

Unlocking Insight : A Comprehensive Customer Segmentation Study

An Empirical Study of Enhancing Decision Making through Segmentation Findings

Dharma Setiawan

伍铠佑

Josep Ariel Christopher Teja

黃熙

Darren Boesono

吴凡



WORKFLOW ML

Stage 1



Introduction

Stage 2



Data Collection
& EDA

Stage 3



Modeling &
Results

Stage 4



Implications

Stage 5



Recommendation
& Conclusion

Background

Change or Die?



Customer Segmentation

What is customer segmentation?

- Practice of dividing customers into different groups based on different marketing methods and specific perspectives.
- Used to predict customer behavior patterns. Once identified, it is possible to make predictions about the groups' responses to different situations, to align marketing strategies and types of policy, and to allow more creative and better targeted policies to emerge.
- Most commonly used attributes: location, age, gender, income, lifestyle, and previous purchasing behavior.



Background

Why customer segmentation is important?

Pareto principle: 20% of customers contribute more to the company's revenue than the rest. (Srivastava, 2016)

Delivering more targeted information

Marketing personnel can directly target different customer groups, delivering products and marketing information that better suit their interests and needs to resonate more.

Reasonably allocate marketing resources

After segmenting customers, brands can effectively allocate marketing resources to specific customer groups and maximize cross selling and upward sales opportunities.

Maintain better customer relationships

Help marketers better understand their customers, providing personalized customer service experiences to improve customer loyalty and retention.

Adjusting marketing budget

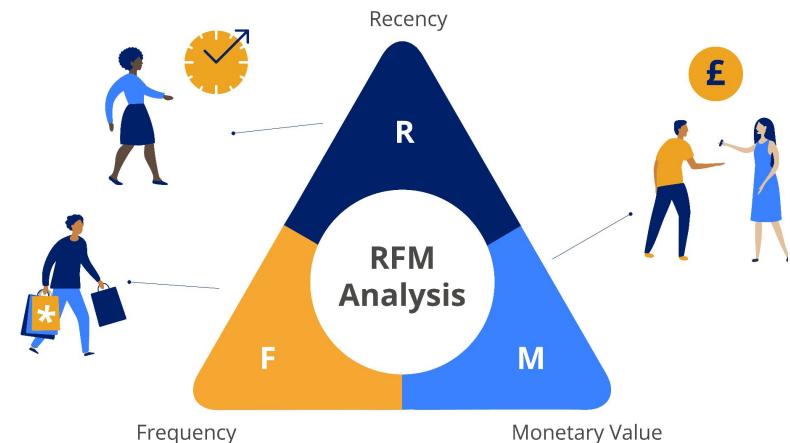
By analyzing its revenue and costs, the profit potential of each segmented market can be determined, helping the company determine the profit level of the segmented market, so that they can adjust their marketing budget accordingly.

RFM Analysis

RFM is a strategic approach employed in the analysis and estimation of a customer's worth, predicated upon the evaluation of three crucial data points, namely: **Recency, Frequency, and Monetary Value.**

Limitations of RFM

- It does not take into account key factors such as customer demographics or purchased items.
- More importantly, it depends entirely on the customer's history data, which means it may not be able to accurately predict the future activities of customers.
- Therefore, it is necessary to combine more advanced prediction methods.



Research Questions

Questions

Description



Question #1

Which is the best approach between RFM analysis, Clustering combined with classification, and a standalone classification model based on customer response to previous campaigns, in terms of their abilities to segment customers and provide actionable business insights for optimizing marketing strategies?



Question #2

Given our clusters, what segments do we get? How much the potential Gross Merchandise Value (GMV) and optimization cost? What is the personalized approach marketing strategy for each segment?

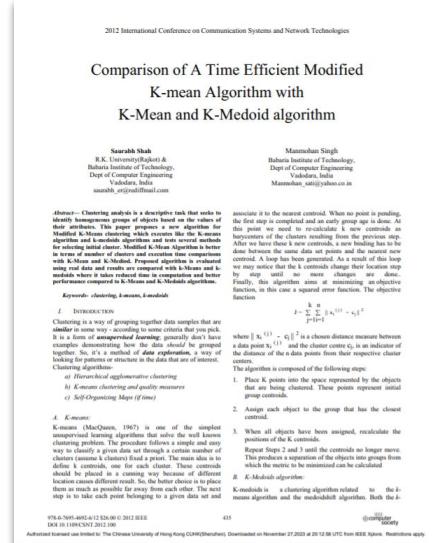
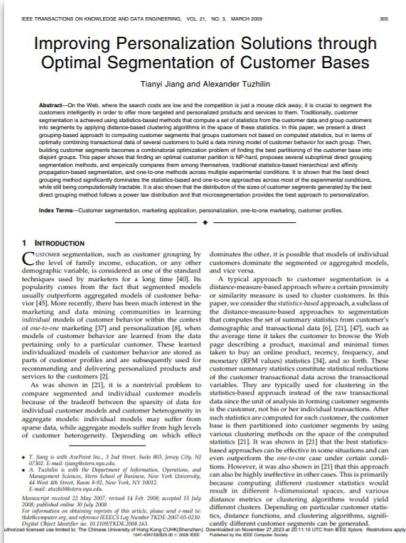


Question #3

How might a company implement a customer segmentation service in their internal stack?

Literature Review

Jiang and Tuzhilin (2009) proposed the K-Classifiers Segmentation algorithm, and found that this approach yielded far superior results compared to the traditional statistical approach. Shah and Singh (2012) proposed a new clustering algorithm which executes similar to the K-means algorithm and K-medoids algorithms. But the proposed algorithm does not provide an optimal solution in all cases. Zahrotun (2017) used the customer data from online to identify the finest customer using Customer Relationship Management (CRM).



WORKFLOW ML

Stage 1



Introduction

Stage 2



Data Collection
& EDA

Stage 3



Modeling &
Results

Stage 4



Implications

Stage 5



Recommendation
& Conclusion

Data Collection

<https://www.kaggle.com/datasets/rodsaldanha/marketing-campaign/data>

2240 Rows

29 Features

0 Duplicate

24 Missing Values

20 Numerical

1 Target

2 Categorical

Year_Birth

Recency

Teen home

Mnt Wines

Num Deals Purchases

...

Response

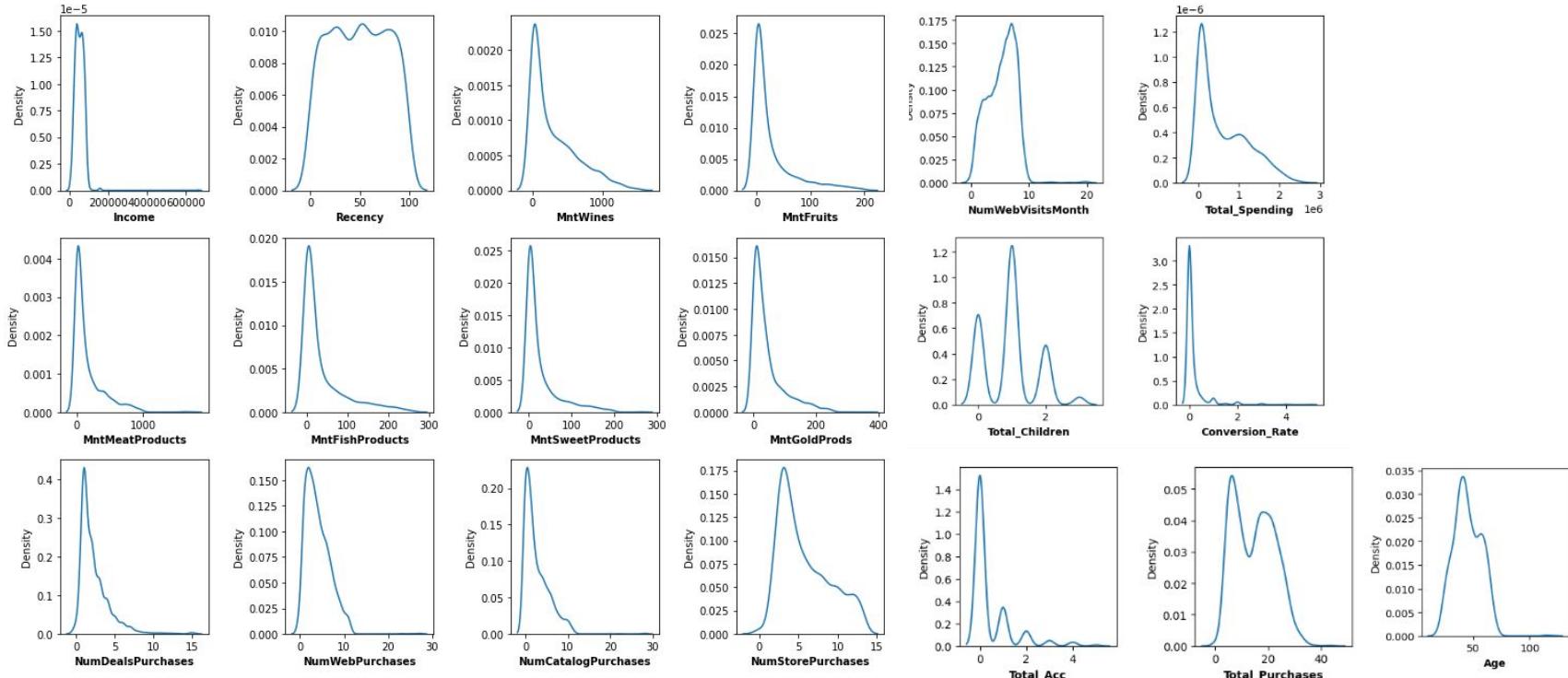
Educa
tion

Marital_Status

EDA in General

Univariate Analysis: Distplot Analysis

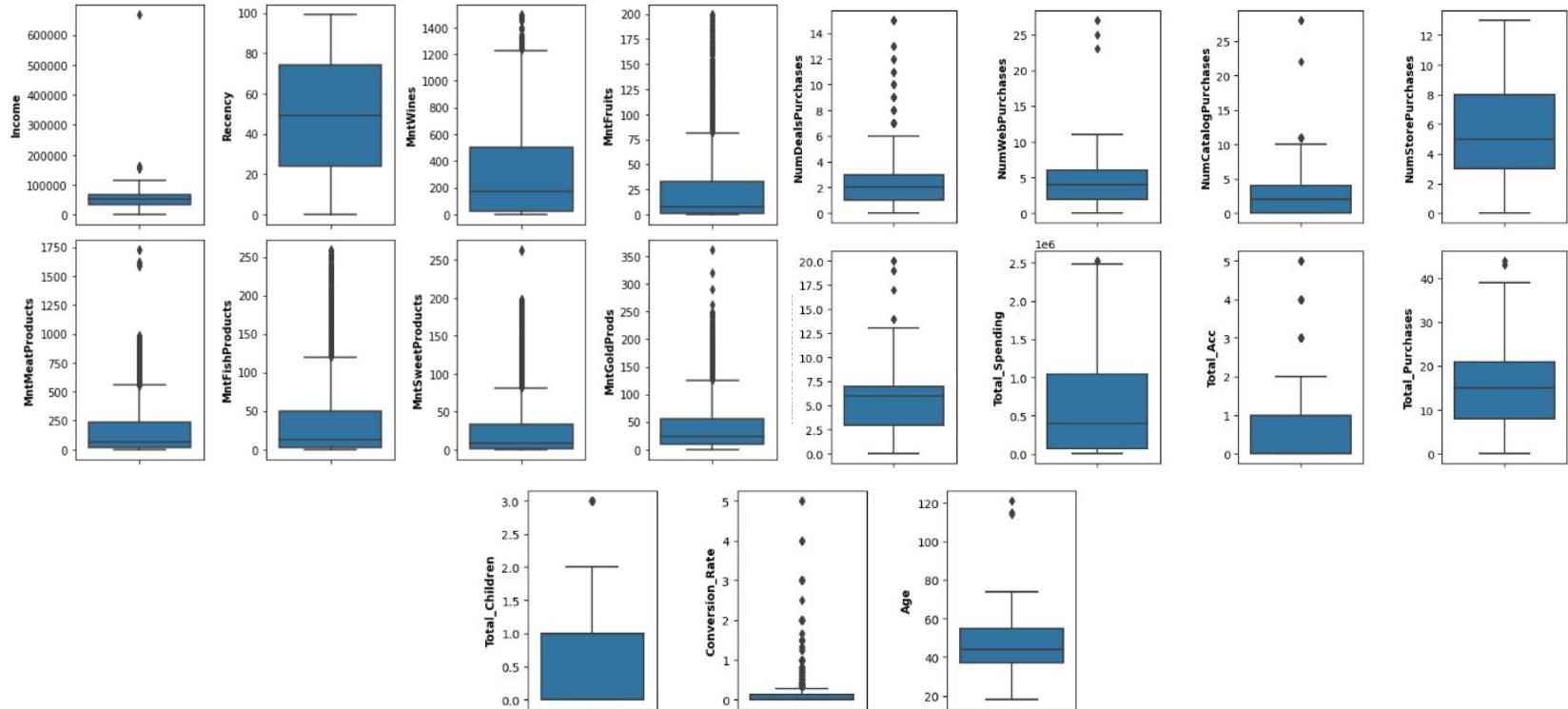
Most of the dataset has highly positively skewed distribution



EDA in General

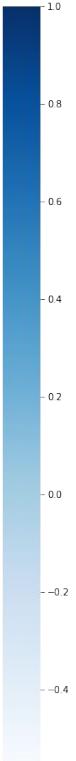
Univariate Analysis: Box Plot Analysis

There are some outliers in each features and most features are not symmetric



EDA in General

Income	1.00	-0.00	0.58	0.43	0.58	0.44	0.44	0.33	-0.08	0.39	0.59	0.53	-0.55	0.67	0.29	0.57	-0.29	0.33	0.16
Recency	-0.00	1.00	0.02	-0.00	0.02	0.00	0.02	0.02	-0.00	-0.01	0.03	0.00	-0.02	0.02	-0.09	0.01	0.02	-0.05	0.02
MntWines	0.58	0.02	1.00	0.39	0.56	0.40	0.39	0.39	0.01	0.54	0.64	0.64	-0.32	0.89	0.49	0.71	-0.35	0.43	0.16
MntFruits	0.43	-0.00	0.39	1.00	0.54	0.59	0.57	0.39	-0.13	0.30	0.49	0.46	-0.42	0.61	0.17	0.46	-0.39	0.26	0.02
MntMeatProducts	0.58	0.02	0.56	0.54	1.00	0.57	0.52	0.35	-0.12	0.29	0.72	0.48	-0.54	0.84	0.33	0.55	-0.50	0.41	0.03
MntFishProducts	0.44	0.00	0.40	0.59	0.57	1.00	0.58	0.42	-0.14	0.29	0.53	0.46	-0.45	0.64	0.18	0.47	-0.43	0.26	0.04
MntSweetProducts	0.44	0.02	0.39	0.57	0.52	0.58	1.00	0.37	-0.12	0.35	0.49	0.45	-0.42	0.60	0.20	0.47	-0.38	0.28	0.02
MntGoldProds	0.33	0.02	0.39	0.39	0.35	0.42	0.37	1.00	0.05	0.42	0.44	0.38	-0.25	0.52	0.20	0.49	-0.27	0.22	0.06
NumDealsPurchases	-0.08	-0.00	0.01	-0.13	-0.12	-0.14	-0.12	0.05	1.00	0.23	-0.01	0.07	0.35	-0.07	-0.09	0.36	0.44	-0.17	0.06
NumWebPurchases	0.39	-0.01	0.54	0.30	0.29	0.29	0.35	0.42	0.23	1.00	0.38	0.50	-0.06	0.52	0.21	0.78	-0.15	0.09	0.15
NumCatalogPurchases	0.59	0.03	0.64	0.49	0.72	0.53	0.49	0.44	-0.01	0.38	1.00	0.52	-0.52	0.78	0.35	0.74	-0.44	0.36	0.12
NumStorePurchases	0.53	0.00	0.64	0.46	0.48	0.46	0.45	0.38	0.07	0.50	0.52	1.00	-0.43	0.67	0.17	0.82	-0.32	0.20	0.13
NumWebVisitsMonth	-0.55	-0.02	-0.32	-0.42	-0.54	-0.45	-0.42	-0.25	0.35	-0.06	-0.52	-0.43	1.00	-0.50	-0.13	-0.31	0.42	-0.34	-0.12
Total_Spending	0.67	0.02	0.89	0.61	0.84	0.64	0.60	0.52	-0.07	0.52	0.78	0.67	-0.50	1.00	0.46	0.75	-0.50	0.47	0.11
Total_Acc	0.29	-0.09	0.49	0.17	0.33	0.18	0.20	0.20	-0.09	0.21	0.35	0.17	-0.13	0.46	1.00	0.26	-0.25	0.77	-0.01
Total_Purchases	0.57	0.01	0.71	0.46	0.55	0.47	0.47	0.49	0.36	0.78	0.74	0.82	-0.31	0.75	0.26	1.00	-0.25	0.21	0.17
Total_Children	-0.29	0.02	-0.35	-0.39	-0.50	-0.43	-0.38	-0.27	0.44	-0.15	-0.44	-0.32	0.42	-0.50	-0.25	-0.25	1.00	-0.31	0.09
Conversion_Rate	0.33	-0.05	0.43	0.26	0.41	0.26	0.28	0.22	-0.17	0.09	0.36	0.20	-0.34	0.47	0.77	0.21	-0.31	1.00	-0.02
Age	0.16	0.02	0.16	0.02	0.03	0.04	0.02	0.06	0.06	0.15	0.12	0.13	-0.12	0.11	-0.01	0.17	0.09	-0.02	1.00



Multivariate Analysis: Correlation Heatmap

- Customers are more likely to buy wine and meat product
- Customers tend to buy products directly in stores
- Number of visits to company's web site and number of children in customers' household contribute little to the response

Feature Engineering

Method

**customer_year
&
customer_month**

Description

```
#Convert date to date object
df["Dt_Customer"] = pd.to_datetime(df["Dt_Customer"])

#Create Month and Year variables from Dt_Customer
df["Customer_year"] = df["Dt_Customer"].dt.year
df["Customer_month"] = df["Dt_Customer"].dt.month

#Removing ID, Date from further analysis
df.drop(["ID"], axis=1, inplace=True)
df.head()
```

**age &
join_at_age**

```
# age_range
df['age'] = 2023 - df['Year_Birth']
df['join_at_age'] = df['Dt_Customer'].dt.year - df['Year_Birth']
df.loc[(df['age'] >= 0) & (df['age'] < 12), 'age_range'] = "child"
df.loc[(df['age'] >= 12) & (df['age'] < 18), 'age_range'] = "teens"
df.loc[(df['age'] >= 18) & (df['age'] < 36), 'age_range'] = "young_adults"
df.loc[(df['age'] >= 36) & (df['age'] < 55), 'age_range'] = "middle_aged_adults"
df.loc[(df['age'] >= 55), 'age_range'] = "older_adults"
```

**total_children &
is_parents**

```
df['total_children'] = df['Kidhome'] + df['Teenhome']
df['is_parents'] = np.where(df['total_children'] > 0, 1, 0)
```

Data Preprocessing Feature Transformation

Method	Description	Example
Normalization	<ul style="list-style-type: none">Hard feature scaling: change the range of the feature with the valid or exact rangeDoesn't change the distribution of the data	<pre>scaler= MinMaxScaler() scaled = scaler.fit_transform(df_enc) df_enc_scaled= pd.DataFrame(scaled, columns= df_enc.columns) df_enc_scaled.head()</pre>
Standardization	<ul style="list-style-type: none">Soft feature scaling: change the range of the features with not exact rangeChange the distribution of the data to the normal distributionTransformation will result in mean = 0 with sd = 1	<pre>from sklearn.preprocessing import StandardScaler df_standard = df.copy() for i in numerical_cols: df_standard[i] = StandardScaler().fit_transform(df_standard[i].values.reshape(len(df_standard), 1)) df_standard</pre>

Data Preprocessing Feature Encoding

Method	Description	Example
Label Encoding	Label Encoding is a categorical feature change to numeric by giving a different number for each value is unique	<pre>#label encoder mapping_education = { 'Basic' : 0, 'Graduation' : 1, '2n Cycle': 2, "Master":2, "PhD":3 } df['Education_mapped'] = df['Education'].map(mapping_education)</pre> <pre>le= LabelEncoder() enc_cols= ["Marital_Status", 'age_range', 'is_parents'] for c in enc_cols: df_enc[c]= le.fit_transform(df_enc[c]) df_enc.drop(["Education"], axis= 1, inplace=True)</pre>

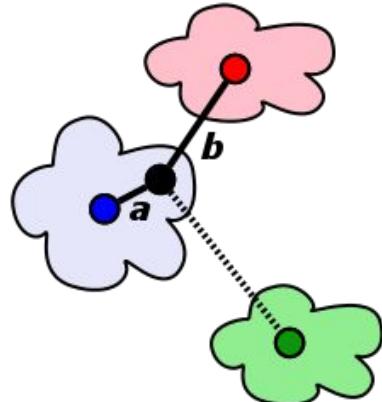
Model Evaluation (Silhouette Score)

Silhouette coefficient or silhouette score is a metric used to calculate the goodness of a clustering technique. Its value ranges from -1 to 1.

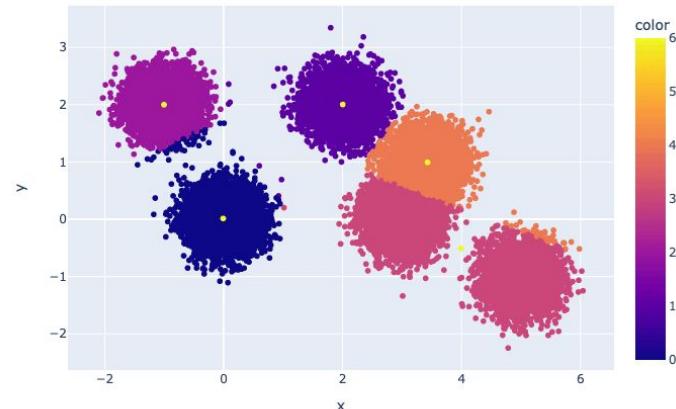
1: Means clusters are well apart from each other and clearly distinguished.

0: Means clusters are indifferent, or we can say that the distance between clusters is not significant.

-1: Means clusters are assigned in the wrong way.



$$SSI_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$



WORKFLOW ML

Stage 1



Introduction

Stage 2



Data Collection
& EDA

Stage 3



Modeling &
Results

Stage 4



Implications

Stage 5



Recommendation
& Conclusion

RFM Architecture

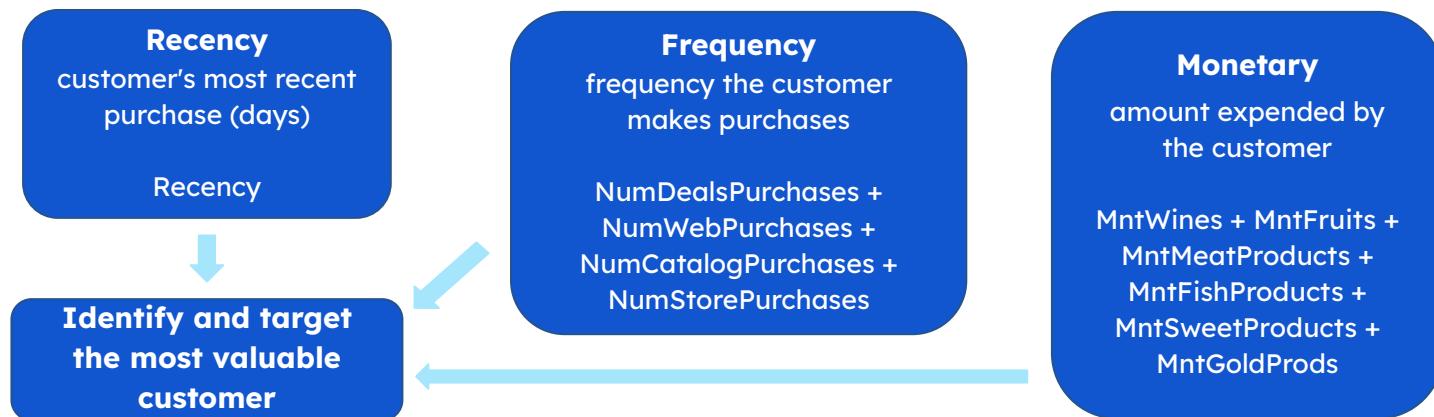
Motivation

RFM is a straightforward and highly effective way to comprehend and assess consumer behavior predicated on purchase.

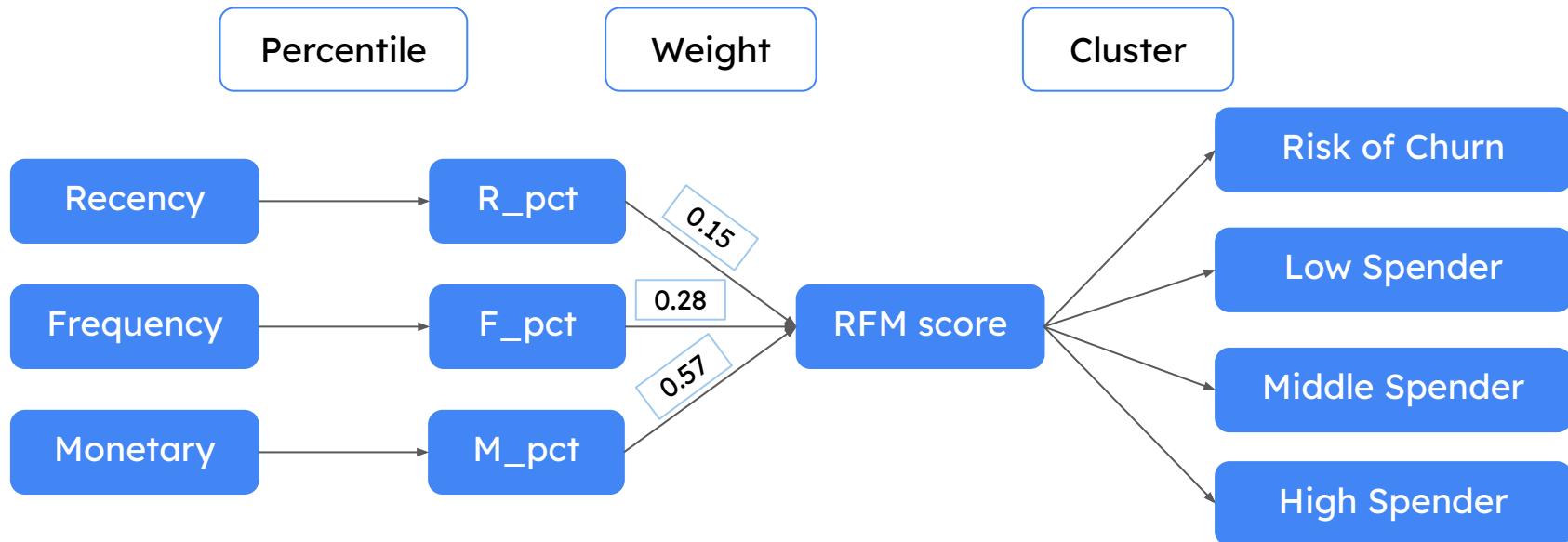
Idea

It functions by quantitatively categorizing and classifying patrons based on the recency, frequency, and monetary total of their most recent transactions, with the end goal of identifying and targeting the most valuable customers for the purposes of performing focused and precision-targeted marketing campaigns.

Implementation



RFM Scores

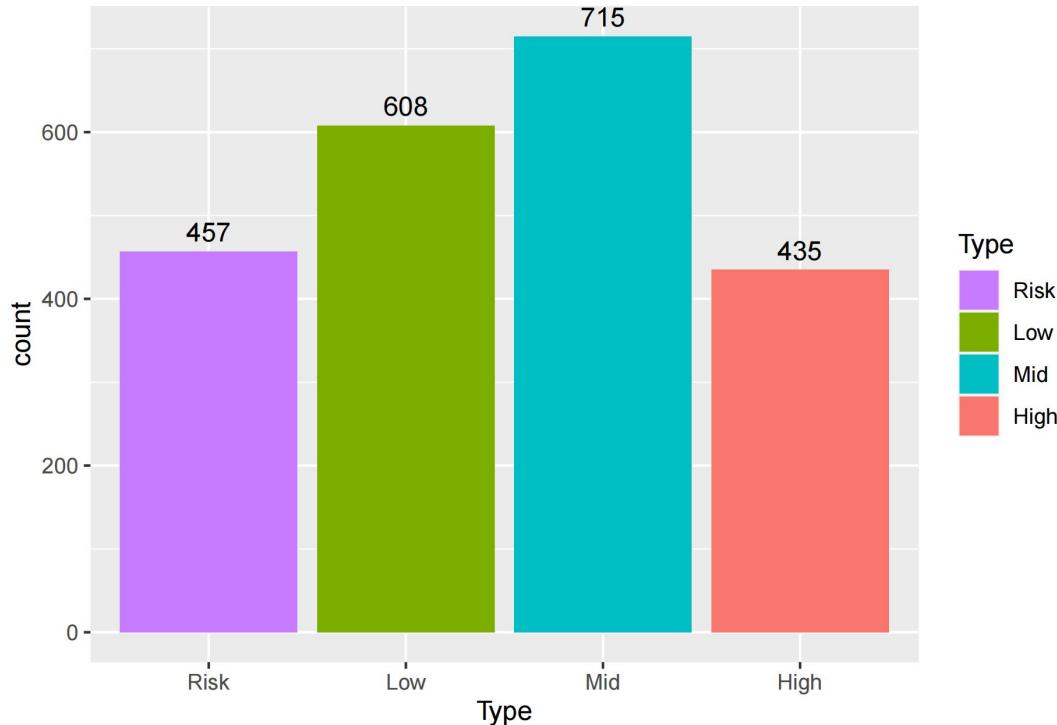


For clustering convenience, we assume $\text{RFM score} = \text{RFM score} * 4$

Cluster criteria:

- $\text{RFM score} \leq 1$, Risk of Churn
- $1 < \text{RFM score} \leq 2$, Low Spender
- $2 < \text{RFM score} \leq 3$, Middle Spender
- $3 < \text{RFM score}$, High Spender

RFM Results



After calculating the RFM score and clustering, we have

457 risk of churn

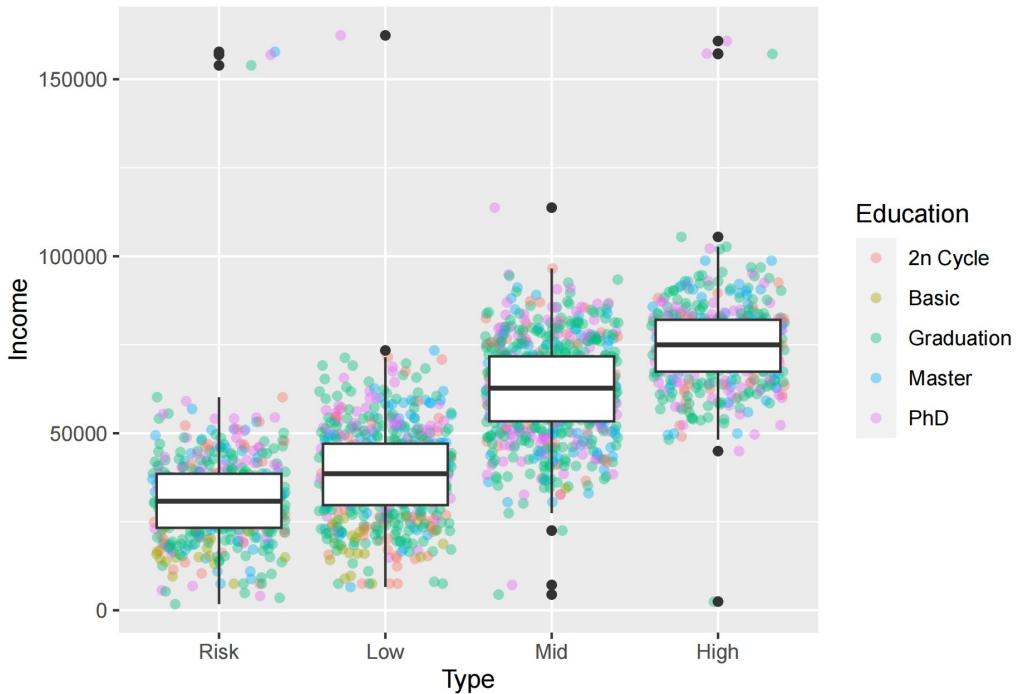
608 low spenders

715 middle spenders

435 high spenders

Silhouette score: 0.528

RFM Cluster Interpretation

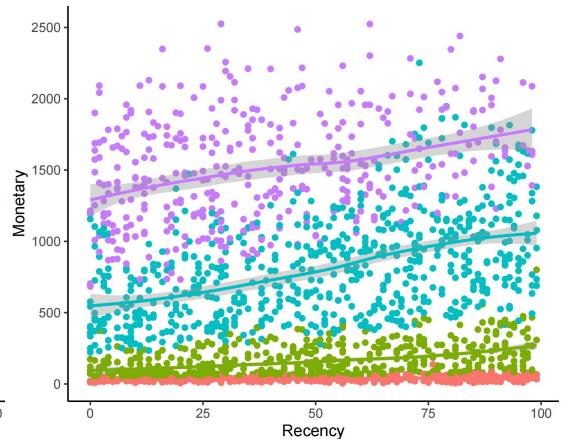
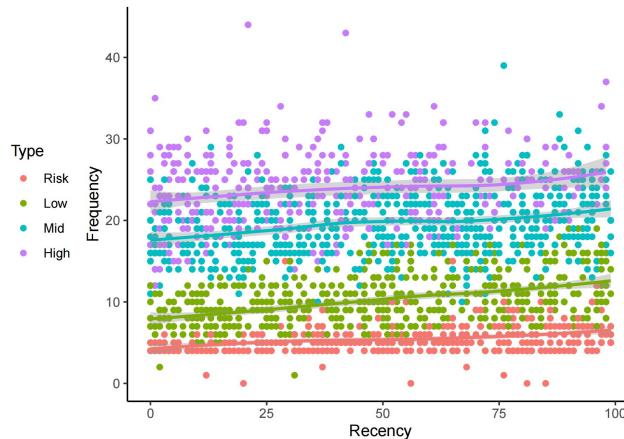
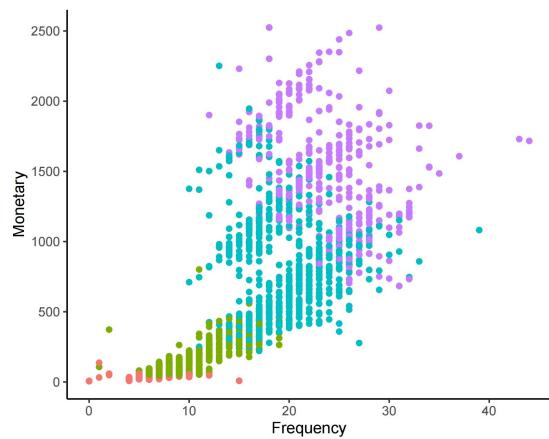


Education

- 2n Cycle
- Basic
- Graduation
- Master
- PhD

- High spenders typically have higher income, while low spenders usually have lower income.
- Type of customer is education independent.

RFM Cluster Interpretation



Frequency

Corr: 0.007
 High: **0.167*****
 Low: **0.437*****
 Mid: **0.250*****
 Risk: **0.298*****

Monetary

Corr: 0.020
 High: **0.334*****
 Low: **0.435*****
 Mid: **0.472*****
 Risk: **0.404*****

Recency

Corr: 0.756***
 High: **-0.359*****
 Low: **0.782*****
 Mid: **0.044**
 Risk: **0.452*****

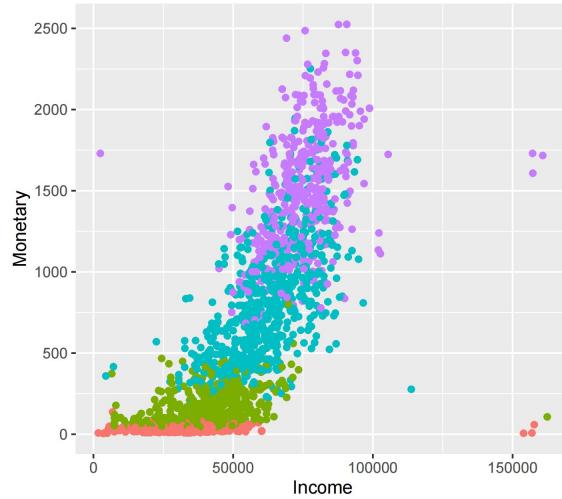
Frequency

To avoid Simpson's Paradox, we stratify the data with customer segments

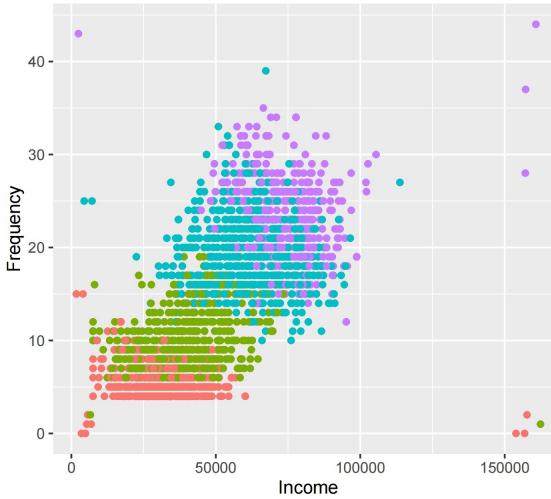
- High spenders typically have higher frequency and monetary, but recency can be varied.
- Frequency and recency; monetary and recency are significantly positively correlated within customer segments.
- Although monetary and frequency seem to be positively correlated, they are significantly negatively correlated within high spender segment.

*** p < 0.001 ** p < 0.01 * p < 0.05 . p < 0.10 . otherwise

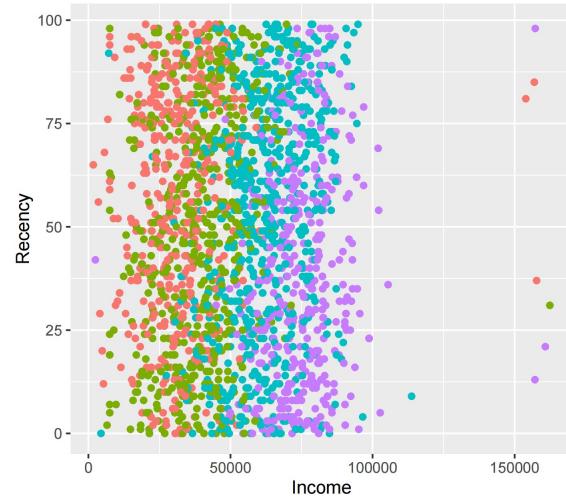
RFM Customer Interpretation



- Customers with higher income may not necessarily spend a lot on purchases
- But low monetary customers generally have lower incomes



- Generally, customers with higher income purchase products more frequently



- The correlation between recency and income is not significant
- High spenders own higher income

Clustering + Classification Architecture

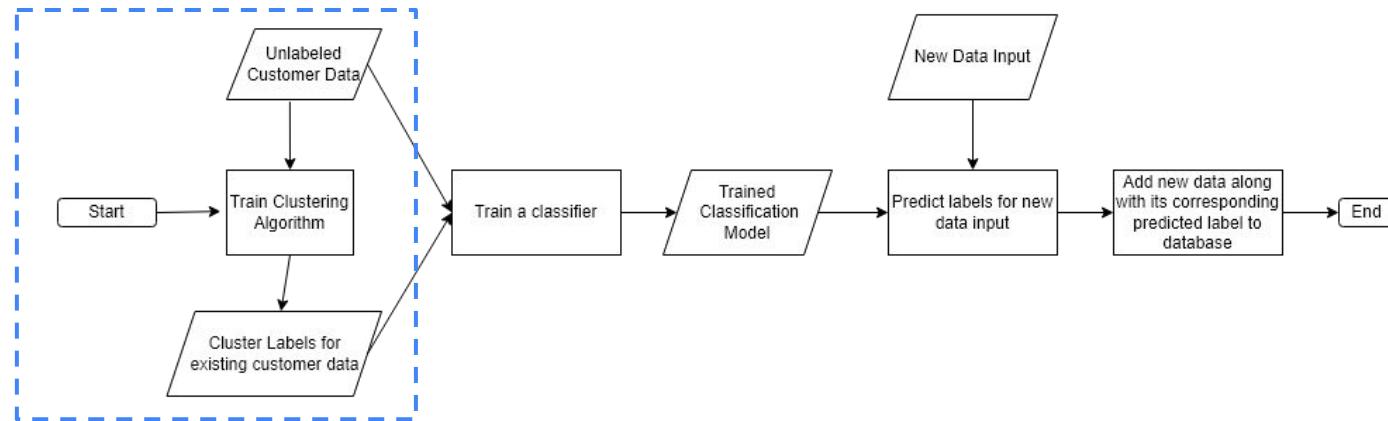
Motivation

Although clustering can be done to segment existing customers, not all clustering algorithms are able to efficiently cluster new data points (future customer data).

Idea

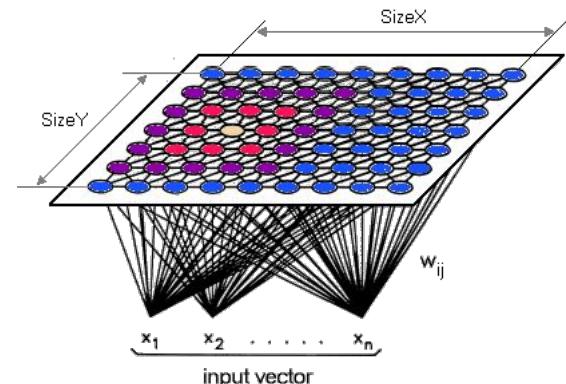
Train a robust classifier using the clustered customer data to efficiently segment new customers.

Implementation



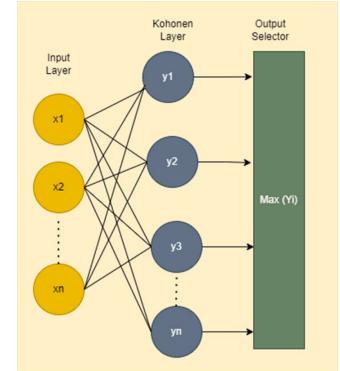
Self-Organizing Map (SOM), (Kohonen)

- It is an **unsupervised neural network** that is trained using unsupervised learning techniques to produce a low dimensional, discretized representation from the input space of the training samples, known as a map and is a method to reduce data dimensions.
- Unlike other neural networks, SOMs apply **competitive learning technique** instead of error-correction learning methods like backpropagation with gradient descent.
- Unlike other clustering techniques, SOM uses a neighbourhood function to **preserve the topological properties** within the input space
- Intuitively, you can think of a SOM as something like k-means where the k-centroids are constrained in a 2D manifold.



SOM as a Neural Network

- There are two layers: 1. **Input Layer** and 2. **Kohonen (Competitive) Layer** which acts as an output node
- The **Input Layer** is responsible for receiving input data and transmitting it to the second layer.
- The **Kohonen Layer** acts as the output layer where each neuron competes to become the winner.
 - Neurons in this layer are interconnected, forming a competitive network by restricting feedback.
 - Neurons have **excitatory links to immediate neighbors**, fostering local competition, and **inhibitory connections to distant neurons**, preventing excessive influence.
 - The layer uses a "**winner takes all**" approach, where the most responsive neuron to a given input set becomes the sole winner, generating a binary output (1 for the winner, 0 for others).
- Unlike feed forward and back propagation neural networks, SOM does not use threshold or activation function, nor does it have any hidden layers.



SOM: Training Process

- Parameters: Map size, Neighbourhood size, learning rate, iterations
- SOMs are trained with neighbourhood size, update parameters and weight.

Training Steps:

1. **Network Initialization:** Initialize $w_{ij}(t)$ ($0 \leq i \leq n - 1$) is the weight from i^{th} input node to j^{th} node at iteration t . the n input nodes are assigned with the weights randomly with smaller value. The neighborhood radius around the node j is $N_j(0)$.

2. **New Input Initialization:** input $X_0(t), X_1(t), X_2(t), \dots, X_{n-1}(t)$ where $X_i(t)$ is the input from node to node i at iteration t .

3. **Distance calculation:** the distance d_j between input and j^{th} output node is computed using the equation

$$d_j = \sum_{i=1}^{n-1} (X_i(t) - w_{ij}(t))^2$$

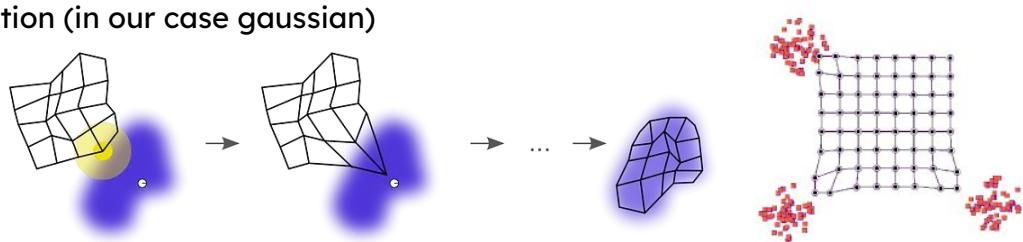
Minimum distance is selected as the winning node with the comparison and it is selected as the output node j' .

4. **Weight update:** the output node j' weight and its neighborhood node weights are updated using the following equation

$$w_{ij}(t+1) = w_{ij}(t) + \eta(t) \cdot h_{ij}(t) [X_i(t) - w_{ij}(t)]$$

This process is repeated for j in $N_{j'}(t)$ ($0 \leq i \leq n - 1$), where $N_{j'}(t)$ is neighborhood size of output neuron j' , η is learning rate in the range $[0,1]$ which decreases over the time and the N value also decreases over time with the localization of maximum activity. h_{ij} is the neighbourhood function (in our case gaussian)

5. Until maximum iteration repeat step 2-4



SOM: Methods & Evaluation Metrics

Treatments

- There are two ways in which we used SOM for our case
 1. Directly using it as a **clustering algorithm**.
 2. Uses it as a **dimensionality reduction technique** followed by the hierarchical agglomerative clustering algorithm for clustering into 4 clusters.
- Due to unsatisfactory initial results, we explored the impact of reducing feature dimensionality on algorithm performance. We repeated both methods, limiting features to 4 which are:
 - Income, Recency, Total transaction, Total amount spent.

Parameters

- Neighbourhood Size = 0.5 ; Learning Rate = 0.5 ; Neighbourhood Function = Gaussian ; Maximum number of iterations = 1000
- Map Size: Treatment 1 = (4,1) ; Treatment 2 = (20,20)

Evaluation Metrics

- We used three evaluation metrics to measure model performance
 1. **Quantization error:** Measures the distance from each data vector to its best matching unit.
 2. **Topographic error:** Measures how well the topographic structure of the data is preserved on the map.
 3. **Silhouette score:** Used to get a reference for comparison with the other algorithms

SOM: Results

Treatment	Quantization Error	Topographic Error	Silhouette Score
SOM + No Feature Trimming	1.084	0.69	0.16
SOM + Agglomerative + No Feature Trimming	0.597	0.97	0.143
SOM + Feature Trimming	0.229	0.84	0.387
SOM + Agglomerative + Feature Trimming	0.038	0.98	0.372

SOM: Remarks

1. Contrary to literature, SOM seems to perform quite badly on this particular dataset.
2. The poor performance of SOM in our case is likely due to the large number of categorical variables and outliers in our dataset.
3. Feature reduction increases silhouette score and reduces quantization error.
4. Having a larger map size reduces quantization error but increases topographic error.

Modeling

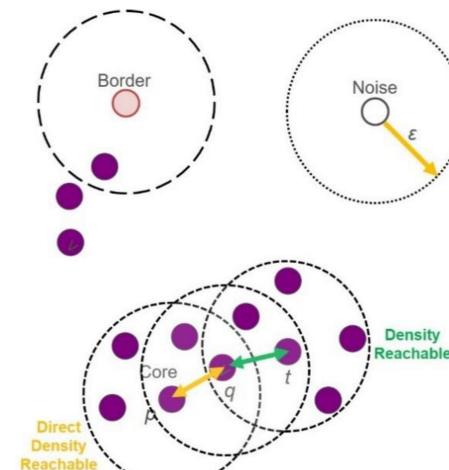
Model	Description	Silhouette Score
K-Means	K-Means is a centroid-based clustering algorithm that partitions data into K clusters, assigning each data point to the cluster with the nearest centroid.	0.1666
K-Means PCA	K-Means with PCA helps to handle high-dimensional data and improve clustering performance.	0.5786
K-Means t-SNE	K-Means with t-SNE incorporates t-SNE dimensionality reduction to visualize and cluster data points in a lower-dimensional space, preserving local structures.	0.6723
Mini Batch K-Means	Mini Batch K-Means is a variant of K-Means that processes random subsets (mini-batches) of the data at each iteration, making it computationally more efficient for large datasets.	0.6702
Agglomerative Clustering	Agglomerative Clustering is a hierarchical clustering algorithm that starts with individual data points and progressively merges them into clusters based on similarity, forming a tree-like structure.	0.6681
DBSCAN	DBSCAN is a density-based clustering algorithm that identifies clusters based on the density of data points, discovering clusters of varying shapes and handling noise effectively.	0.6651
OPTICS	OPTICS (Ordering Points To Identify the Clustering Structure) is a density-based clustering algorithm that extends DBSCAN by providing a hierarchical view of clusters and identifying varying densities within the data.	0.7199
BIRCH	BIRCH is a hierarchical clustering algorithm designed for large datasets, using a tree structure to efficiently group data points and reduce memory requirements.	0.6714

DBSCAN

The DBSCAN algorithm is based on this intuitive notion of “clusters” and “noise”. The key idea is that for each point of a cluster, the neighborhood of a given radius has to contain at least a minimum number of points.

Algorithm:

1. Find all the neighbor points within **eps** (defines the neighborhood around a data point) and identify the core points or visited with more than **MinPts** (minimum number of neighbors within eps radius) neighbors.
2. For each core point if it is not already assigned to a cluster, create a new cluster.
3. Find recursively all its **density-connected points** and assign them to the same cluster as the core point.
4. A point a and b are said to be density connected if there exists a point c which has a sufficient number of points in its neighbors and both points a and b are within the eps distance. This is a **chaining process**. So, if b is a neighbor of c, c is a neighbor of d, and d is a neighbor of e, which in turn is neighbor of a implying that b is a neighbor of a.
5. Iterate through the remaining unvisited points in the dataset. Those points that do not belong to any cluster are **noise**.



DBSCAN vs K-Means

DBSCAN

DBSCAN



In DBSCAN we do not need to specify the number of clusters

Clusters formed in DBSCAN can be of any arbitrary shape

DBSCAN can work well with datasets having noise and outliers

In DBSCAN two parameters are required for training the model

K-Means

k-means



K-Means is very sensitive to the number of clusters so it needs to be specified

Clusters formed in K-Means are spherical or convex in shape

K-Means does not work well with outliers data. Outliers can skew the clusters in K-Means to a very large extent

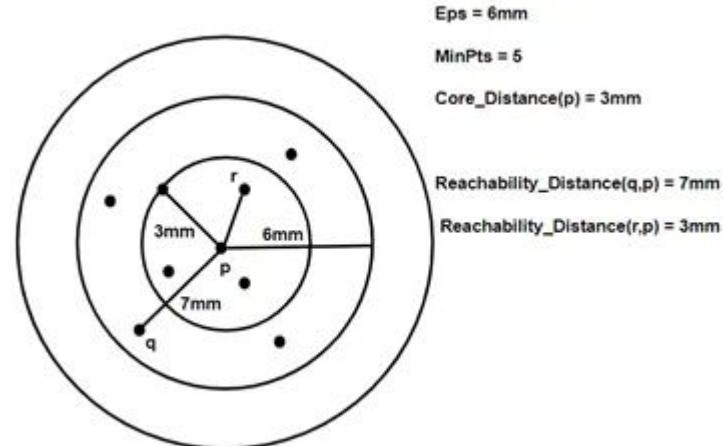
In K-Means only one parameter is required for training the model

OPTICS

Ordering Points To Identify Cluster Structure (OPTICS) is a density-based clustering technique that allows partitioning data into groups with similar characteristics (clusters). It addresses one of the DBSCAN's major weaknesses. The problem of detecting meaningful clusters in data of varying density. In a density based clustering, clusters are defined as dense regions of data points separated by low-density regions.

OPTICS adds two more terms to the concepts of DBSCAN clustering:

1. **Core Distance:** it is the minimum value of radius required to classify a given point as a core point. if the given point is not a Core point, then its Core Distance is undefined.
2. **Reachability Distance:** it is defined with respect to another data point 'q'. The Reachability distance between a point p and q

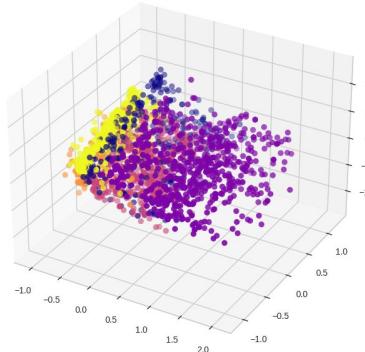


Best Model

```
OP = OPTICS(eps=0.4,min_samples = 250)
```

- OPTICS produces a hierarchical clustering result, providing a more **comprehensive view** of the data structure. This hierarchical information can be valuable for understanding varying densities and shapes of clusters within the dataset.
- **Robust to noise and outliers** in the data. It can identify clusters of varying shapes and sizes while gracefully handling data points that do not belong to any cluster.
- Well-suited for datasets with **clusters of varying densities**. It can adapt to clusters that have different levels of density, making it effective in scenarios where traditional methods like K-Means may struggle

OPTICS Clustering

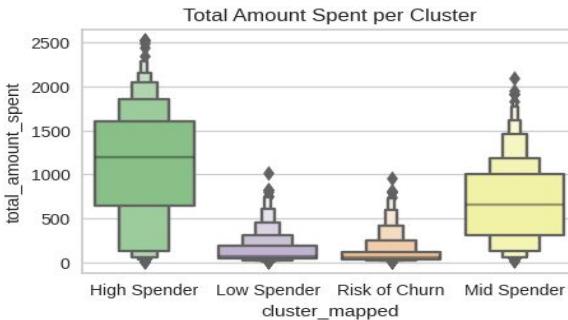
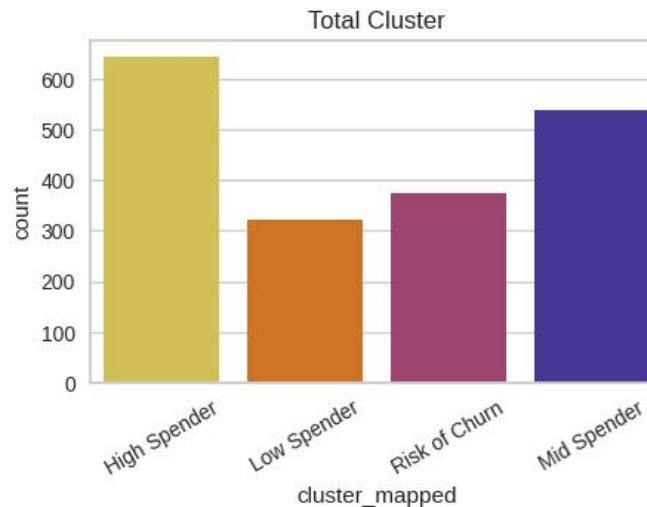


Best Model

```
OP = OPTICS(eps=0.4,min_samples = 250)
```

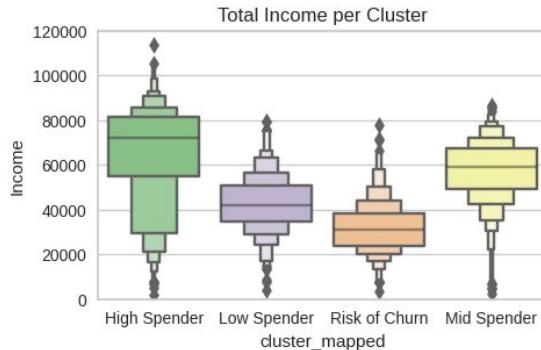
- eps: We set the number of estimators to 0.4. This parameter controls the maximum distance between two samples for one to be considered as in the neighborhood of the other. In OPTICS, it is often referred to as the "reachability distance." The eps parameter influences the sensitivity of the algorithm to the density of clusters. Smaller values of eps result in more fine-grained clusters. **We chose a relatively smaller number of estimators to ensure that the model has better separation in cluster whilst still ensuring that we are not making it too small which will cause most of the data points to be treated as noise. Setting eps as 0.4 results in only about 10% of data points to be treated as noise.**
- min_samples: We set the min_samples to 250. This parameter specifies the minimum number of samples in a neighborhood for a data point to be considered a core point. A core point is a point that has at least min_samples within its epsilon neighborhood. Increasing min_samples can lead to a more robust identification of dense regions, filtering out smaller clusters and noise. **We chose a moderate value of 250 to balance the model's capacity to capture complex patterns while avoiding overfitting.**

OPTICS Cluster Interpretation

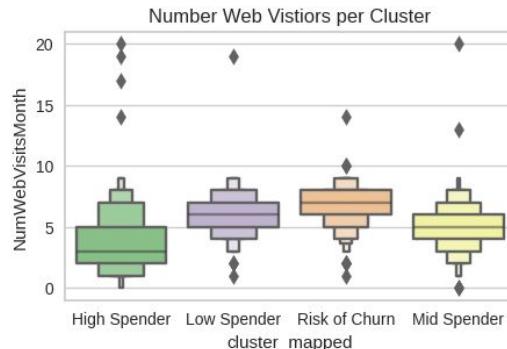


cluster_optics	total_amount_spent								
	count	mean	std	min	25%	50%	75%	max	
0	643.0	1113.525661	645.533234	6.0	647.00	1193.0	1599.00	2525.0	
1	539.0	679.024119	453.989709	15.0	305.00	653.0	1001.00	2092.0	
2	322.0	141.897516	170.814355	8.0	38.25	68.0	185.75	1019.0	
3	374.0	118.919786	158.494218	5.0	37.00	57.0	119.75	960.0	

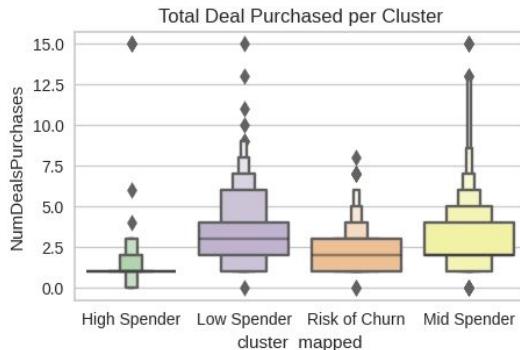
OPTICS Cluster Interpretation



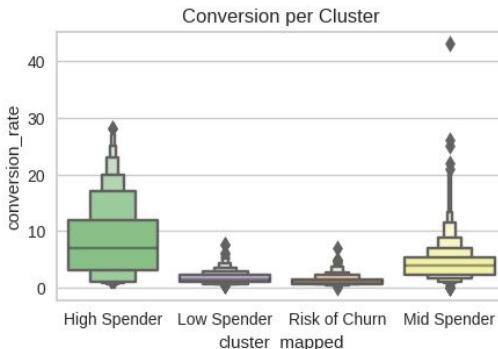
High spender customers have the highest income as expected



However, High spender customers rarely visit the web.

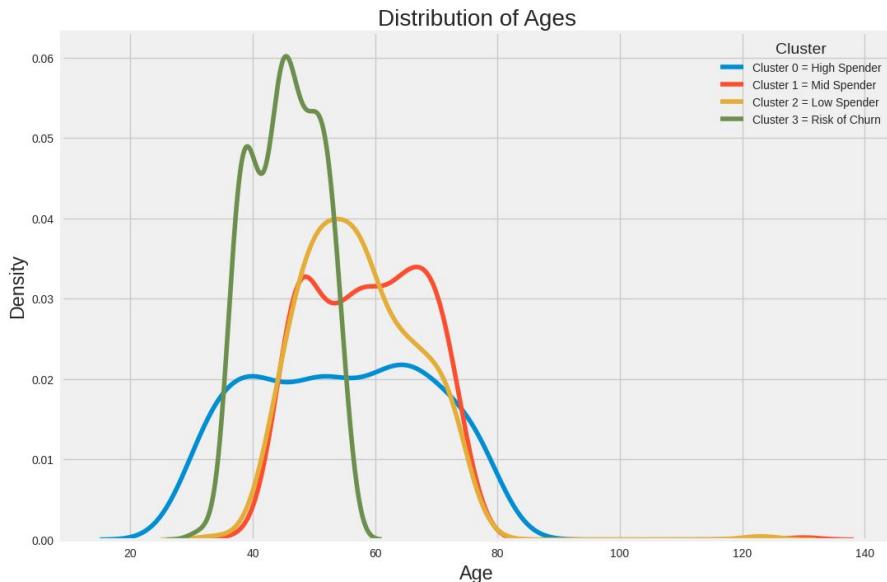


High spender customers rarely make purchases with discounts.

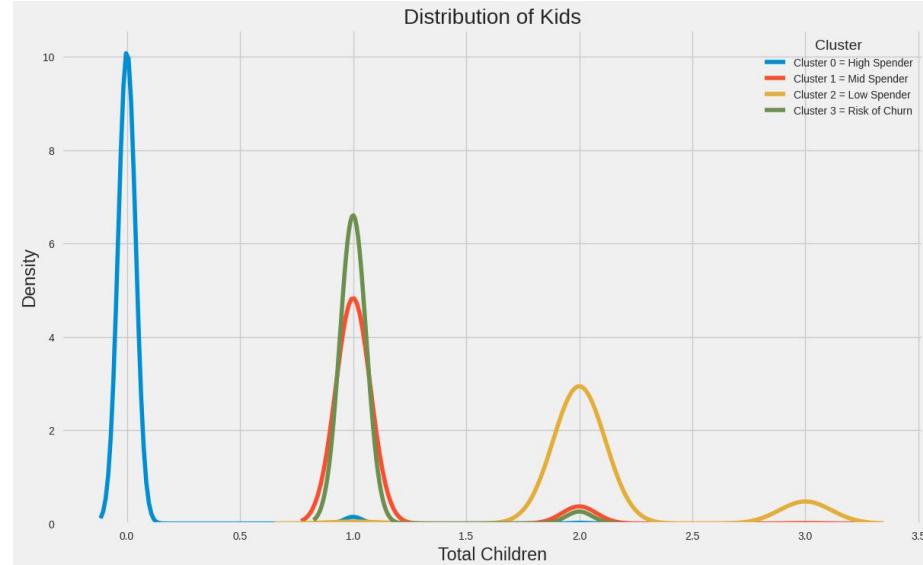


High spender customers also have high conversion rate.

OPTICS Cluster Interpretation



High spender customers are highly distributed within their age range. The most important pattern is that young generations (Gen Y) are most likely to churn. When they don't like the product anymore, they will go.



Almost all the high spender customers do not have kids. On the other side, low spender customers have from 2-3 kids.

Clustering + Classification Architecture

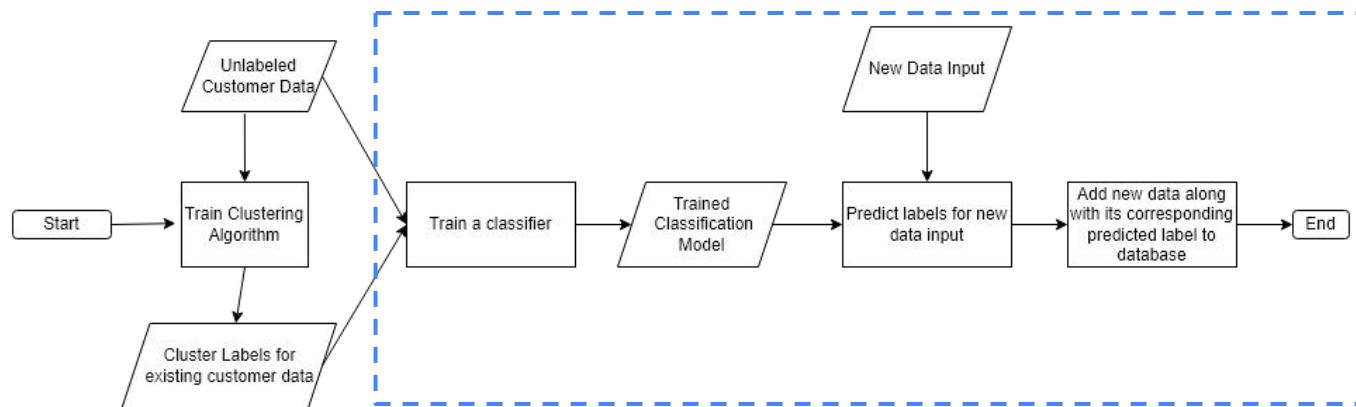
Motivation

Although clustering can be done to segment existing customers, not all clustering algorithms are able to efficiently cluster new data points (future customer data).

Idea

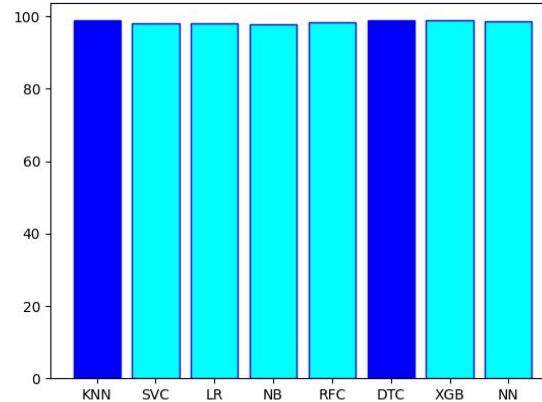
Train a robust classifier using the clustered customer data to efficiently segment new customers

Implementation



Classification for New Data Point Results

Model Names	Accuracy
K-Nearest Neighbours	98.91%
Support Vector Machine (Linear Kernel)	98.23%
Logistic Regression	98.07%
Naive Bayes	97.92%
Random Forest Classifier	98.43%
Decision Tree Classifier	98.91%
XGBoost Classifier	98.89%
Sequential Neural Network (1 Hidden Layer)	98.77%



```
Model: "sequential"
Layer (type)      Output Shape       Param #
dense (Dense)    (None, 15)          570
dense_1 (Dense)   (None, 15)          240
dense_2 (Dense)   (None, 4)           64
=====
Total params: 874 (3.41 KB)
Trainable params: 874 (3.41 KB)
Non-trainable params: 0 (0.00 Byte)
```

Neural Network Structure

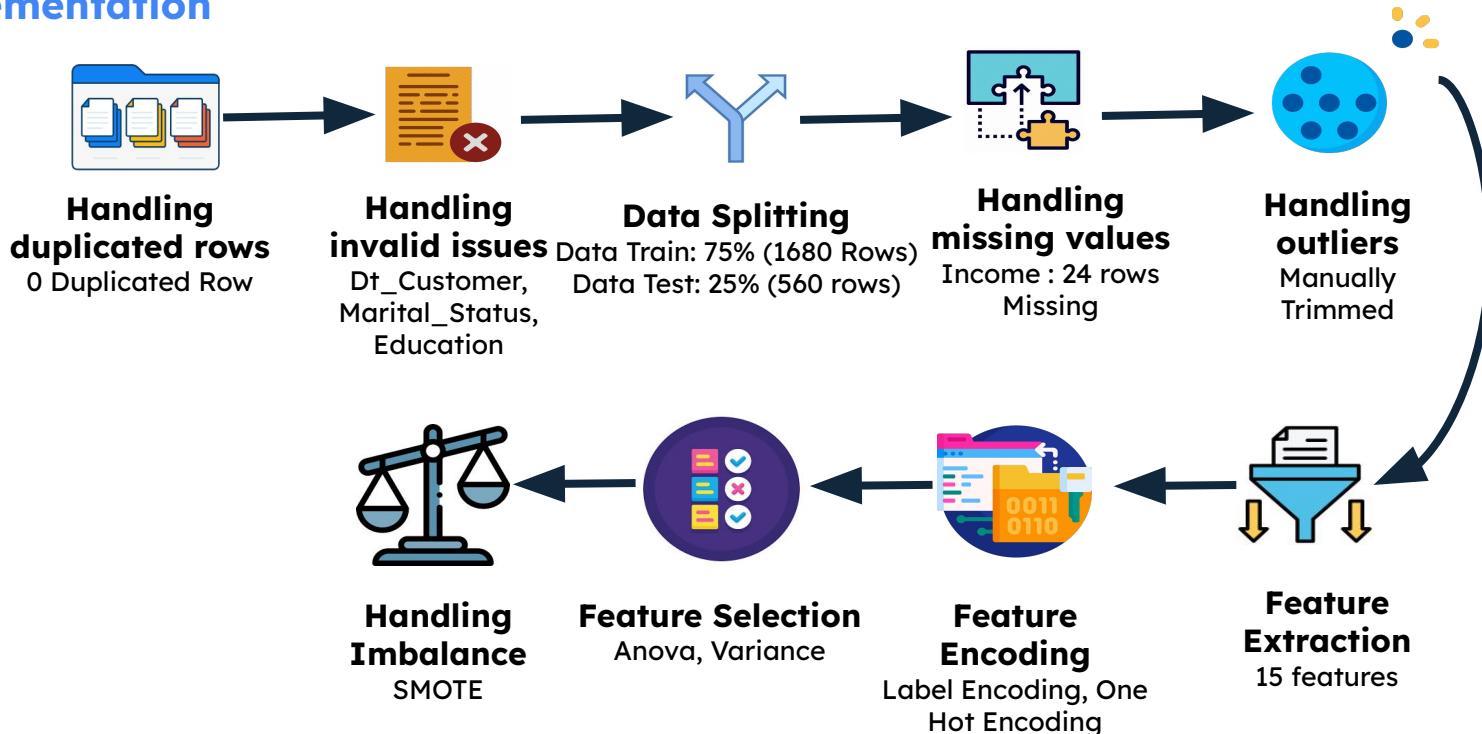
- **Input Layer:** 15 neurons; Rectified Linear Unit (ReLU) activation function; 37 input features.
- **Hidden Layer:** 15 neurons; ReLU activation function.
- **Output Layer:** 4 neurons; Softmax activation function.
- **Compilation:**
 - Loss function: categorical cross entropy
 - Optimizer: Adaptive Moment Estimation (Adam)
 - Evaluation Metric: Accuracy

Classification based on Customer Response to Past Campaigns Architecture

Motivation

After getting the cluster, we still need to predict whether the customers will join in the next campaign or not. We will use supervised learning binary classification methods to segment the customer.

Implementation



Classification Attributes

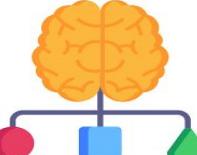
AcceptedCmp1
Accepted mp2
AcceptedCmp3
AcceptedCmp4
Age ★
Education
Income
Lifetime ★

Married
MntGoldProds
NumCatalogPurchases
NumDealsPurchases
NumWebVisitsMonth
Recency
Recency_sgmt ★
Teenhome
Total_cmp ★

**There are 17 features that will be used for modeling.
The stars is the result of Feature Engineering**

Classification Modeling & Evaluation Metrics

Positive = Response Campaign
Negative = No Response Campaign

	 Model Prediction	 Actual / Reality
False Positive Main target to be reduced	 Response	 No Response
False Negative Second target to be reduced	 No Response	 Response
	FALSE	TRUE

Classification Modeling & Evaluation Metrics

1.

Precision

Primary Metrics
Evaluation

Reducing False Positive

Customers are predicted to respond ✓
but in reality they don't ✗

↑ Increase Response Rate
Minimize Cost Marketing ↓

2.

Recall

Secondary Metrics
Evaluation

Reducing False Positive

Customers are predicted not to respond ✗
but are interested in responding campaign ✓

Optimizing Revenue Rates



3.

F1 Score

Combined Metrics
Evaluation

Necessary for our problem, we need a

Trade-off between Precision and Recall

F1 score is defined as the harmonic average of
Precision and Recall

Classification Results

Model	Accuracy	Precision	Recall	F1 Score	Cross Val F1 (k=5)	ROC AUC	Cross Val ROC AUC (k=5)
LogisticRegression	0.861000	0.523000	0.675000	0.589000	0.510000	0.900000	0.885000
LinearSVC	0.859000	0.520000	0.614000	0.564000	0.448000	0.895000	0.884000
MLPClassifier	0.846000	0.483000	0.518000	0.500000	0.558000	0.830000	0.887000
GradientBoostingClassifier	0.846000	0.482000	0.494000	0.488000	0.487000	0.725000	0.752000
KNeighborsClassifier	0.845000	0.482000	0.639000	0.549000	0.400000	0.842000	0.810000
DecisionTreeClassifier	0.838000	0.456000	0.494000	0.474000	0.488000	0.699000	0.718000

The **Precision**, **Recall**, and **F1 Score** values on the best evaluation are produced by **Logistic Regression** and **Support Vector Machine**

WORKFLOW ML

Stage 1



Introduction

Stage 2



Data Collection
& EDA

Stage 3



Modeling &
Results

Stage 4



Implications

Stage 5



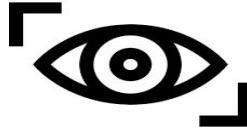
Recommendation
& Conclusion

Implications



Theoretical Implications

Which is the best approach between RFM analysis, clustering combined with classification, and a standalone classification model based on customer response to previous campaigns, in terms of their abilities to segment customers and provide actionable business insights for optimizing marketing strategies?



Marketing Implications

Given our clusters, what segments do we get? How much the potential GMV and optimization cost? What is the personalized approach marketing strategy for each segment?



Managerial Implications

How might a company impact a customer segmentation microservices in their internal stack?

Theoretical Implications

Approach	Performance (silhouette score - accuracy)	Results Interpretability	Method	Computational Cost
RFM Analysis	0.528 - N/A	Resulting clusters can be interpreted and turned into actionable insights	Data Analysis (No Machine Learning)	Low
Clustering + Classification (OPTICS + KNN/DTC)	0.72 - 98.91%	Resulting clusters can be interpreted and turned into actionable insights. Hard to find feature importance.	Unsupervised + Supervised Machine Learning	High
Classification based on customer response to past campaigns (Logistic Regression)	N/A - 86.1%	Low model interpretability, challenging to derive actionable insights.	Supervised Machine Learning	Moderate

Marketing Implications

AVERAGE TOTAL SPEND MID & HIGH CLUSTER



HIGH SPENDER (N=643)

1113,5



MID SPENDER (N=539)



679

POTENTIAL REVENUE INCREASE PER PERSON MID-> HIGH

FORMULA

MEAN TOTAL AMOUNT HIGH-
MEAN TOTAL AMOUNT MID



434

ASSUMING 15% CONVERSION RATE FROM MID-> HIGH, IT WILL GET US
POTENTIAL ADDITIONAL REVENUE



FORMULA

NUMBER OF CONVERSION X
POTENTIAL REVENUE INCREASE PER PERSON

\$ 35,154

Assuming \$40 budget per person, we will get

PROFIT



FORMULA

POTENTIAL ADDITIONAL REVENUE
-TOTAL PROMOTIONAL COST

\$ 13,594

Marketing Campaign for Each Clusters



High Spender “VIP Shoppers”

- Personalized Concierge Service
- Exclusive Preview Sales
- Elite Loyalty Program
- Premium Product Bundles



Middle Spender “Smart Shoppers”

- Upgrade & Save Campaign
- Flash Sales & Limited Offers
- Referral Rewards Program
- Personalized Product Recommendation



Marketing Campaign for Each Clusters



Low Spender “Value Seekers”

- First Purchase Discounts
- Frequent Shopper Discounts
- Bundle and Save
- Social Media Contests

The collage includes:
1. A dark banner with the text "10% OFF" and "First Time Purchase Discount".
2. An illustration of five people holding shopping bags, with one person pointing to the text "Frequent Buyer Programs".
3. A screenshot of a website showing "SUPER SAVER BUNDLES" with various travel accessories like wallets and pouches.



Risk of Churn “Win-Back Warriors”

- Reactivation Email Campaign
- Survey and Feedback Collection
- Exclusive Offers to Win-Back Customers



Managerial Implications

Goal: Building an automated, scalable, efficient, and accurate customer segmentation service.

- Data Processing Considerations: Batch Processing vs **Real-Time Stream Processing**
 - Batch Processing: Processing data in batches in regularly scheduled intervals
 - Stream Processing: Real-time processing of data streams as it is continuously generated or received
- If we use batch processing at long intervals, reclustering the entire customer database might be feasible, albeit inefficient
- If we use stream processing, reclustering the entire customer database is infeasible

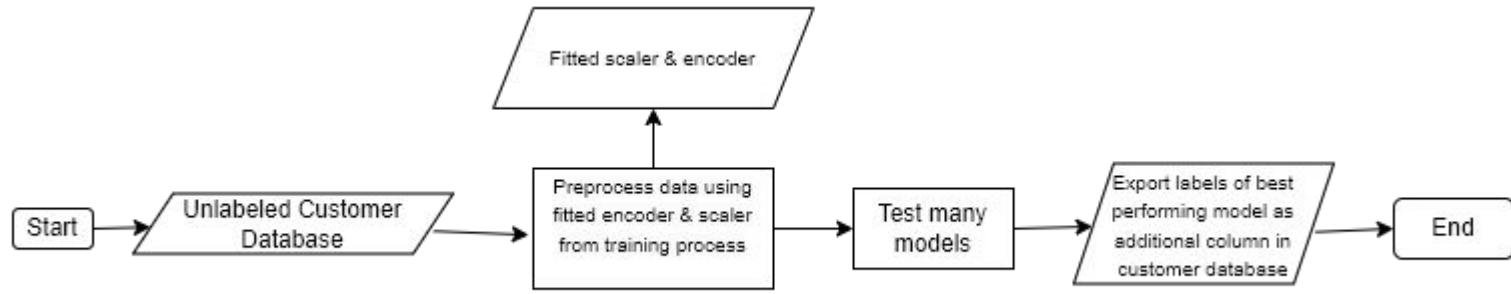
Solution

4 Step Process:

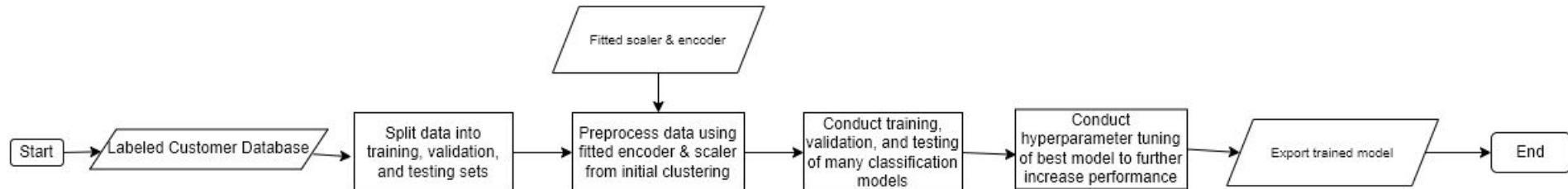
1. Segmentation of existing customers using clustering methods
2. Train a classification model using the cluster labels from the previous step as ground truth label for efficiency and scalability of future inference
3. Segment new customers using the trained classification model
4. Do model monitoring & maintenance to ensure reliability of results

Implementation of a Clustering + Classification Solution For Customer Segmentation

1. Segmentation of existing customers using clustering methods

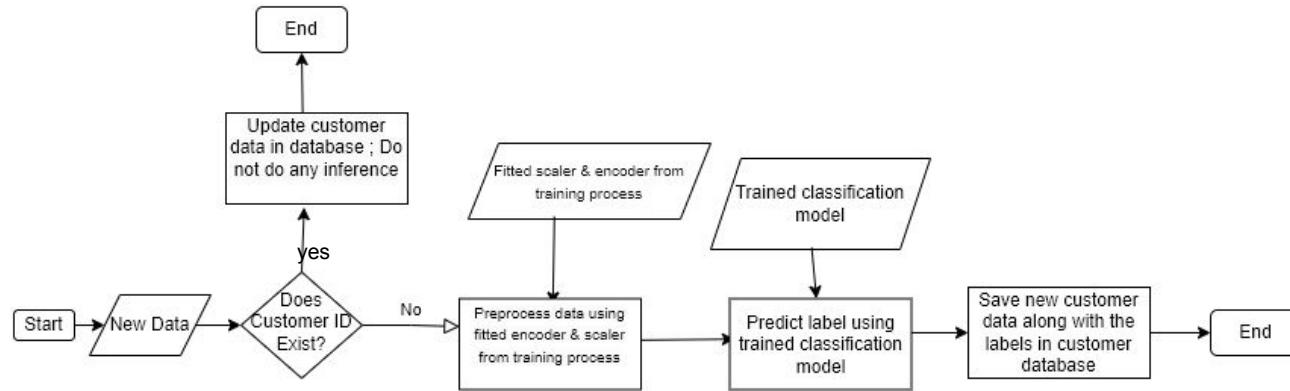


2. Train a classification model using the cluster labels from the previous step as ground truth label for efficiency and scalability of future inference



Implementation of a Clustering + Classification Solution For Customer Segmentation

3. Segment new customers using the trained classification model



4. Do model monitoring & maintenance to ensure reliability of results. Examples:

- **Data Drift Monitoring:** Monitor Changes in statistical properties of input features. If changes are too drastic, reclustering the entire database might be needed to ensure accuracy of customer segmentation.
 - Example Tools: Tensorflow Data Validation, Apache Drift
- **Cluster Stability Monitoring:** Monitor stability metrics of our clusters over time. In this case, an example is silhouette score. If cluster stability degrades significantly, reclustering is necessary.
- **Classification Model Performance Monitoring:** Monitor the performance of the classification model over time. Using performance metrics such as accuracy, Gini and KS -Statistics, to ensure model reliability. If metrics degrades, then retraining our classification model is necessary.

WORKFLOW ML

Stage 1



Introduction

Stage 2



Data Collection
& EDA

Stage 3



Modeling &
Results

Stage 4



Implications

Stage 5



Recommendation
& Conclusion

Novelty & Significance

Novelty

1

Employ different architecture

In this project we employed many different architectures. We used 10+ different models and tune several models to have better silhouette score result

2

Pipeline for Customer Segmentation problem

We also created a pipeline of customer segmentation that can be used to deploy the dashboard in the next stage.

Significance

1

Low Computational Cost

Our model can be converged very fast within just several iterations

2

Offering Advanced Method to solve bigger problem

OPTICS can be used to solve bigger problem of customer segmentation with high silhouette score.

Limitation & Future Work

Limitation

1

Limited Time

We were only given a month to develop a comprehensive end to end project.

2

Model is still under development

As can be seen in the performance report, our models' capabilities are still towards better development. We will use explore more hyperparameter tuning

3

Limited Data

As can be seen in the data collection, the data that we have is relatively small.

Future of Work

1

Enhancing model

Create much more varieties of models and maybe some combinations between models

2

Trying on different dataset

The results from different dataset will make the results more generalized and applicable

3

Making a better pipeline

A user friendly dashboard must be made so that the users can use the dashboard easily.



Q&A
ASK AWAY!

THANK YOU!

For Being part of our Capstone Journey!