# Customer Segmentation and Forecasting Sales Quantity

Kalbe Nutritionals Data Science

Presented by
Tito

# Case Study

You are currently getting a new project from the **inventory team** and **marketing team**.
- From the **marketing team** you are asked to create customer clusters/segments based on several criteria.
- From the **inventory team**, you are asked to help predicting the number of sales (quantity) from the total all Kalbe products.

# Dataset

Two types of customer segmentation are conducted, **Behavioral Segmentation** and **Demographic Segmentation**.

| CustomerID | Frequency | Total_Quantity | Total_Amount_Sales | Recency | Age | Gender | Marital Status | Income |
|---|---|---|---|---|---|---|---|---|
| 1 | 17 | 60 | 623300 | 53 | 55 | Man | Married | 5.12 |
| 2 | 13 | 57 | 392300 | 97 | 60 | Man | Married | 6.23 |
| 3 | 15 | 56 | 446200 | 10 | 32 | Man | Married | 9.17 |
| 4 | 10 | 46 | 302500 | 5 | 31 | Man | Married | 4.87 |
| 5 | 7 | 27 | 268600 | 91 | 58 | Man | Married | 3.57 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 443 | 16 | 59 | 485100 | 63 | 33 | Man | Married | 9.28 |
| 444 | 18 | 62 | 577700 | 38 | 53 | Woman | Married | 15.31 |
| 445 | 18 | 68 | 587200 | 42 | 51 | Woman | Married | 14.48 |
| 446 | 11 | 42 | 423300 | 11 | 57 | Woman | Married | 7.81 |
| 447 | 13 | 42 | 439300 | 26 | 54 | Man | Married | 20.37 |

Behavioral data
- Number of transactions.
- Total Quantity.
- Total Amount Sales.
- Number of days since last transaction.

Demographic data
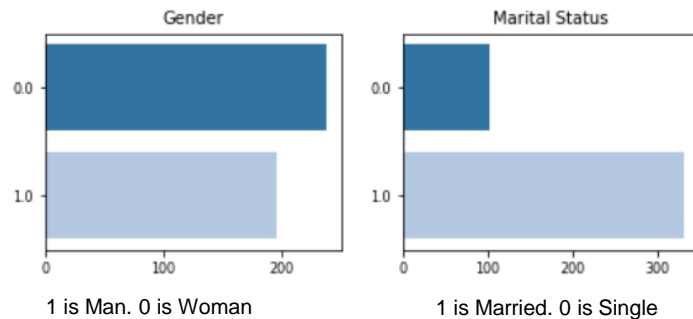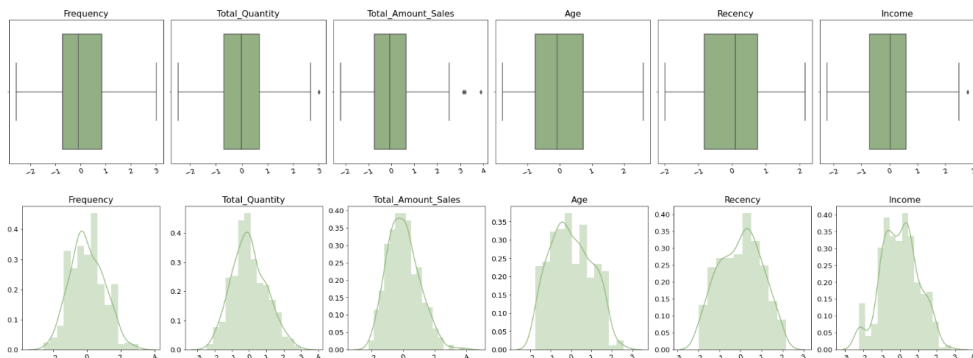- Age
- Gender
- Marital Status
- Income

KMeans

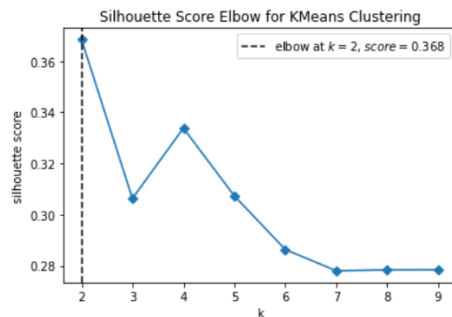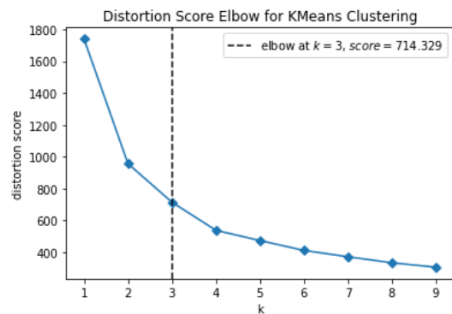KPrototypes

# Data Cleaning and Preprocessing



- Removing outliers
- Data transformation to make the distribution close to normal
- Standardization to make the features space having the same scale
- Label Encoding for categorical variables

| | Frequency | Total_Quantity | Total_Amount_Sales | Age | Recency | Income | Marital Status | Gender |
|---|---|---|---|---|---|---|---|---|
| 0 | 1.773511 | 1.482019 | 2.083088 | 1.218149 | 0.989894 | -0.492632 | 1.0 | 1.0 |
| 1 | 0.528691 | 1.245034 | 0.213660 | 1.621540 | 1.783630 | -0.245929 | 1.0 | 1.0 |
| 2 | 1.151101 | 1.166039 | 0.649860 | -0.637451 | -0.703350 | 0.323560 | 1.0 | 1.0 |
| 3 | -0.404924 | 0.376089 | -0.513070 | -0.718129 | -1.214207 | -0.551401 | 1.0 | 1.0 |
| 4 | -1.338539 | -1.124817 | -0.787415 | 1.460184 | 1.694733 | -0.880941 | 1.0 | 1.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 430 | 1.462306 | 1.403024 | 0.964668 | -0.556773 | 1.206197 | 0.343043 | 1.0 | 1.0 |
| 431 | 2.084716 | 1.640010 | 1.714058 | 1.056792 | 0.596408 | 1.278902 | 1.0 | 0.0 |
| 432 | 2.084716 | 2.113980 | 1.790939 | 0.895436 | 0.711674 | 1.162541 | 1.0 | 0.0 |
| 433 | -0.093719 | 0.060109 | 0.464536 | 1.379505 | -0.624566 | 0.072791 | 1.0 | 0.0 |
| 434 | 0.528691 | 0.060109 | 0.594020 | 1.137471 | 0.182779 | 1.928694 | 1.0 | 1.0 |





1 is Man. 0 is Woman

1 is Married. 0 is Single

# Modeling – Elbow Method

Behavioral segmentation
(Kmeans)

Demographic segmentation
(KPrototypes)



- The recommended number of clusters are 3 or 4 for both models based on the the elbow method.
- 4 clusters are chosen since the Silhouette score is higher compared to 3 clusters.
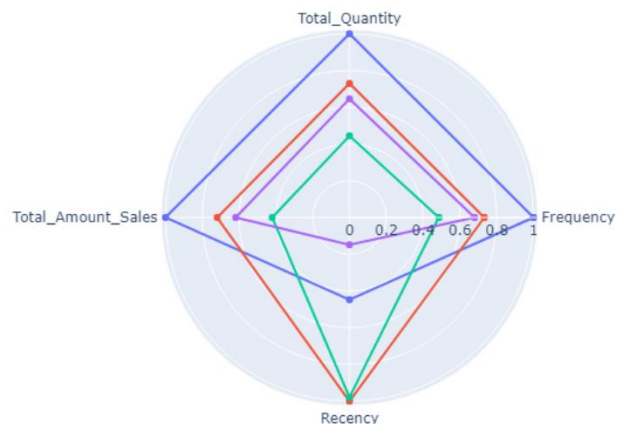
# Modeling – Interpretation

Behavioral segmentation
(Kmeans)


Radar Plot of the resulting 4 clusters

| Behavioral_Segmentation | Recency mean | Frequency mean | Total_Quantity mean | Total_Amount_Sales mean | count |
|---|---|---|---|---|---|
| Diamond | 20.42 | 15.48 | 57.92 | 524915.38 | 104 |
| Gold | 45.63 | 11.34 | 42.26 | 377484.96 | 133 |
| Potential | 6.78 | 10.52 | 37.38 | 324564.08 | 103 |
| Lost | 44.65 | 7.52 | 25.74 | 220411.58 | 95 |

- Cluster 0: Diamond customer segment. They have spent the most in term of quantity and total amount sales, the most frequent buyers, and have purchased quite recently. **This is the best customers so far.**
- Cluster 1: Gold customer segment. They have spent the second most quantity and total amount sales, the second most frequent buyers, BUT have purchased a long time ago. They are basically near Diamond customers level but haven't purchased anything for a long time. **This customers almost got churned.**
- Cluster 2: Lost customer segment. They have spent the least in term of quantity and total amount sales. The least frequent buyers, and also haven't purchased anything long time ago. **This is most likely the churned customers**.
- Cluster 3: Potential customer segment. They have quite good spending in term of quantity and total amount sales, quite purchase frequently, and the most recent buyers. **This customers have a lot of potential to improve.**

# Modeling – Interpretation

Demographic segmentation
(KPrototypes)

| Demographic_Segmentation | Age mean | Income mean | Marital Status mean | Gender mean | count |
|---|---|---|---|---|---|
| Senior | 57.03 | 7.82 | 0.92 | 0.62 | 78 |
| Late-Adult | 48.51 | 15.86 | 0.94 | 0.30 | 99 |
| Middle-Adult | 35.16 | 6.96 | 0.83 | 0.38 | 172 |
| Young-Adult | 23.95 | 2.46 | 0.29 | 0.62 | 86 |



Radar Plot of the resulting 4 clusters

Clusters
- cluster 0
- cluster 1
- cluster 2
- cluster 3

- Cluster 0: Senior customer segment. They are elder customers having been married with medium income. Most of them are Man. Most likely they are retired worker.
- Cluster 1: Middle-Adult customer segment. Their age is within 30-40 years old, having been married with medium income. Most of them are Woman.
- Cluster 2: Young-Adult customer segment. They are most likely teenager with range 20-30 years old, which are most likely single and low income. Most of them are Man.
- Cluster 3: Late-Adult customer segment. They are in mature age within 40-55 years old, having been married with high income. Most of them are Woman.

# Dashboard

## Sales Analysis by Customer Segmentation

KALBE Nutritionals

2 types of customer segmentation were built, i.e. *Behavioral Segmentation* and *Demographic Segmentation*, using KMeans and KPrototypes algorithms respectively.

4 clusters were made in *Behavioral Segmentation*, i.e. Diamond, Gold, Lost, and Potential.

4 clusters were made in *Demographic Segmentation*, i.e. Senior, Late-Adult, Middle-Adult, and Young-Adult.

| Behavioral Segmentation | Number of Customer | Avg. Frequency | Avg. Recency | Avg. Total Revenue | Avg. Total Quantity |
|---|---|---|---|---|---|
| Diamond | 104 | 15 times | 20 days | Rp 524,9K | 58 Units |
| Gold | 137 | 11 times | 48 days | Rp 376,0K | 42 Units |
| Lost | 102 | 7 times | 52 days | Rp 217,3K | 26 Units |
| Potential | 104 | 11 times | 7 days | Rp 324,8K | 37 Units |

| Demographic Segmentation | Number of Customer | Avg. Age | Avg. Income | Avg. Gender | Avg. Marital Status |
|---|---|---|---|---|---|
| Senior | 80 | 57 years | Rp 7,8M | 63% Man | 93% Married |
| Late-Adult | 105 | 49 years | Rp 17,1M | 31% Man | 93% Married |
| Middle-Adult | 171 | 35 years | Rp 7,0M | 37% Man | 83% Married |
| Young-Adult | 91 | 23 years | Rp 2,6M | 64% Man | 32% Married |

Rp 10.630M — Rp 25.294M

© 2023 Mapbox © OpenStreetMap

### Revenue by Behavioral Segment
Rp 162,043M
- Potential 20,85%
- Diamond 33,69%
- Lost 13,68%
- Gold 31,79%

### Revenue Over Time in 2022 Period
Rp 15,446M
Rp 12,095M

### Revenue by Demographic Segment
Rp 162,043M
- Young-Adult 19,93%
- Senior 17,98%
- Middle-Adult 39,34%
- Late-Adult 22,75%

### Revenue Top 5 product
- Cheese Stick — Rp 27,615M
- Choco Bar
- Coffee Candy — Rp 19,712M
- Yoghurt — Rp 19,630M
- Oat

### Quantity Top 5 product
- Thai Tea — 2.853 Units
- Ginger Candy
- Choco Bar — 2.408 Units
- Coffee Candy — 2.097 Units
- Yoghurt

### Quantity by Behavioral Segment
18.296 Units
- Potential 21,22%
- Diamond 32,93%
- Lost 14,22%
- Gold 31,64%

### Quantity Over Time in 2022 Period
1.753 Units
1.409 Units
Jan    Apr    Jul    Oct

### Quantity by Demographic Segment
18.296 Units
- Young-Adult 19,69%
- Senior 17,67%
- Middle-Adult 39,68%
- Late-Adult 22,96%

### Revenue Top 5 customers
- ID: 365 — Rp 0,847M
- ID: 89
- ID: 392 — Rp 0,751M
- ID: 184 — Rp 0,676M
- ID: 13

### Quantity Top 5 customers
- ID: 89 — 79 Units
- ID: 365 — 75 Units
- ID: 189 — 73 Units
- ID: 245 — 72 Units
- ID: 37

### Revenue Bottom 5 customers
- ID: 208 — Rp 0,121M
- ID: 441
- ID: 439 — Rp 0,114M
- ID: 259 — Rp 0,108M
- ID: 25 — Rp 0,092M

### Quantity Bottom 5 customers
- ID: 255 — 16 Units
- ID: 262
- ID: 441 — 13 Units
- ID: 371 — 13 Units
- ID: 25 — 10 Units

Note: you can filter the dashboard by clicking the plot.

You can explore further using a dashboard created on Tableau Public.

https://public.tableau.com/app/profile/tito5892/viz/shared/CG5HZHQRK

# Strategy

## Behavioral Segmentation

Diamond customer segment. (Focus on increasing frequency and retention)
- Loyalty programs and give rewards/promo to make them feel respected,
- Market most expensive products,
- Offer new products, and cross-selling/up-selling strategy.

Gold customer segment. (Focus on maintaining their loyalty and improve their value )
- Offer them a discount, free trial, or another incentive.
- Make limited time offers.

Lost customer segment. (Focus on reactivating the customer)
- Reactivation strategy such as send them reactivation emails and ask them for feedback.
- Provide support.

Potential customer segment. (Focus on increasing their value)
- Cross-selling/up-selling strategy, give price incentives and new products recomendation.

## Demographic Segmentation

- Adjust pricing strategies based on income levels.
- Adapt products or services to satisfy to the needs and preferences of different demographic segments.

# Dataset

The transaction dataset (5020 rows) are grouped by day and the quantity are summed for each group resulting daily total quantity of all products sold in 2022 within 365 days (365 rows).

Our time-series is already stationary

| Date | Total_Quantity |
|------|----------------|
| 2022-01-01 | 49 |
| 2022-01-02 | 50 |
| 2022-01-03 | 76 |
| 2022-01-04 | 98 |
| 2022-01-05 | 67 |
| ... | ... |
| 2022-12-27 | 70 |
| 2022-12-28 | 68 |
| 2022-12-29 | 42 |
| 2022-12-30 | 44 |
| 2022-12-31 | 37 |



```
adfuller test:
p-values: 0.0
Reject H0 ---> time series is stationary

KPSS test:
p-values: 0.06622747821677664
Accept H0 ---> time series is stationary
```

# ACF and PACF

### Original Time Series



Using the original data, There is no significant autocorrelation. This is called "**White Noise**". In this case our time-series is stationary, yet zero autocorrelations at all lags are found.
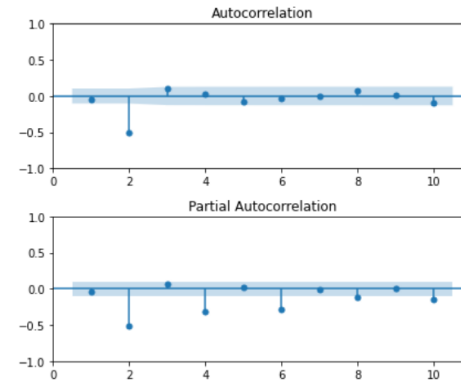
**d = 1**



adfuller test:
p-values: 6.148467265279569e-21
Reject H0 ---> time series is stationary

KPSS test:
p-values: 0.1
Accept H0 ---> time series is stationary

**d = 2**



adfuller test:
p-values: 3.4821335991476044e-11
Reject H0 ---> time series is stationary

KPSS test:
p-values: 0.1
Accept H0 ---> time series is stationary

- With **d** = 1 or **d** = 2, our time-series now have autocorrelation greater than zero.
- From ACF and PACF, It suggest that the candidates for parameters **p, q,** and **d** as follows:
1. **d** can be 1 or 2,
2. **p** can be 1 to 6,
3. **q** can be 1 to 2

# **Build ARIMA Model**

```
Performing stepwise search to minimize aic
 ARIMA(2,2,2)(0,0,0)[0]             : AIC=inf, Time=0.55 sec
 ARIMA(0,2,0)(0,0,0)[0]             : AIC=3731.707, Time=0.02 sec
 ARIMA(1,2,0)(0,0,0)[0]             : AIC=3518.903, Time=0.03 sec
 ARIMA(0,2,1)(0,0,0)[0]             : AIC=inf, Time=0.09 sec
 ARIMA(2,2,0)(0,0,0)[0]             : AIC=3401.324, Time=0.05 sec
 ARIMA(3,2,0)(0,0,0)[0]             : AIC=3319.933, Time=0.07 sec
 ARIMA(4,2,0)(0,0,0)[0]             : AIC=3287.880, Time=0.11 sec
 ARIMA(5,2,0)(0,0,0)[0]             : AIC=3259.857, Time=0.10 sec
 ARIMA(6,2,0)(0,0,0)[0]             : AIC=3249.721, Time=0.19 sec
 ARIMA(6,2,1)(0,0,0)[0]             : AIC=inf, Time=0.82 sec
 ARIMA(5,2,1)(0,0,0)[0]             : AIC=inf, Time=0.78 sec
 ARIMA(6,2,0)(0,0,0)[0] intercept   : AIC=3251.685, Time=0.43 sec

Best model:  ARIMA(6,2,0)(0,0,0)[0]
Total fit time: 3.226 seconds
```
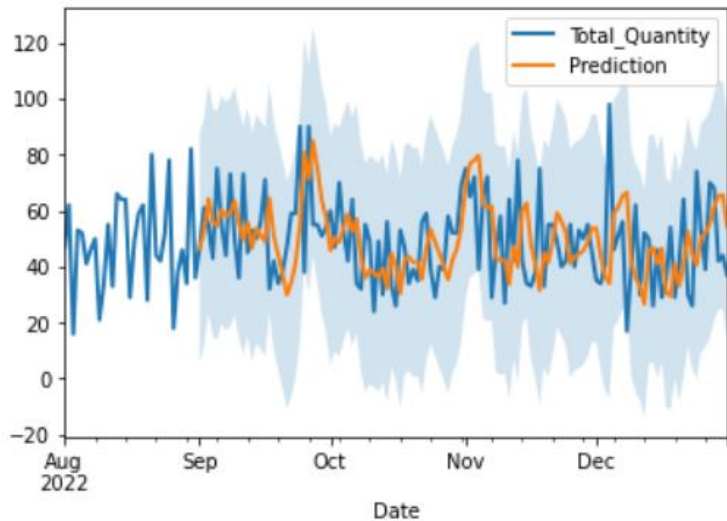
```python
In [50]:  # build ARIMA model
          from statsmodels.tsa.arima.model import ARIMA
          model = ARIMA(df_daily, order = (6,2,0))
          result = model.fit()
```

Using pm.auto_arima library we should use ARIMA(6, 2, 0) in
which gives the optimize evaluation metrics

# Model Evaluation
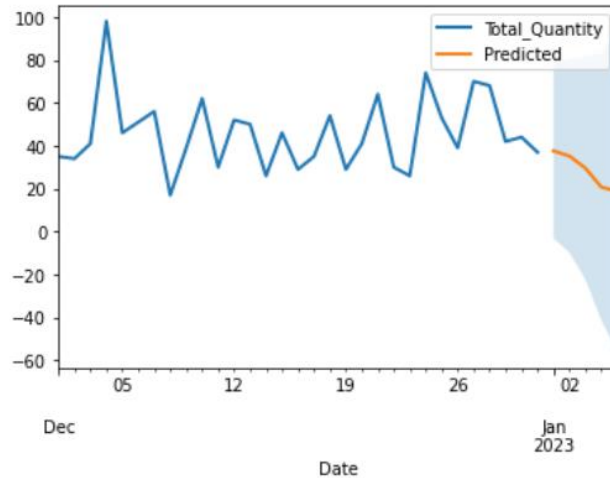


```
In [79]: # calculate performance metric
         from sklearn.metrics import mean_squared_error, mean_absolute_error
         y_true = df_daily.values
         y_pred = result.get_prediction(start = df_daily.index[0], dynamic = False).predicted_mean.values

         MAE = mean_absolute_error(y_true, y_pred)
         RMSE = mean_squared_error(y_true, y_pred, squared = False)
         print('MAE: {:.2f} units'.format(MAE))
         print('RMSE: {:.2f} units'.format(RMSE))

         MAE: 16.84 units
         RMSE: 21.01 units
```

The model is far from good, but it can follow general upward and downward trends

# Forecast The Future



| | Predicted |
|---|---|
| **2023-01-01** | 37.606836 |
| **2023-01-02** | 35.212744 |
| **2023-01-03** | 29.564290 |
| **2023-01-04** | 20.710638 |
| **2023-01-05** | 18.881172 |

The model said that the total quantity of products will be declining
5 days later from the end of December

https://github.com/dstito/Customer-Segmentation-and-Forecasting-Sales-Quantity

# Thank You