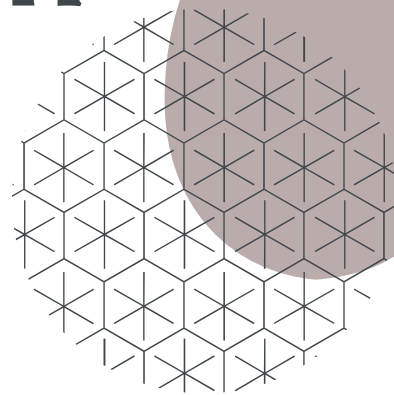This project comes from the inventory team and marketing team of Kalbe Nutritionals.

The inventory team wants to know the estimated quantity of products sold so that the inventory team can create sufficient daily inventory. To overcome this problem, we can use the machine learning method to forecast product quantity time series using the ARIMA model.

The marketing team wants to segment customers based on several criteria, which will later be used to provide personalized promotion and sales treatment. To overcome this problem, we can utilize the machine learning clustering method using the K-Means model.
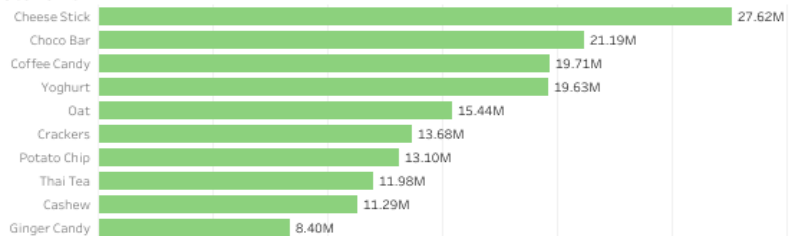
# 02

# Visualization
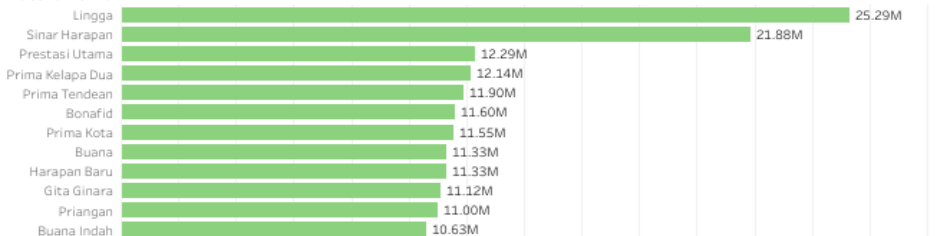
# Sales Performance Dashboard
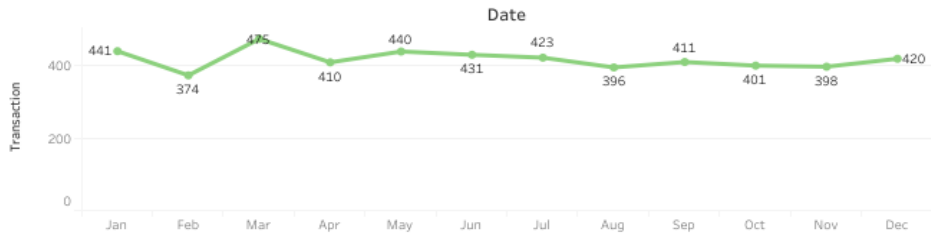
## Total Amount Per Day

### Date

Total Amount values per day (1–31): 6.9M, 6.3M, 5.2M, 4.8M, 5.1M, 4.6M, 5.4M, 4.4M, 5.8M, 5.0M, 6.0M, 4.9M, 5.2M, 5.2M, 5.7M, 4.7M, 5.4M, 5.3M, 5.1M, 4.3M, 5.4M, 4.9M, 4.9M, 6.1M, 5.7M, 5.8M, 5.9M, 5.7M, 4.4M, 4.7M, 3.2M

## Total Amount By Product

| Product Name1 | Total Amount |
|---|---|
| Cheese Stick | 27.62M |
| Choco Bar | 21.19M |
| Coffee Candy | 19.71M |
| Yoghurt | 19.63M |
| Oat | 15.44M |
| Crackers | 13.68M |
| Potato Chip | 13.10M |
| Thai Tea | 11.98M |
| Cashew | 11.29M |
| Ginger Candy | 8.40M |

## Total Amount By Store Name

| Store Name | Total Amount |
|---|---|
| Lingga | 25.29M |
| Sinar Harapan | 21.88M |
| Prestasi Utama | 12.29M |
| Prima Kelapa Dua | 12.14M |
| Prima Tendean | 11.90M |
| Bonafid | 11.60M |
| Prima Kota | 11.55M |
| Buana | 11.33M |
| Harapan Baru | 11.33M |
| Gita Ginara | 11.12M |
| Priangan | 11.00M |
| Buana Indah | 10.63M |

## Total Transaction By Month

### Date

| Month | Transaction |
|---|---|
| Jan | 441 |
| Feb | 374 |
| Mar | 475 |
| Apr | 410 |
| May | 440 |
| Jun | 431 |
| Jul | 423 |
| Aug | 396 |
| Sep | 411 |
| Oct | 401 |
| Nov | 398 |
| Dec | 420 |

## Total Quantity Per Month

### Date

| Month | Quantity |
|---|---|
| Jan | 1,560 |
| Feb | 1,441 |
| Mar | 1,753 |
| Apr | 1,554 |
| May | 1,589 |
| Jun | 1,592 |
| Jul | 1,532 |
| Aug | 1,492 |
| Sep | 1,499 |
| Oct | 1,453 |
| Nov | 1,422 |
| Dec | 1,409 |

# 03
# Data Preprocessing

# Data Preprocessing

## 1. Check Data Type
"Date" column should be of the datetime data type not object.
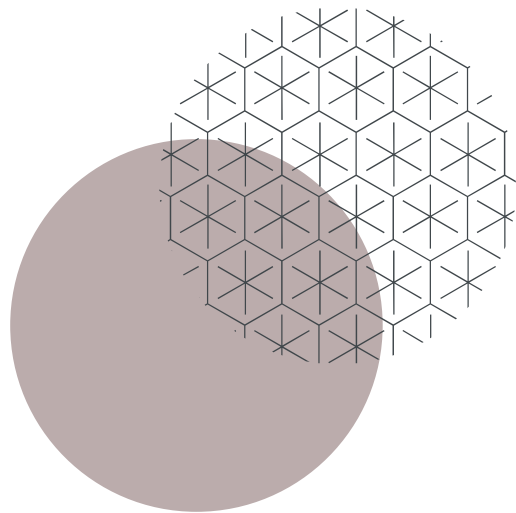
## 2. Check Duplicated Data
There is no duplicated data.

## 3. Check Missing Value
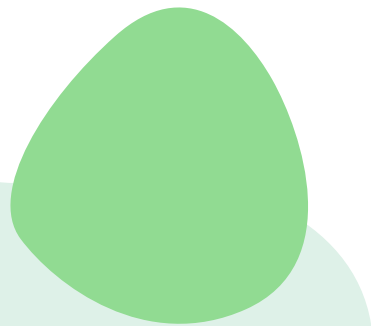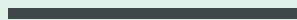There is no missing value.

## 4. Preparing Data for Machine Learning Modeling
- ARIMA Modeling using df_inventory (Date, Qty).
- K-Means Modeling using df_marketing (CustomerID, Total_Transaction, Total_Quantity, Total_Amount)

# 04

# EDA

# Average Customer Age Based on Marital Status

| | ABC Marital Status | 123 average of age |
|---|---|---|
| 1 | | 31.3333333333 |
| 2 | Married | 43.0382352941 |
| 3 | Single | 29.3846153846 |

# Average Customer Age Based on Gender

| | 123 gender | 123 average of age |
|---|---|---|
| 1 | 0 | 40.326446281 |
| 2 | 1 | 39.1414634146 |

# Total Quantity Based On Store Name

| | 123 storeid | ABC storename | 123 total quantity |
|---|---|---|---|
| 1 | 9 | Lingga | 1,439 |
| 2 | 12 | Prestasi Utama | 1,395 |
| 3 | 3 | Prima Kota | 1,358 |
| 4 | 6 | Lingga | 1,338 |
| 5 | 11 | Sinar Harapan | 1,331 |

# Total Amount Based On Product Name

| | ABC productid | ABC Product Name | 123 total amount |
|---|---|---|---|
| 1 | P10 | Cheese Stick | 27,615,000 |
| 2 | P1 | Choco Bar | 21,190,400 |
| 3 | P7 | Coffee Candy | 19,711,800 |
| 4 | P9 | Yoghurt | 19,630,000 |
| 5 | P8 | Oat | 15,440,000 |

# EDA Inventory Data

## Daily



## Monthly

# EDA Inventory Data



There are several outliers
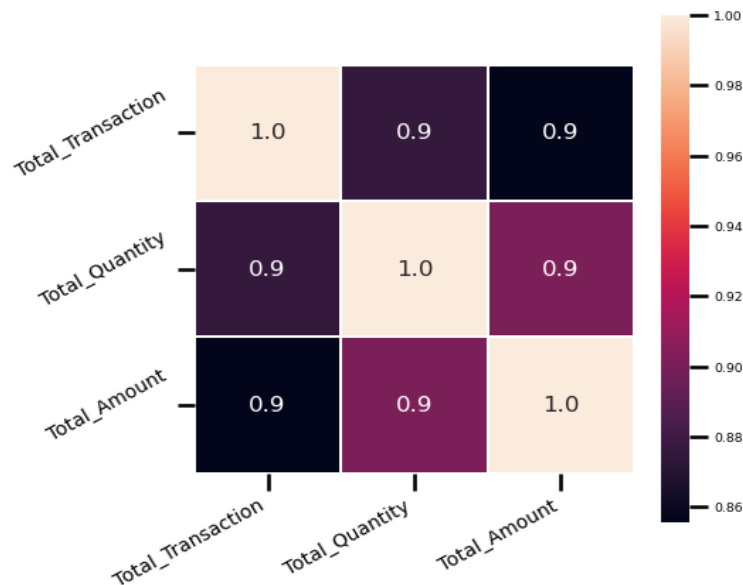


Approximately Symmetrical

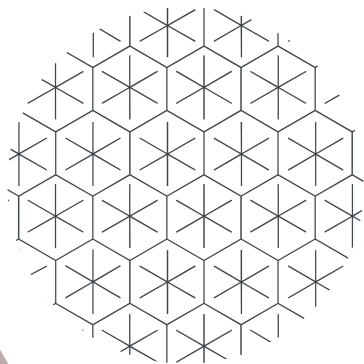# EDA Marketing Data



There is a very strong correlation between these three features, with a correlation of 0.9. It indicates that there is a close relationship between these variables.
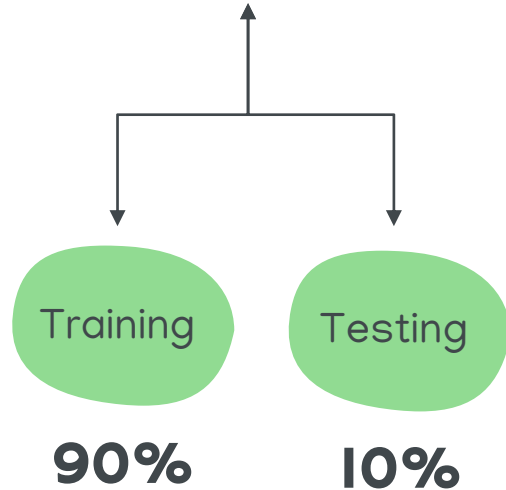
# 05
# Machine Learning

# I. ARIMA Modeling

**Inventory Data**

Training

Testing

**90%**          **10%**
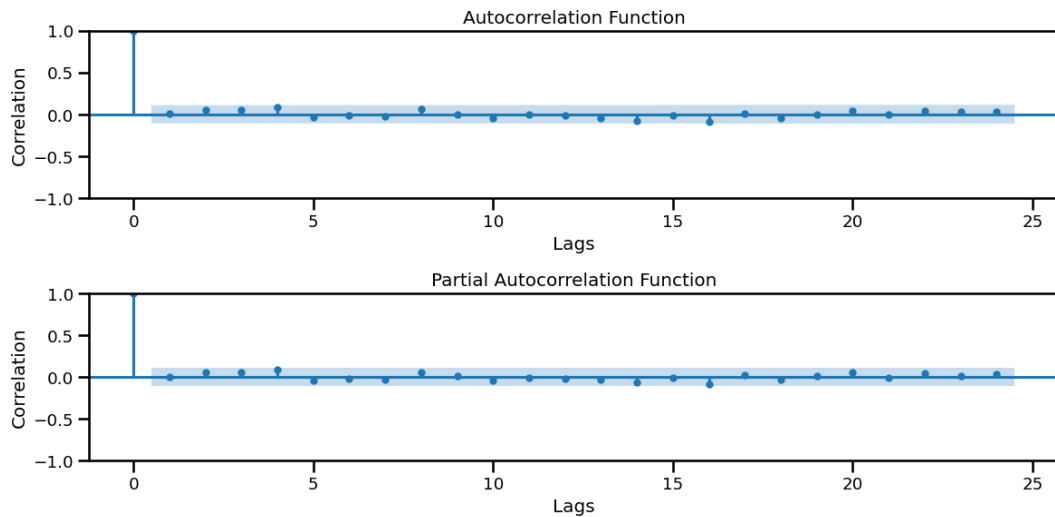
**ADF Test**

```
ADF test Total Quantity
ADF Statistic: -17.978891
p-value: 0.000000
Critical Values:
        1%: -3.451
        5%: -2.870
        10%: -2.572
Conclusion : Stasionary Data
```

# I. ARIMA Modeling
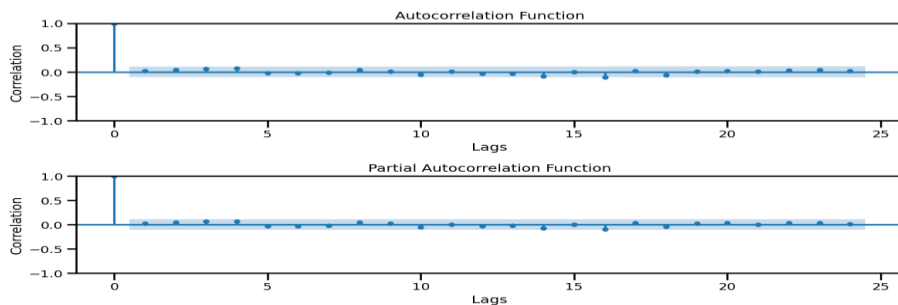
## ACF and PACF Inventory Data

# I. ARIMA Modeling

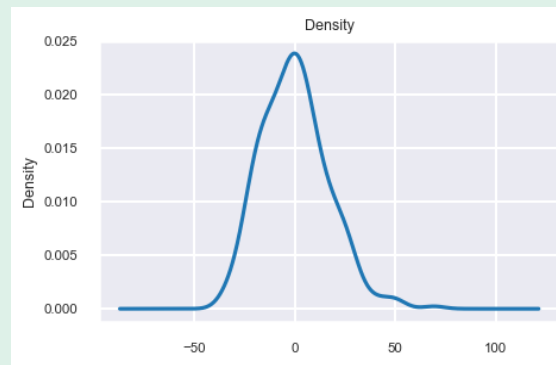| Model | Significance of Parameter |
|---|---|
| ARIMA (1,0,0) | There are parameters that are not significant |
| ARIMA (0,0,1) | There are parameters that are not significant |
| ARIMA (1,0,1) | All parameters are significant |
| ARIMA (2,0,2) | There are parameters that are not significant |
| ARIMA (3,0,3) | There are parameters that are not significant |
| ARIMA (4,0,4) | There are parameters that are not significant |

# I. ARIMA Modeling

## White Noise

### ACF and PACF Residuals



Lower lags barely show any significant 'spikes'. This indicates that the residuals are close to white noise. It can be concluded that ARIMA (4,0,4) meets the white noise assumption.
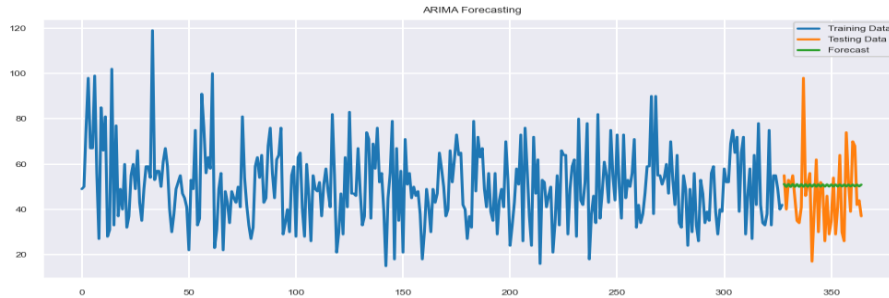
## Normality



The normality test results show that the residuals are not normally distributed

# I. ARIMA Modeling

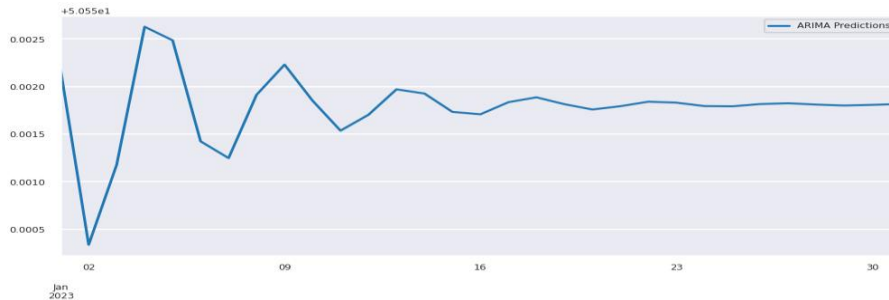## Forecasting Testing Data



```
MAE   : 12.675885484467079
MAPE  : 0.3416714532952726
RMSE  : 16.21730311234308
```

## Forecasting for The next 31 Days



The forecasting results show that for the next period the total quantity per day will be around 50.

# 2. K-Means Modeling

## Marketing Data

| | Total_Transaction | Total_Quantity | Total_Amount |
|---|---|---|---|
| 0 | 17 | 60 | 623300 |
| 1 | 13 | 57 | 392300 |
| 2 | 15 | 56 | 446200 |
| 3 | 10 | 46 | 302500 |
| 4 | 7 | 27 | 268600 |

Data must be standardized because the data scale is different

## marketing data that has been standardized

| | Total_Transaction | Total_Quantity | Total_Amount |
|---|---|---|---|
| 0 | 1.78 | 1.50 | 2.09 |
| 1 | 0.55 | 1.26 | 0.24 |
| 2 | 1.16 | 1.18 | 0.67 |
| 3 | -0.38 | 0.40 | -0.48 |
| 4 | -1.31 | -1.09 | -0.75 |

# 2. K-Means Modeling



Elbow Method of K-means Clustering
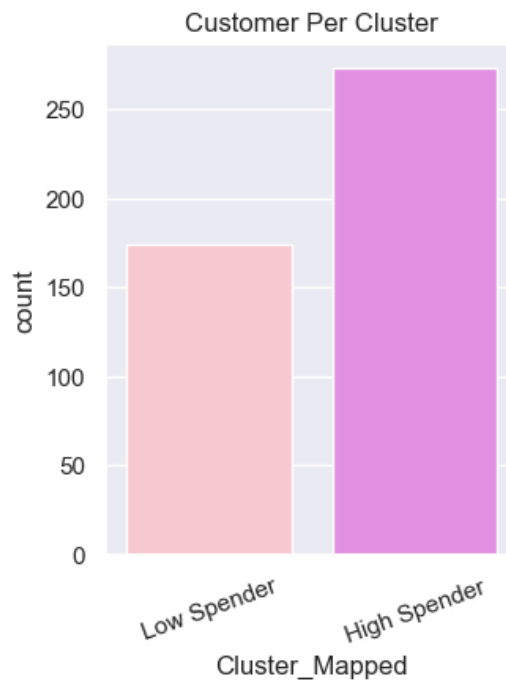


Silhouette Score of K-means Clustering

In the Elbow graph, it can be seen that the elbow fracture occurs when the clusters are equal to 2 and 3.
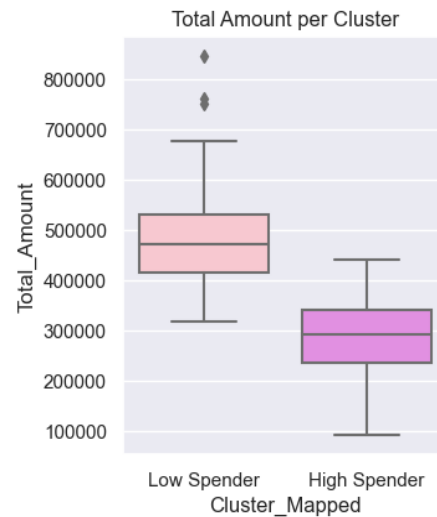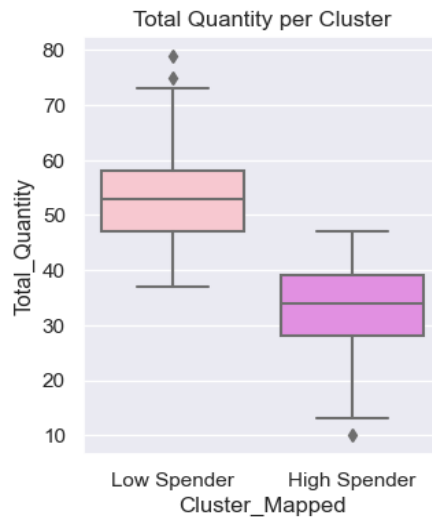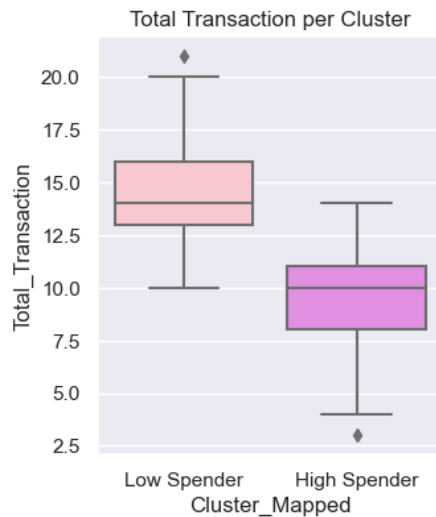
# 2. K-Means Modeling

## K-Means 2 Clusters



High Spender    273
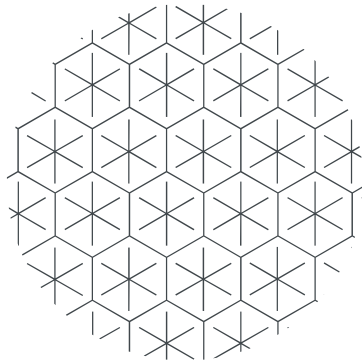Low Spender     174

# 2. K-Means Modeling

## Boxplot of Each Features

# 06

# Insight
# &
# Recommendations

# Insight

## Inventory Data

1. There is interesting pattern in monthly time series plot. This pattern may reflect a seasonal pattern or annual trend that deserves further scrutiny.
2. The best ARIMA modeling result for total quantity is ARIMA (1,0,1). The forecasting results show that for the next period the total quantity per day will be around 50.

## Marketing Data

1. The results from the Heatmap show that there is a very strong correlation between these three features, with a correlation of 0.9. It indicates that there is a close relationship between these variables.
2. The cluster analysis results formed 2 clusters, namely low spenders and high spenders. High spenders have a greater number of customers who are classified as high spenders than the number of customers who are classified as low spenders.

# Recommendations

**Inventory Data**

1. The inventory team can perform more in-depth analysis to identify the cause of the pattern.
2. The existence of outliers in the data is important to pay attention to. The inventory team needs to conduct further analysis of these outliers to understand why.
3. ARIMA (1,0,1) The model has residuals that meet the white noise assumption but do not meet the normality assumption

**Marketing Data**

1. Marketing team can design different marketing strategies for each group
1. Marketing team can develop more premium products or services for the high spender. More affordable products or bundle packages may tempt low spenders.
2. Marketing team can provide reward and loyalty programs to motivate low spender customers to increase their spending.

# Thanks

linkedin.com/in/radishafannis

radishafs@gmail.com