# [Business optimization in predicting customer churn : a machine learning approach]

## [ANN Team]

TURNAMEN
SAINS DATA
NASIONAL
2022

# Table of contents

**01** About us

**02** Background

**03** EDA

**04** Data Pre-Processing

**05** Churn Predict

**06** Survival Analysis

**07** RFM Segmentation

**08** Summary & Recommendation

# Archie Citra Muhammad | archiecm09@gmail.com

TTL : Sragen, 22 Sept 1994
No. Hp : 08112165945
Address : Sragen Tengah, Sragen , Jawa Tengah
Social Media : @archiecm

# Nur Amilah|
# nuramilahnuramilah@gmail.com

TTL            : Tangerang, 16 May 2001
No. Hp         : 08159887509
Address        : Kp. Pagedangan, Kab. Tangerang, Banten.
Social Media   : @nuramilah_16

# Natalia Dinda Sartika Putri | nata.dsptr@gmail.com

TTL            : Tangerang, 09 June 2000
No. Hp         : 085771768020
Address        : Jl. Raya Mauk No.45, Jatiwaringin, Tangerang
                 Regency, Banten.
Social Media   : @nata.dsptr_

# FOREWORD

## Ekonomi Digital?
(Brynjolfsson & McAfee, 2014)

## Business Optimization?
(Apte, 2010)

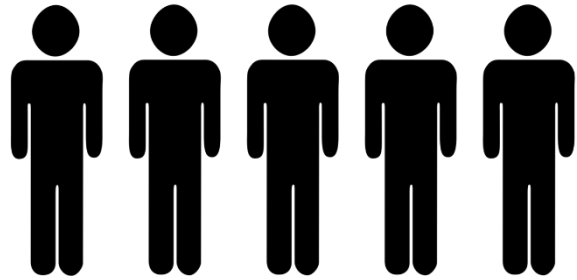## Machine Learning?
(Al-Sahaf et al., 2019)

## Customer Churn?
(Masarifoglu & Buyuklu, 2019)

## Cashback Amount?
(Pinem et al., 2020)

# PROBLEM STATEMENT

China Internet Network Information Center (CNNIC)

E-commerce customer churn rate is up to **80%** compared with the traditional business customer management (Wu & Meng, 2016)

Churn **16.8%**

## Goals

Memprediksi pelanggan churn rate dan memberikan rekomendasi kepada business team agar perusahaan mampu menerapkan strategi customer retention.

## Objective

Membentuk sebuah model machine learning dengan false negative terkecil, mengidentifikasi prediktor/faktor yang berpengaruh terhadap churn rate dan lost opportunity customer churn, Memprediksi customer yang berpotensi churn dengan machine learning model. Serta memberikan insight & rekomendasi untuk mengidentifikasi prediktor/faktor yang berpengaruh terhadap churn rate melalui cashback amount.

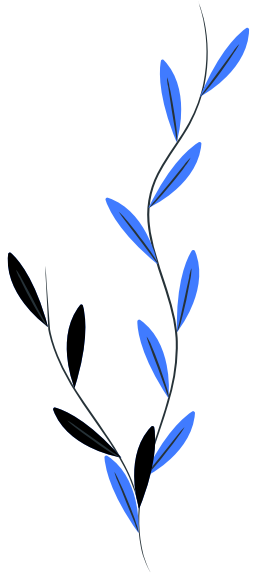## Business Matrix

$$Churn\ Rate = \frac{CUSTOMER\ CHURN}{TOTAL\ CUSTOMERS}$$

$$Lost\ Opportunity = Total\ Customers\ Complain\ \&\ Berpotensi\ Churn \times Average\ Monthly\ Spending\ User$$

# Data Overview

```
 #   Column                      Non-Null Count   Dtype
---  ------                      --------------   -----
 0   CustomerID                  5630 non-null    int64
 1   Churn                       5630 non-null    int64
 2   Tenure                      5366 non-null    float64
 3   PreferredLoginDevice        5630 non-null    object
 4   CityTier                    5630 non-null    int64
 5   WarehouseToHome             5379 non-null    float64
 6   PreferredPaymentMode        5630 non-null    object
 7   Gender                      5630 non-null    object
 8   HourSpendOnApp              5375 non-null    float64
 9   NumberOfDeviceRegistered    5630 non-null    int64
 10  PreferedOrderCat            5630 non-null    object
 11  SatisfactionScore           5630 non-null    int64
 12  MaritalStatus               5630 non-null    object
 13  NumberOfAddress             5630 non-null    int64
 14  Complain                    5630 non-null    int64
 15  OrderAmountHikeFromlastYear  5365 non-null   float64
 16  CouponUsed                  5374 non-null    float64
 17  OrderCount                  5372 non-null    float64
 18  DaySinceLastOrder           5323 non-null    float64
 19  CashbackAmount              5630 non-null    float64
dtypes: float64(8), int64(7), object(5)
```
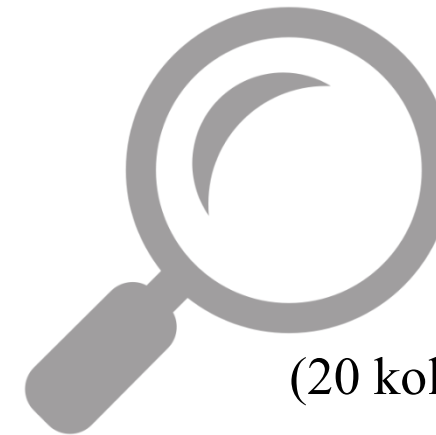
Source : *Kaggle*

**Target Variable :**

*Churn (Classification Model)*

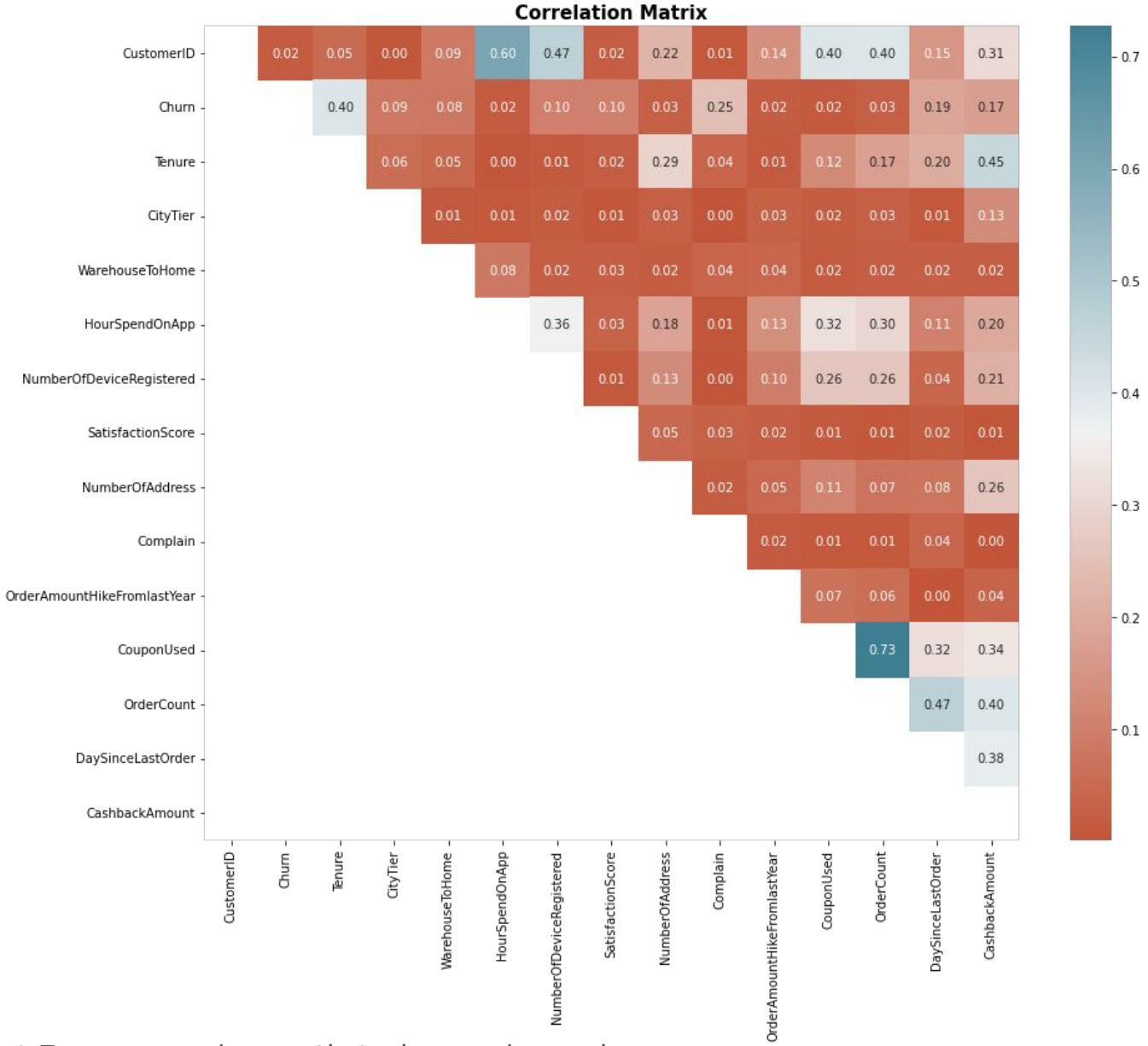*Tenure (Regression Model)*

**Informasi Dataset?**

(20 kolom dan 5630 baris, 19 variabel input, jenis data, 1 var.target)

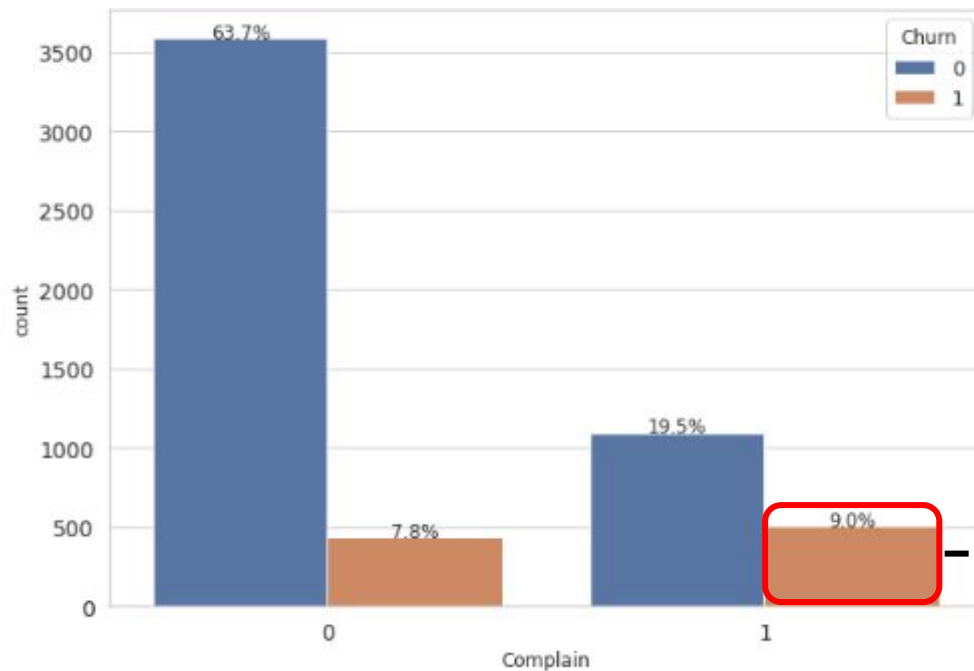**Dtype (Data Type)**

# EXPLORATORY DATA ANALYSIS

| Column | Correlation_ratio |
|---|---|
| Tenure | 0.40 |
| Complain | 0.25 |
| Cashbackamount | 0.15 |
| Daysincelastorder | 0.15 |
| Numberofdeviceregistered | 0.11 |
| Satisfactionscore | 0.11 |
| Citytier | 0.08 |
| Warehousetohome | 0.07 |
| Numberofaddress | 0.04 |
| Ordercount | 0.03 |
| Hourspendonapp | 0.02 |
| Couponused | 0.01 |
| Orderamounthikefromlastyear | 0.01 |

## Correlation Matrix



correlation with target **Response** is worth to be reviewed.

# INSIGHTS
## (Comparison Complain to Churn and Not Churn)



1. Customer dengan **churn tertinggi** sebesar

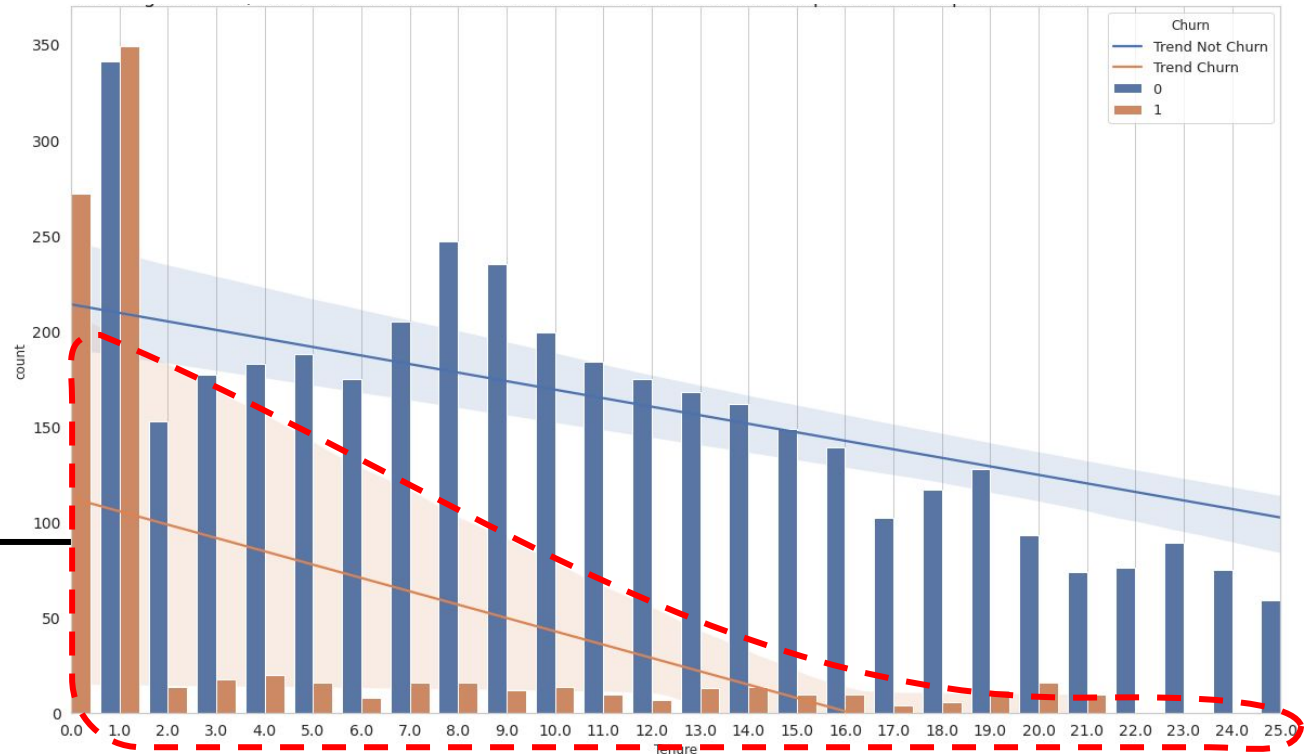   **9.0%** berada pada **customer complain**.


2. Customer dengan **churn terendah** sebesar

   **7.8%** berada pada **customer tiidak complain**.

**Semakin meningkatnya complain customer maka semakin tinggi tingkat churn rate.**
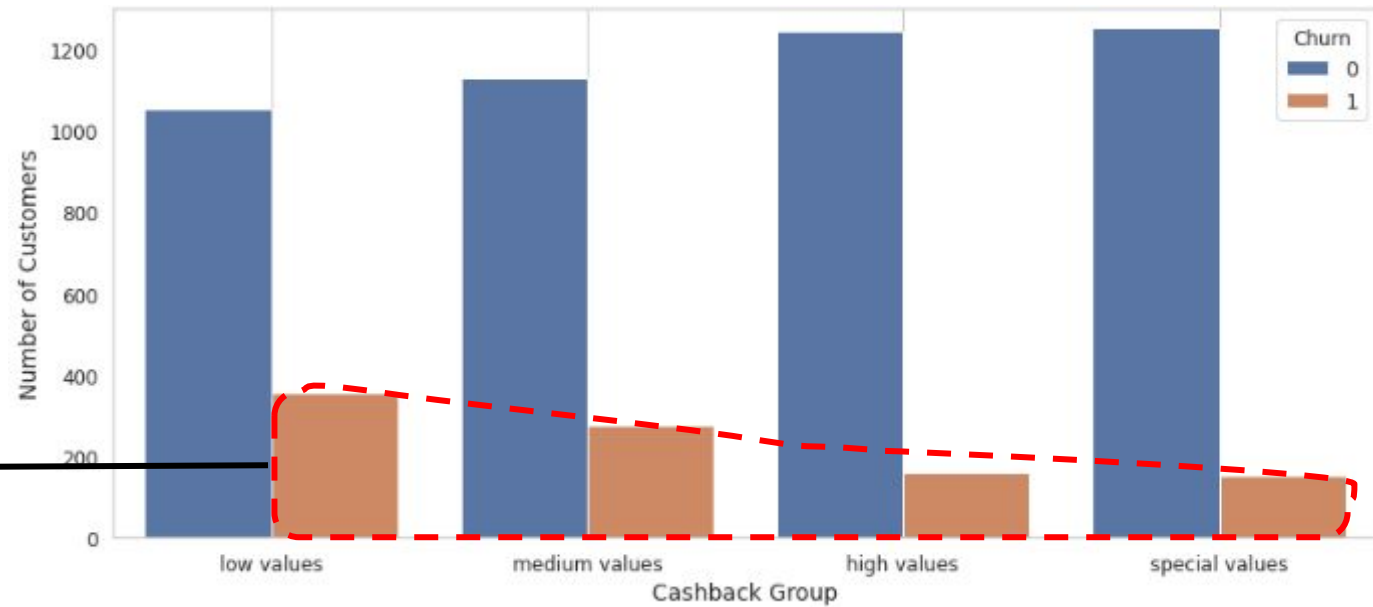
# INSIGHTS
## (Churn and Not Churn)

The longer tenure, the lower number of churns. And not churn has a steeper trend compared to Churn.

# INSIGHTS
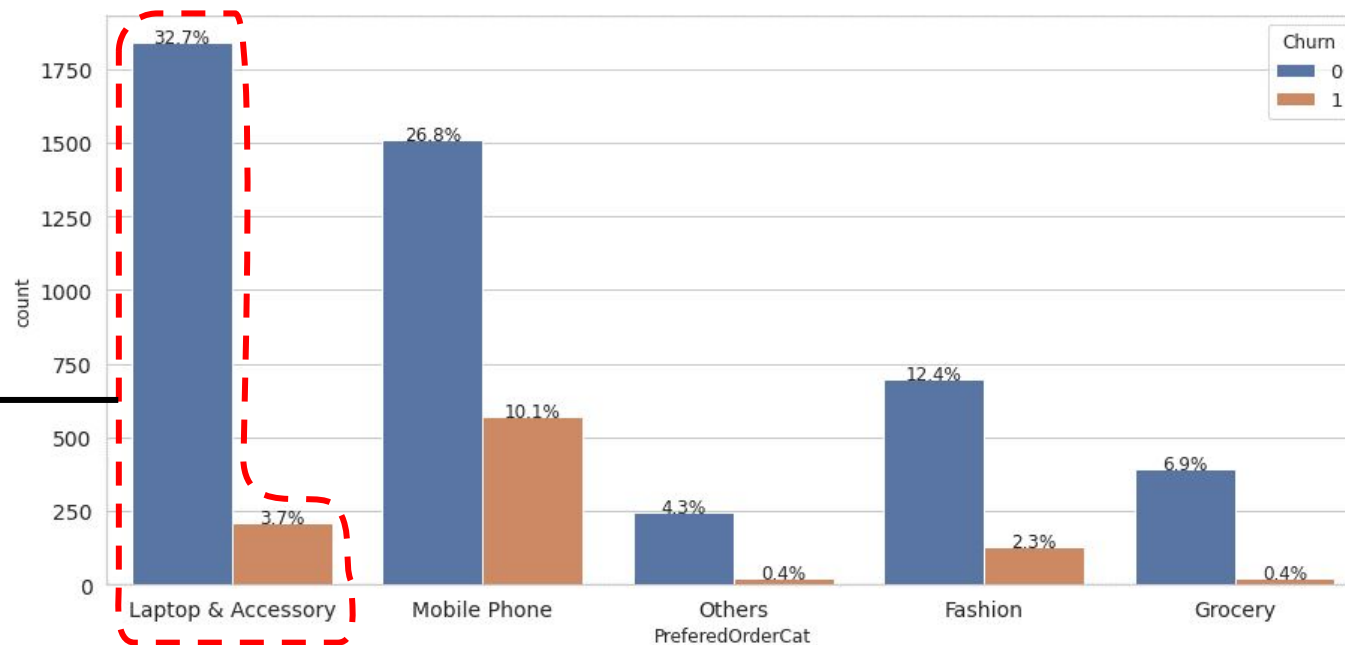## (Distribution of Cashback Customers)



Increase Cashback Amount has trend Positive in Not Churn On the contrary Increase Cashback Amount has trend Negative in Churn
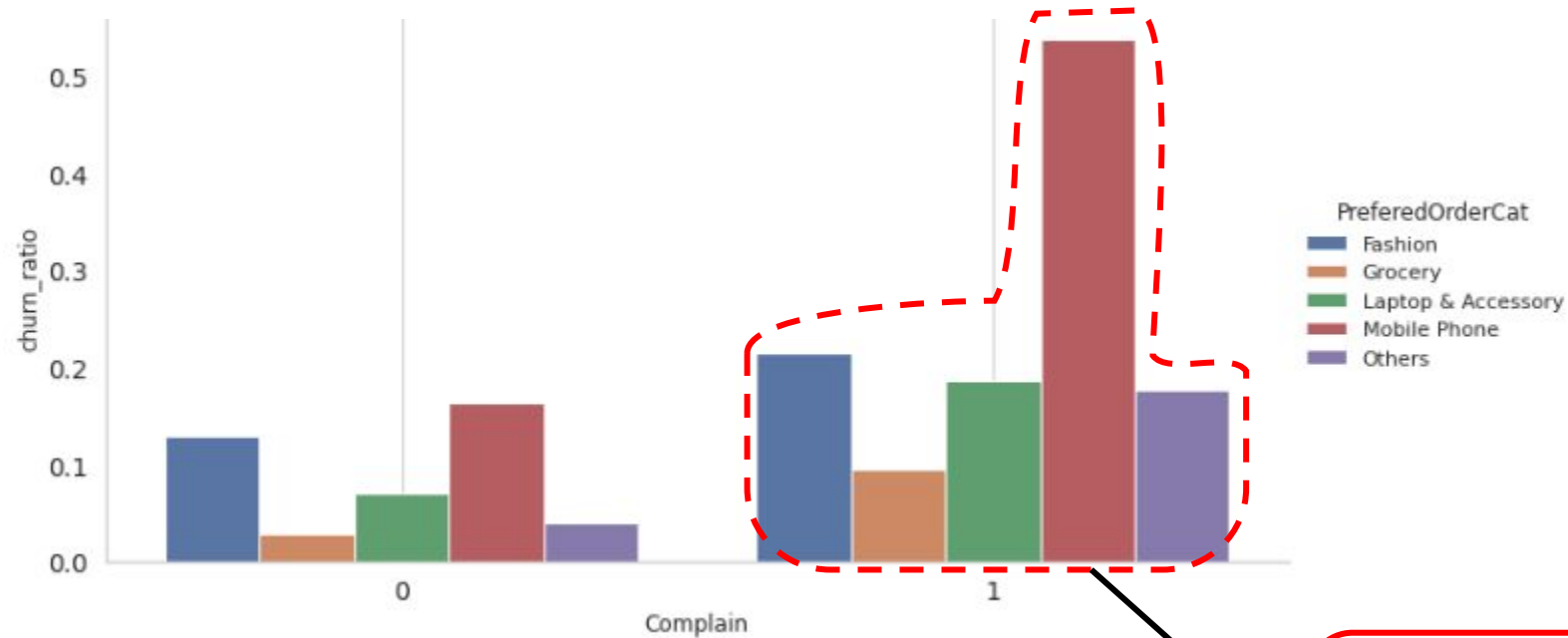
# INSIGHTS
## (Prefered Order Categories Customer)



Customer who ordered Laptop and Accessory has a significant number of Not Churn compared same order category with Churn.
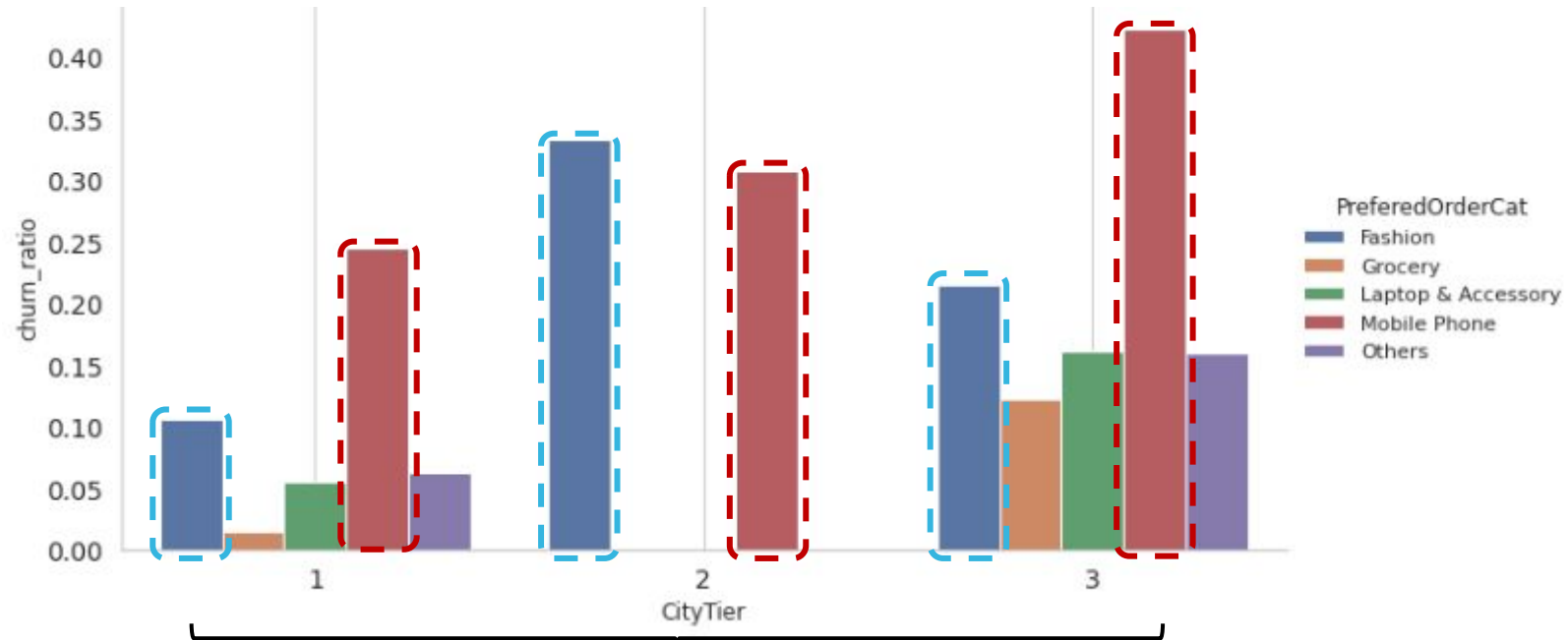
# INSIGHTS
## (Distribution of Complain & Order Categories vs Ratio Churn)

**PreferedOrderCat**
- Fashion
- Grocery
- Laptop & Accessory
- Mobile Phone
- Others

**Customers with complaints have a ratio churn increase in all order categories.**

# INSIGHTS
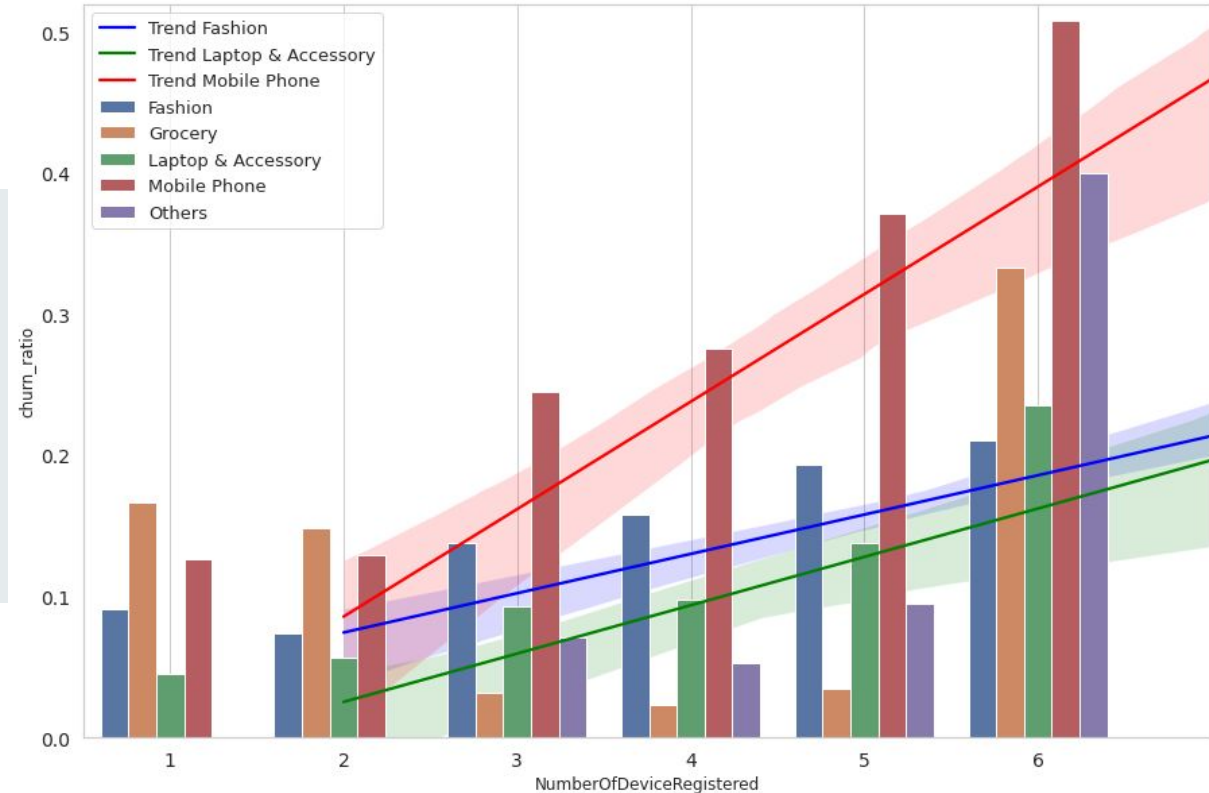## (Distribution of Complain & Order Categories vs Ratio Churn)



More City Tier increase, more ratio churn increase in **Fashion** and **Mobile Phones**

# INSIGHTS
## (Distribution of Complain & Order Categories vs Ratio Churn)

> " More Number Of Device Registered increased and more ratio churn increased in **Fashion, laptops & accessories, and Mobile Phones**.
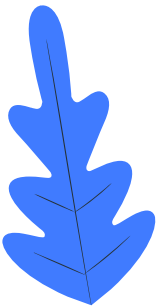
# Data Pre-Processing

## Data Cleaning

Check Irrelevant Data

Check Missing Data

Check Duplicate

Check Outlier

## Feature Encoding

One Hot Encoder
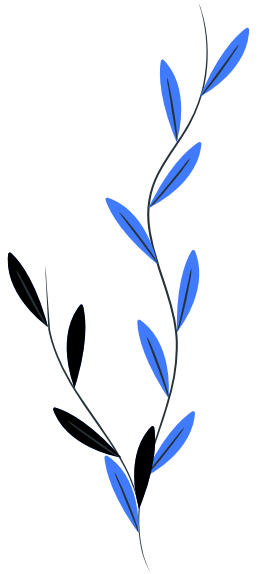
Simple Imputer

Iterative Imputer

## Transforming

Pipeline

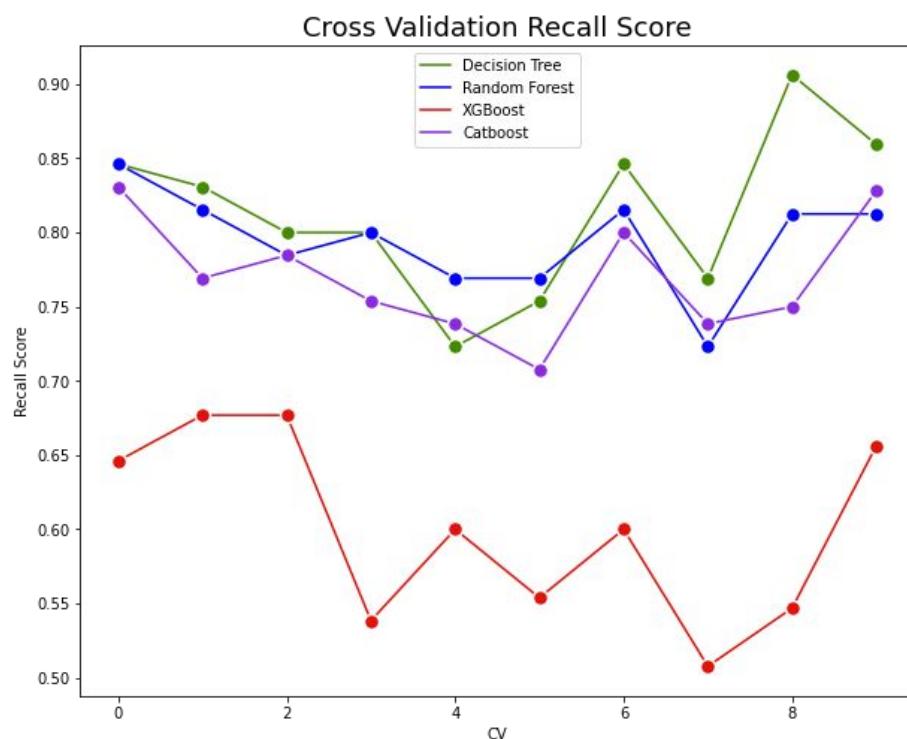Robust Scaler

Standard Scaler

# Predict Churn

Selection Models & Cross Validation, Handling Imbalance, Hyperparameter Tuning, Feature Importance with SHAP

# Model Selection and Cross-Validation



Cross Validation Recall Score

| Models | Mean | Standar Deviasi | Recall |
|---|---|---|---|
| Decision Tree | 0.803951 | 0.038707 | 0.863095 |
| Catboost | 0.780246 | 0.047106 | 0.809524 |
| Random Forest | 0.760799 | 0.048173 | 0.797619 |
| Xgboost | 0.609570 | 0.066454 | 0.553571 |
| Logistic Regression | 0.530575 | 0.051384 | 0.476190 |

**NB : Due to an imbalance dataset**

# Handling Imbalance

**NB : Due to imbalance dataset**

### CatBoost

| | Without | Undersampling | Oversampling |
|---|---|---|---|
| **Train Recall** | 0.953360 | 0.996913 | 0.999826 |
| **Test Recall** | 0.784038 | 0.928990 | 0.915144 |

### DECISION TREE

| | Without | Undersampling | Oversampling |
|---|---|---|---|
| **Train Recall** | 1.000000 | 1.000000 | 1.000000 |
| **Test Recall** | 0.836538 | 0.881202 | 0.819519 |

Catboost have best fit in undersampling and oversampling. But we choose undersampling because it has gap (train-test) smaller than other.

# CatBoost Classifier + Undersampling

## Confusion Matrix



```
classification_report before tuning:
              precision     recall    f1-score     support

          0       0.99       0.91        0.95         800
          1       0.68       0.94        0.79         162

   accuracy                              0.91         962
  macro avg       0.83       0.93        0.87         962
weighted avg       0.94       0.91        0.92         962
```

Recall 0.944444444444444

**Recall** ; How many customers did we correctly predict to take an interest with our product compared to all customers which are truly churn? **94%**

# CatBoost Classifier + Undersampling + Tuning



**Confusion Matrix**

|  | not churn | churn |
|---|---|---|
| not churn | 731 | 69 |
| churn | 3 | 159 |

```
classification_report after tuning:
              precision    recall  f1-score   support

           0       1.00      0.91      0.95       800
           1       0.70      0.98      0.82       162

    accuracy                           0.93       962
   macro avg       0.85      0.95      0.88       962
weighted avg       0.95      0.93      0.93       962
```
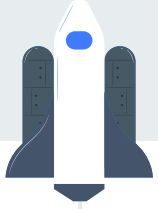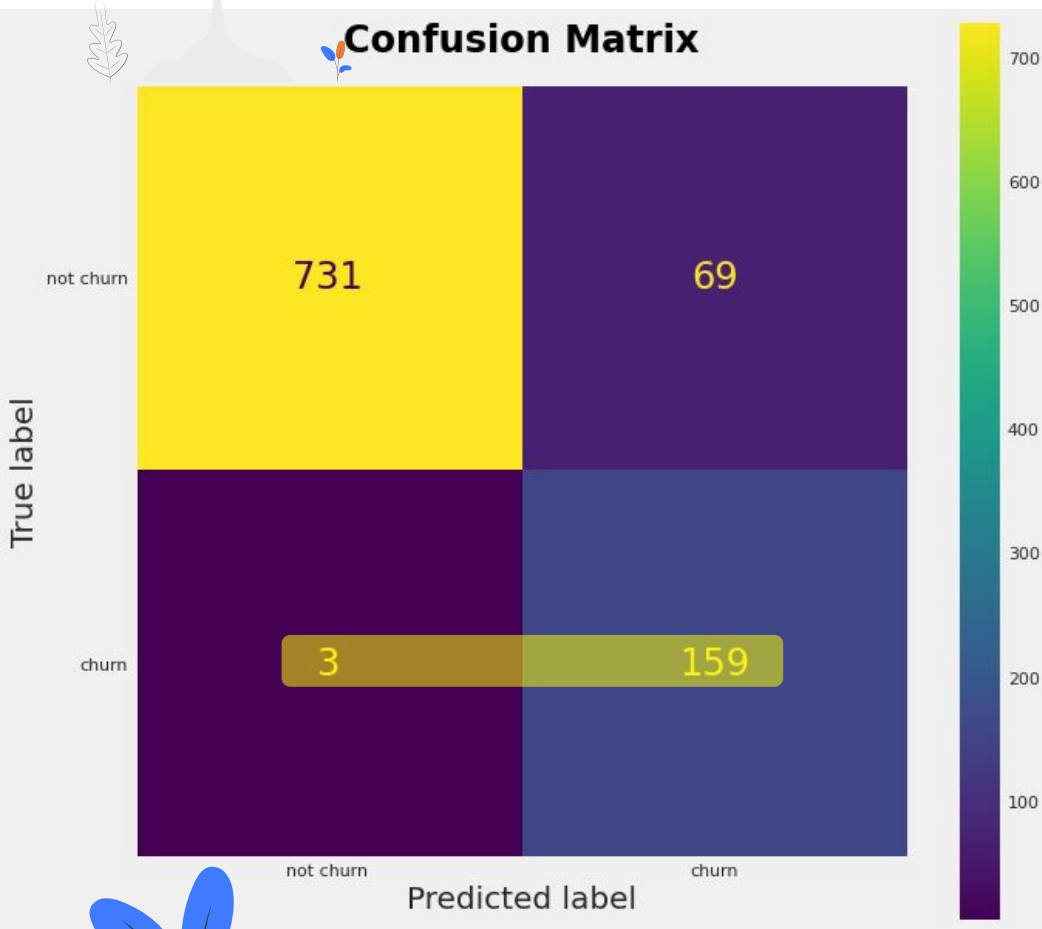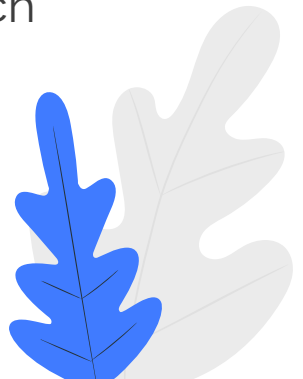
Recall 0.9814814814814815

**Recall** ; How many customers did we correctly predict to take an interest with our product compared to all customers which are truly churn? **98%**
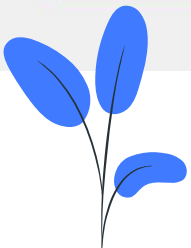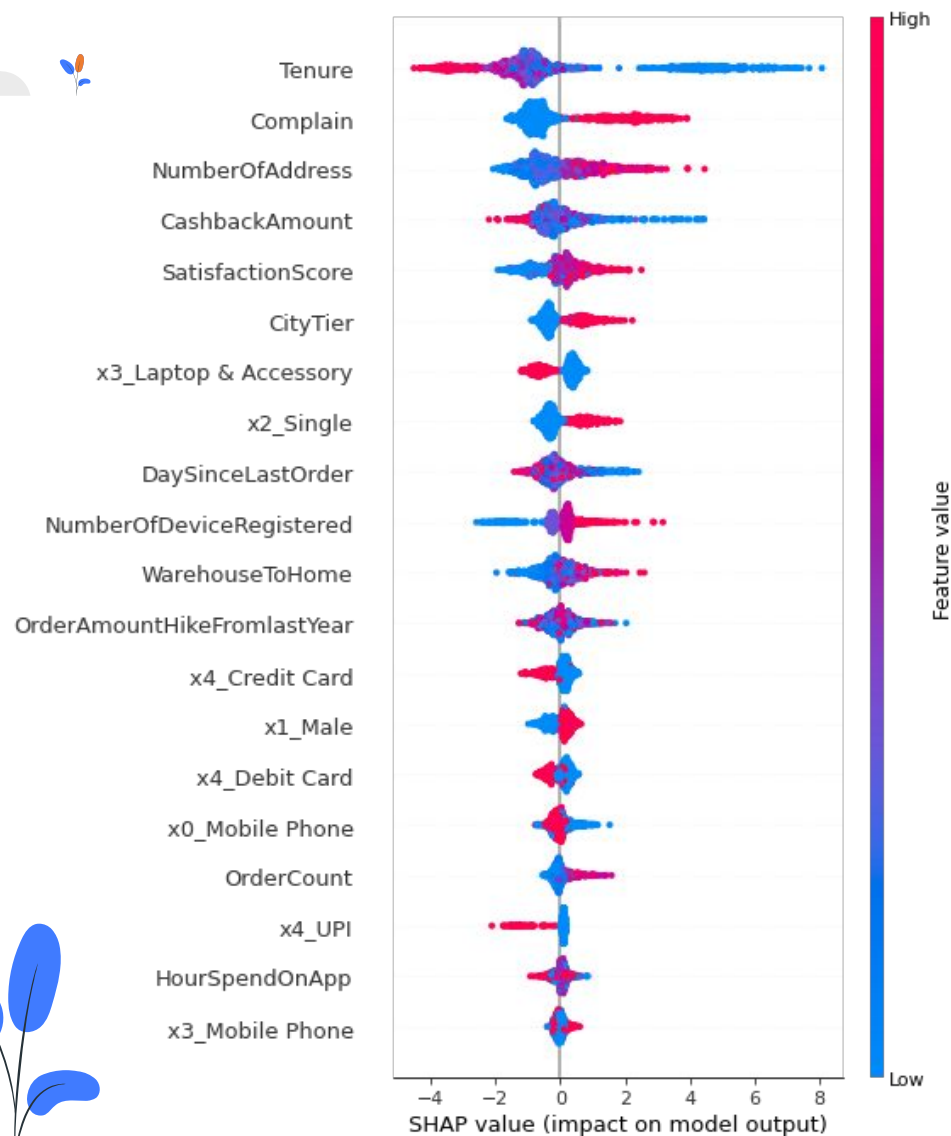
# CatBoost Classifier + Tuning + Undersampling



**Feature yang menghasilkan churn sebagai berikut**
- Tenure dengan nilai rendah
- Complain = 1
- Number of Address dengan nilai tinggi
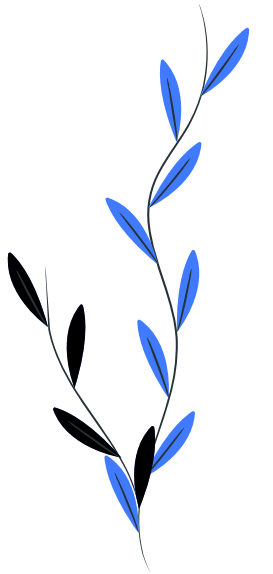- Cashback dengan nilai rendah

**Feature yang menghasilkan retention**
- Tenure dengan nilai tinggi
- Complain = 0
- Number of Address dengan nilai rendah
- Cashback dengan nilai tinggi

# Survival Analyst

Using Kaplan-Meier (KM) and COx Proportional Hazard (CPH) Model

# Kaplan-Meier(KM) Survival Curve



Kaplan-Meier Survival Curve — All Customers

- Pada **21** bulan pertama terjadi churn sebesar **841** customer **(16.6%)** & **tidak** terjadi **churn sampai 60 bulan**

- Pada **20 bulan pertama** terdapat **683** customer **at risk** yang artinya customer tersebut tidak terindikasi churn

- Pada 20 bulan pertama juga terdapat 3559 customer censored artinya customer tersebut terindikasi akan churn namun belum melakukannya

# COx Proportional Hazard (CPH) Model

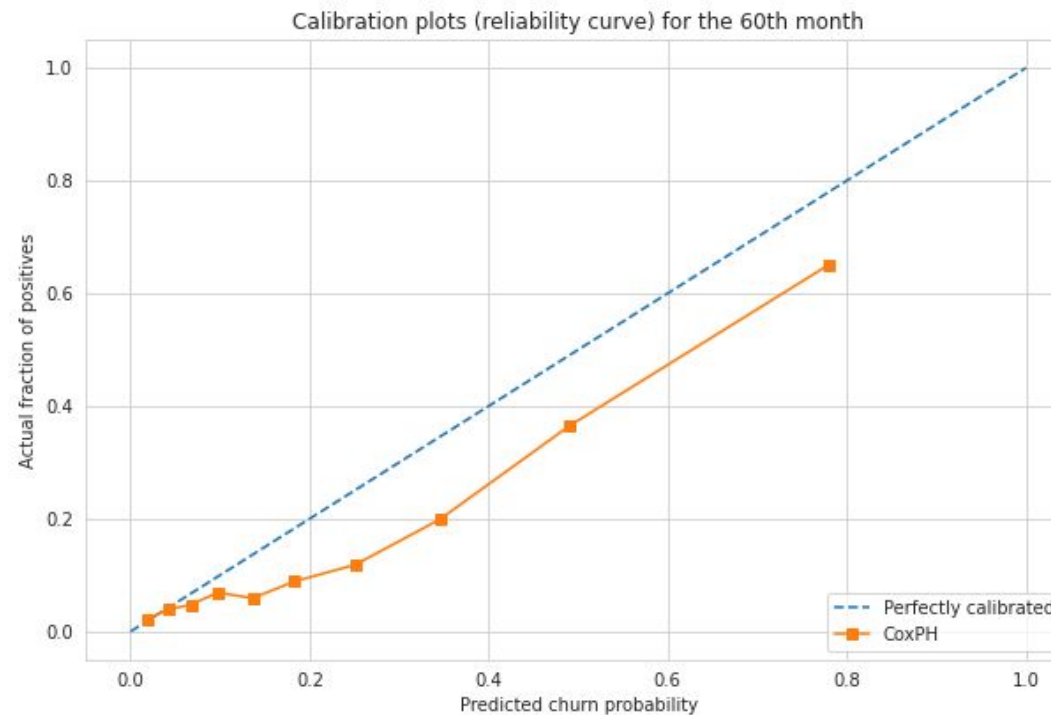| model | lifelines.CoxPHFitter |
|---|---|
| duration col | 'Tenure' |
| event col | 'Churn' |
| baseline estimation | breslow |
| number of observations | 5073 |
| number of events observed | 841 |
| partial log-likelihood | -6296.226 |
| time fit was run | 2022-11-19 08:47:34 UTC |
| model | base model |

| Concordance | 0.829 |
|---|---|
| Partial AIC | 12640.452 |
| log-likelihood ratio test | 1223.310 on 24 df |
| -log2(p) of ll-ratio test | 805.834 |



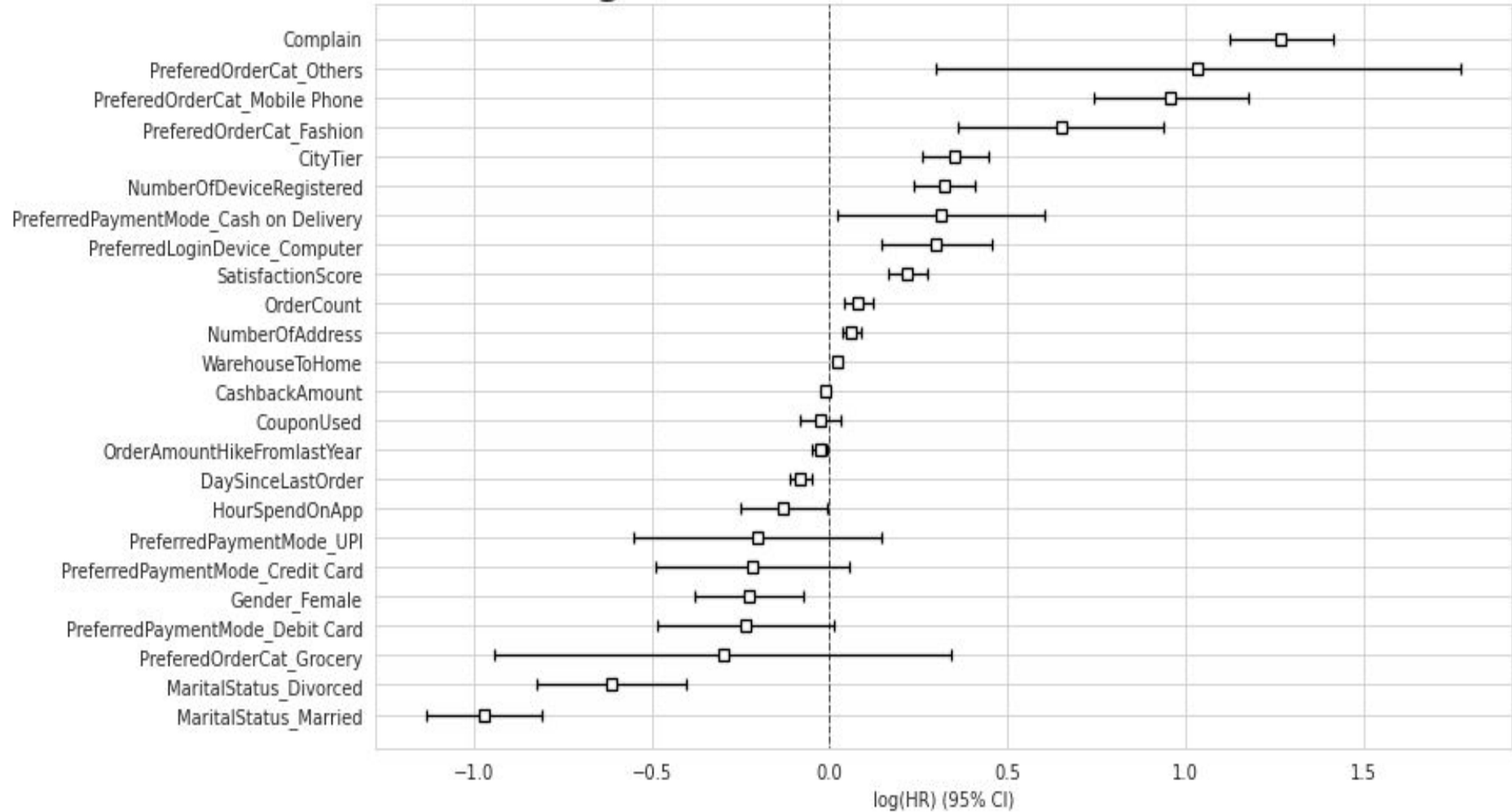Calibration plots (reliability curve) for the 60th month

The Brier Score of our CPH Model is 0.11 at the end of 60 months

- Concordance 0,829 ditafsirkan serupa dengan AUC-ROC regresi logistik

- Brier Score 0.11 at the end of 60 month menandakan bahwa prediksi dari model sampai 60 bulan masih mendekati nilai sebenarnya.

**Insights**

### Survival Regression: Coefficients and Confidence Intervals



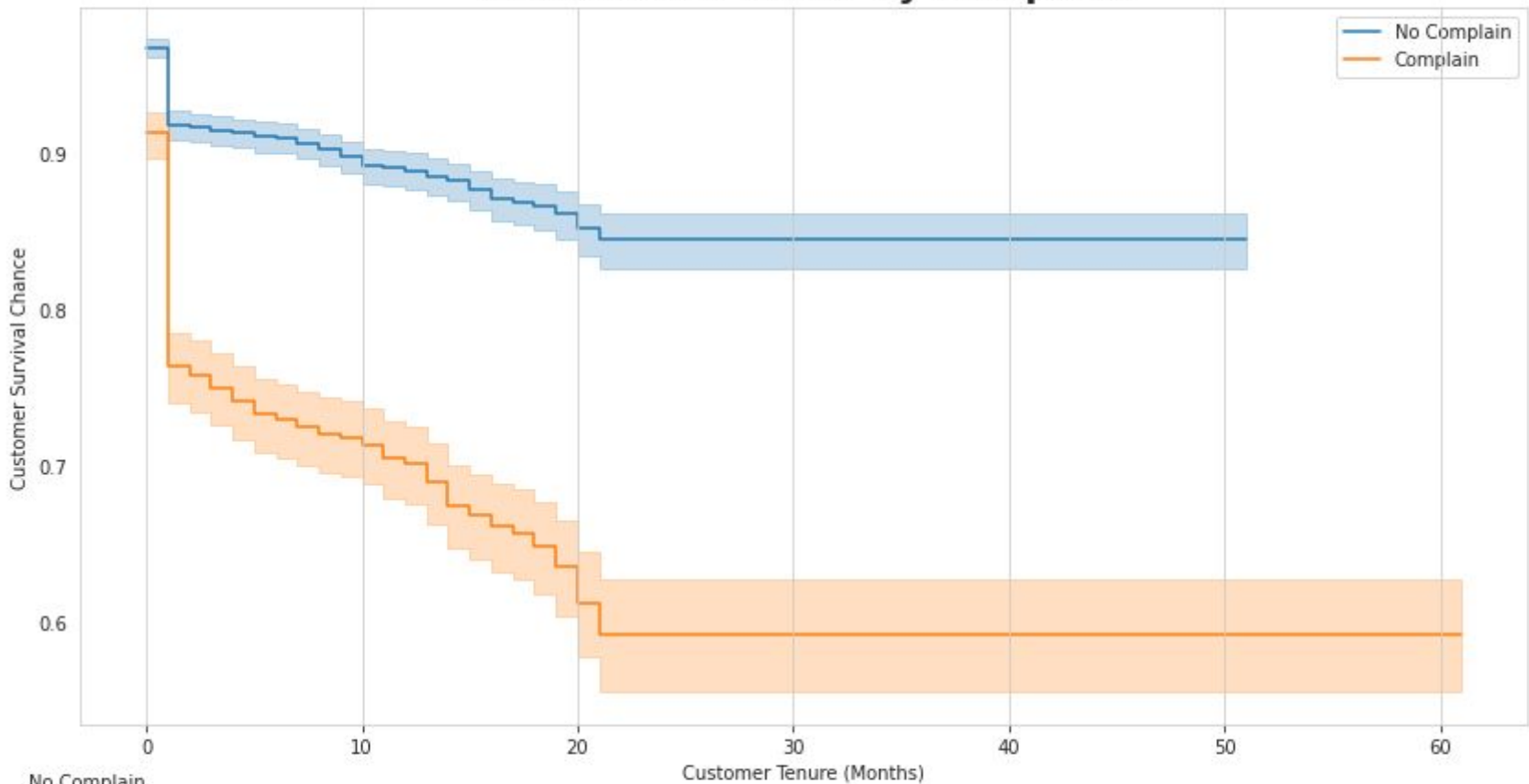Feature yang menghasilkan churn sebagai berikut

- Complain
- Order Category Other
- Order Category Fashion
- Order Category Mobile_Phone

Feature yang menghasilkan retention

- Marital Status Married
- Marital Status Divorced
- Order Category Grocery
- Payment Mode Debit Card

**Insights**



## KM Survival Curve by Complain

- Customer Survival chance with no complain memiliki 89% dan with complain memiliki 68%
- pada 20 bulan pertama customer yang no complain :
  - event (sudah churn) sebesar 388 orang (10%)
  - censored (terindikasi churn tp belum churn) sebesar 2754 orang (75%)
  - at risk(not churn) sebesar 497 orang (15%)

| No Complain | | | | | | |
|---|---|---|---|---|---|---|
| At risk | 3356 | 1460 | 497 | 31 | 2 | 1 | 0 |
| Censored | 165 | 1831 | 2754 | 3216 | 3245 | 3246 | 3247 |
| Events | 118 | 348 | 388 | 392 | 392 | 392 | 392 |

| Complain | | | | | | |
|---|---|---|---|---|---|---|
| At risk | 1260 | 586 | 186 | 19 | 2 | 2 | 1 |
| Censored | 49 | 460 | 805 | 966 | 983 | 983 | 984 |
| Events | 125 | 388 | 443 | 449 | 449 | 449 | 449 |

# Churn Prediction and Prevention

## KM Survival Curve by MaritalStatus



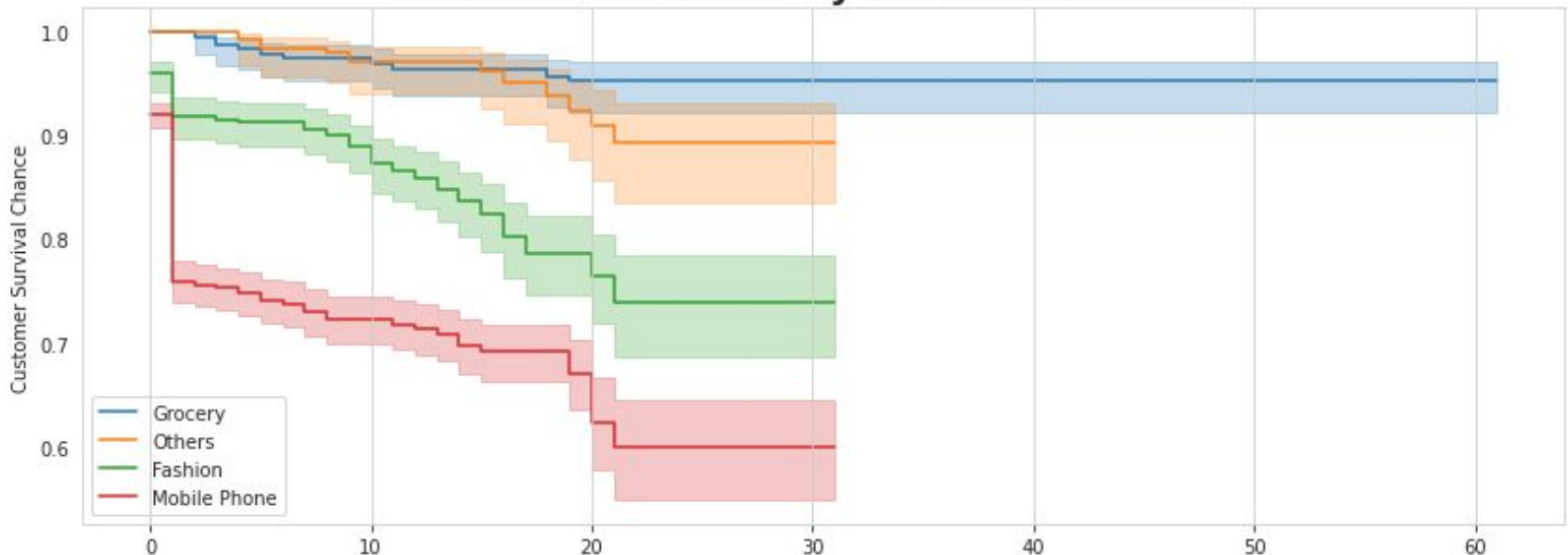| | | | | | | |
|---|---|---|---|---|---|---|
| **Married** | | | | | | |
| At risk | 2494 | 1161 | 377 | 29 | 2 | 2 | 1 |
| Censored | 101 | 1254 | 1997 | 2341 | 2368 | 2368 | 2369 |
| Events | 77 | 257 | 298 | 302 | 302 | 302 | 302 |
| **Divorced** | | | | | | |
| At risk | 778 | 368 | 141 | 10 | 2 | 1 | 0 |
| Censored | 36 | 377 | 585 | 714 | 722 | 723 | 724 |
| Events | 34 | 103 | 122 | 124 | 124 | 124 | 124 |
| **Single** | | | | | | |
| At risk | 1344 | 517 | 165 | 11 | 0 | 0 | 0 |
| Censored | 77 | 660 | 977 | 1127 | 1138 | 1138 | 1138 |
| Events | 132 | 376 | 411 | 415 | 415 | 415 | 415 |

- Customer Survival chance yang Marital Status
  - Married memiliki 88%
  - Divorced memiliki 85%
  - Single memiliki 73%
- pada 20 bulan pertama customer yang Marital Status Married :
  - event (sudah churn) sebesar 377 orang (14%)
  - censored (terindikasi churn tp belum churn) sebesar 1997 orang (75%)
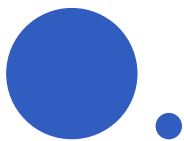  - at risk(not churn) sebesar 298 orang (15%)

# Churn Prediction and Prevention



**KM Survival Curve by PreferedOrderCat**
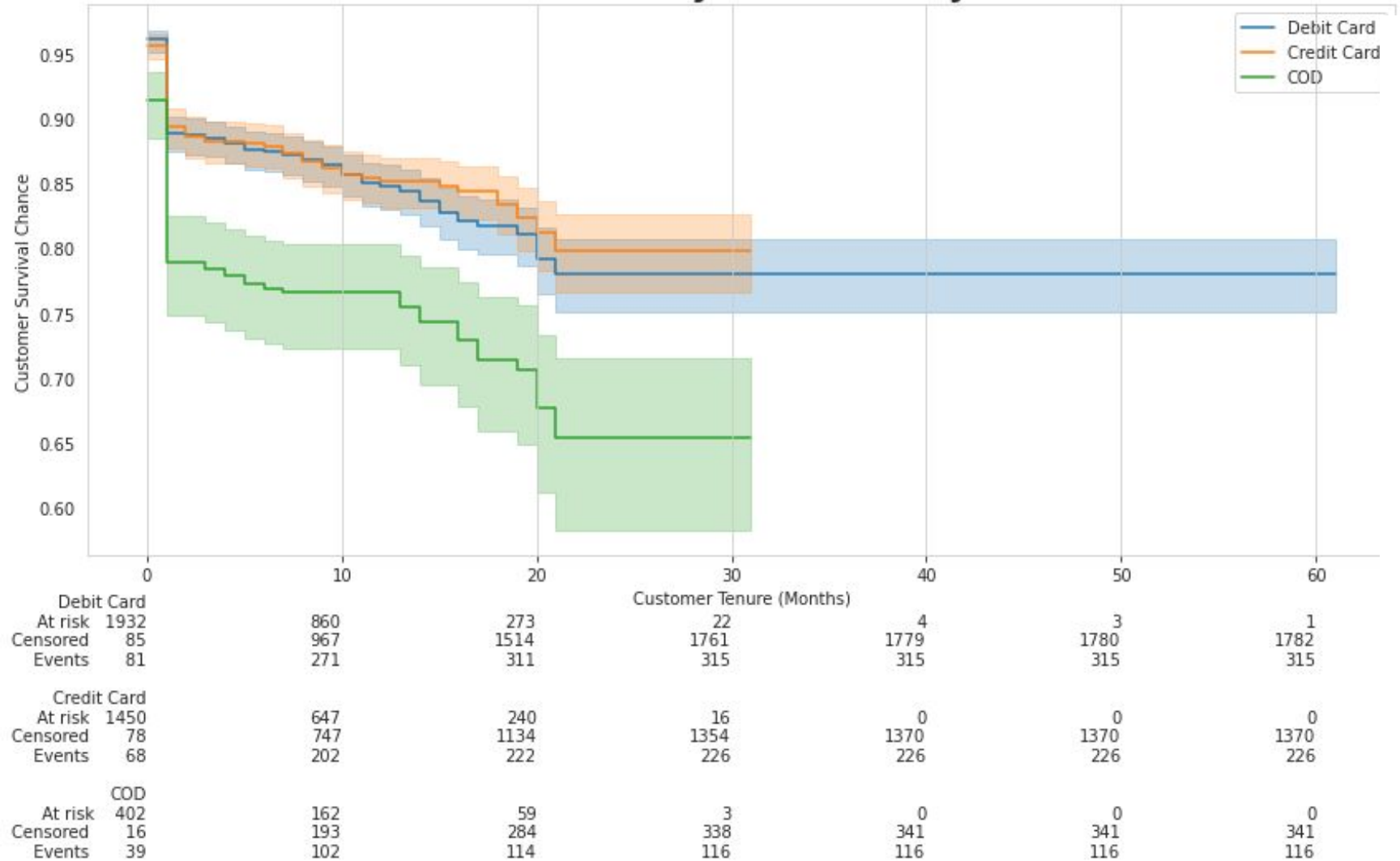
- Customer Survival chance yang Prefered Order Category
  - Grocery memiliki 95%
  - Others memiliki 92%
  - Fashion memiliki 83%
  - Mobile Phone memiliki 73%
- pada 20 bulan pertama customer yang Prefered Order Category Grocery :
  - event (sudah churn) sebesar 16 orang (5%)
  - censored (terindikasi churn tp belum churn) sebesar 151 orang (41%)
  - at risk(not churn) sebesar 199 orang (54%)

## KM Survival Curve by PreferredPaymentMode



- Customer Survival chance yang Prefered Payment Mode
  - Debit Card memiliki 84%
  - Credit Card memiliki 86%
  - COD memiliki 74%
- pada 20 bulan pertama customer yang Prefered Prefered Payment Mode Credit card :
  - event (sudah churn) sebesar 222 orang (14%)
  - censored (terindikasi churn tp belum churn) sebesar 1134 orang (71%)
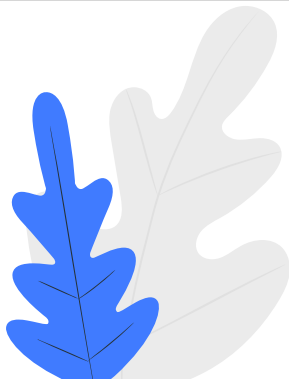  - at risk(not churn) sebesar 240 orang (15%)

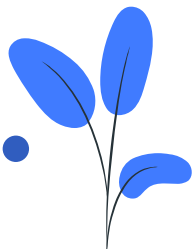| | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Debit Card** | | | | | | | |
| At risk | 1932 | 860 | 273 | 22 | 4 | 3 | 1 |
| Censored | 85 | 967 | 1514 | 1761 | 1779 | 1780 | 1782 |
| Events | 81 | 271 | 311 | 315 | 315 | 315 | 315 |
| **Credit Card** | | | | | | | |
| At risk | 1450 | 647 | 240 | 16 | 0 | 0 | 0 |
| Censored | 78 | 747 | 1134 | 1354 | 1370 | 1370 | 1370 |
| Events | 68 | 202 | 222 | 226 | 226 | 226 | 226 |
| **COD** | | | | | | | |
| At risk | 402 | 162 | 59 | 3 | 0 | 0 | 0 |
| Censored | 16 | 193 | 284 | 338 | 341 | 341 | 341 |
| Events | 39 | 102 | 114 | 116 | 116 | 116 | 116 |

# Calculate Expected Loss & Estimated Revenue Uplift

Let's now drill down a bit more and focus on censored subjects, i.e. those who have not churned yet. We will predict the future survival function of our censored (not churned) customers - the new timeline is the remaining duration of the customer, i.e. normalized back to starting at 0.

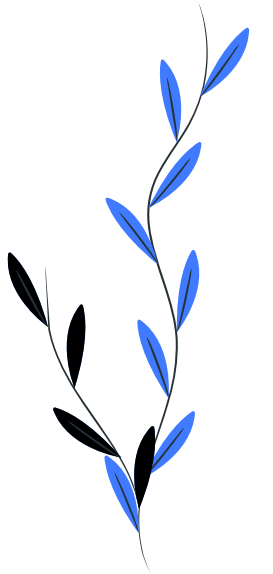| CustomerID | Cashback Amount | Exp_Churn_Month | Exp_Loss | baseline | OrderCat_Grocery_Uplift | PaymentMode_Credit Card_Uplift | PaymentMode_Debit Card_Uplift |
|---|---|---|---|---|---|---|---|
| 50046 | 130.58 | 11.00 | 1,436.38 | 11.00 | 16.00 | 14.00 | 15.00 |
| 50048 | 120.88 | 19.00 | 2,296.72 | 19.00 | 20.00 | 19.00 | 20.00 |
| 50177 | 112.00 | 15.00 | 1,680.00 | 15.00 | 20.00 | 15.00 | 20.00 |
| 50194 | 124.78 | 14.00 | 1,746.92 | 14.00 | 19.00 | 17.00 | 14.00 |
| 50230 | 147.36 | 14.00 | 2,063.04 | 14.00 | 17.00 | 16.00 | 17.00 |

CustomerID 50046 diprediksi akan

- Churn pada bulan ke 11
- Expected Loss sebesar $14,363
- Estimated Revenue Uplift Jika bisa dialirkan ke order category grocery $160 dan payment cc $ 140 atau payment Debit Card $ 150 sehingga akan tidak churn
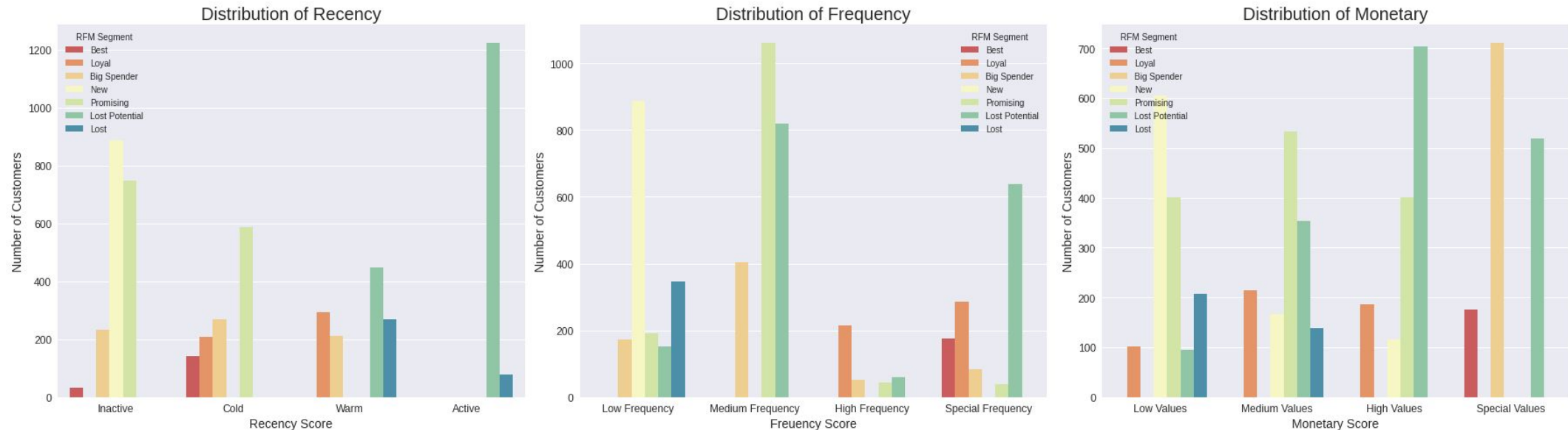
# Segmentation of Customer

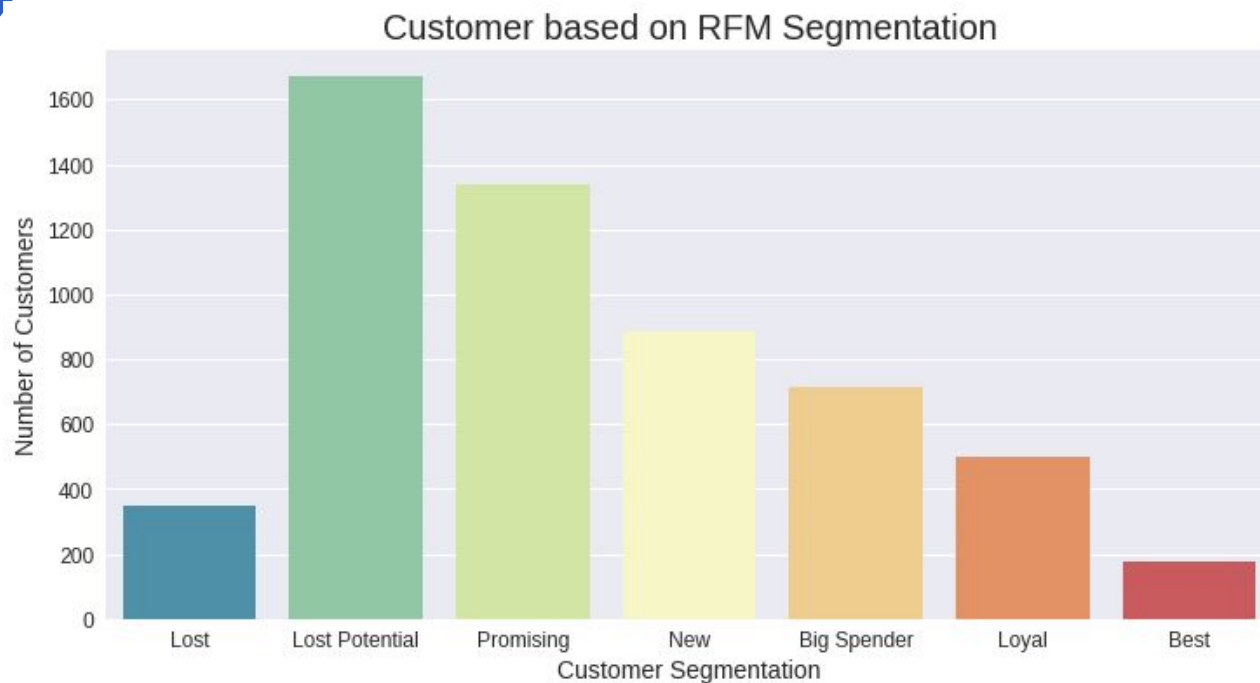Using **RFM Segmentation**, **K-Means**, and **Gaussian**

# RFM Segmentation



Distribution of Recency · Distribution of Frequency · Distribution of Monetary

- Kolom yang digunakan
  - Kolom "DaySinceLastOrder" sebagai "recency"
  - Kolom "OrderCount" sebagai "frequency"
  - Kolom 'CashbackAmount' sebagai 'monetary'

- Kolom "recency" dibagi menjadi 4 segment
  - 'active','warm','cold','innactive'
- Kolom "frequency" dibagi menjadi 4 segment
  - 'special','high','medium','low'
- Kolom 'monetary' dibagi menjadi 4 segment
  - 'low values','medium values','high values','special values'
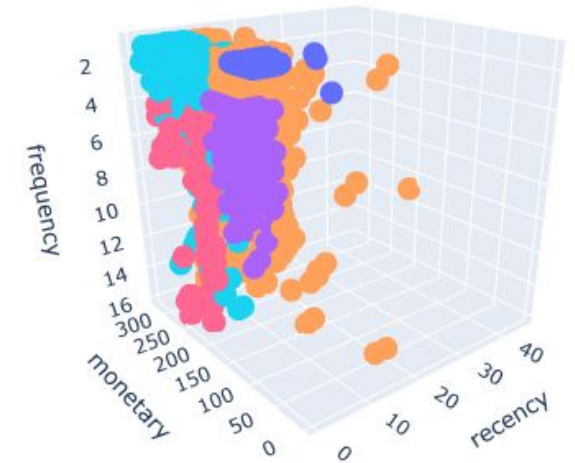
# RFM Segmentation


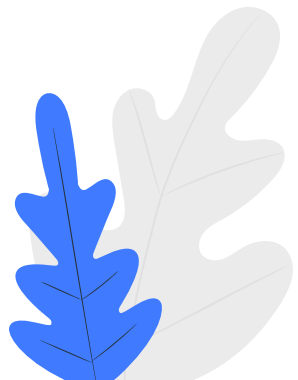
Customer based on RFM Segmentation

- RFM segment berdasarkan score dari distribusi Recency, Frequency, Monetary
- RFM membagi 7 customer segment
  - ['Best', 'Loyal', 'Big Spender', 'New', 'Promising', 'Lost Potential', 'Lost']
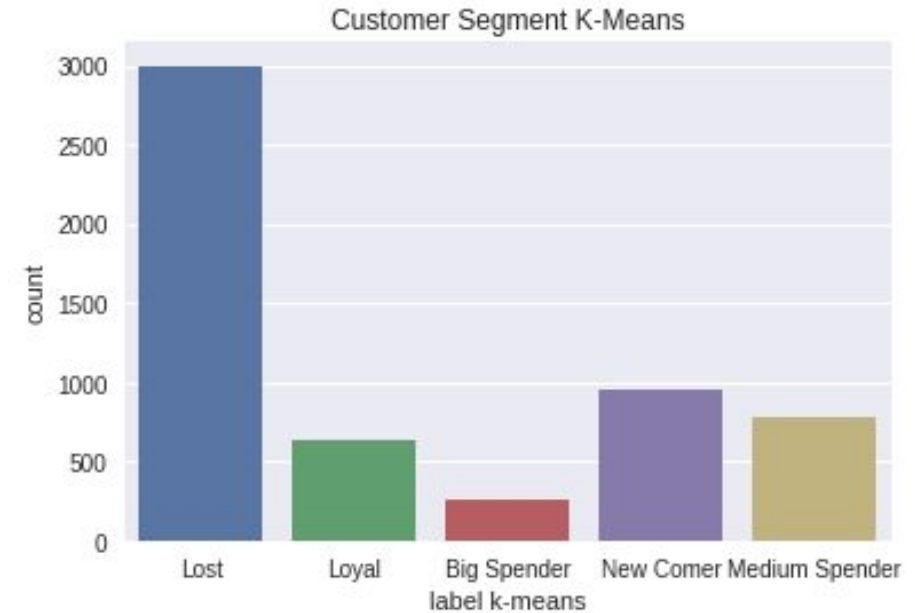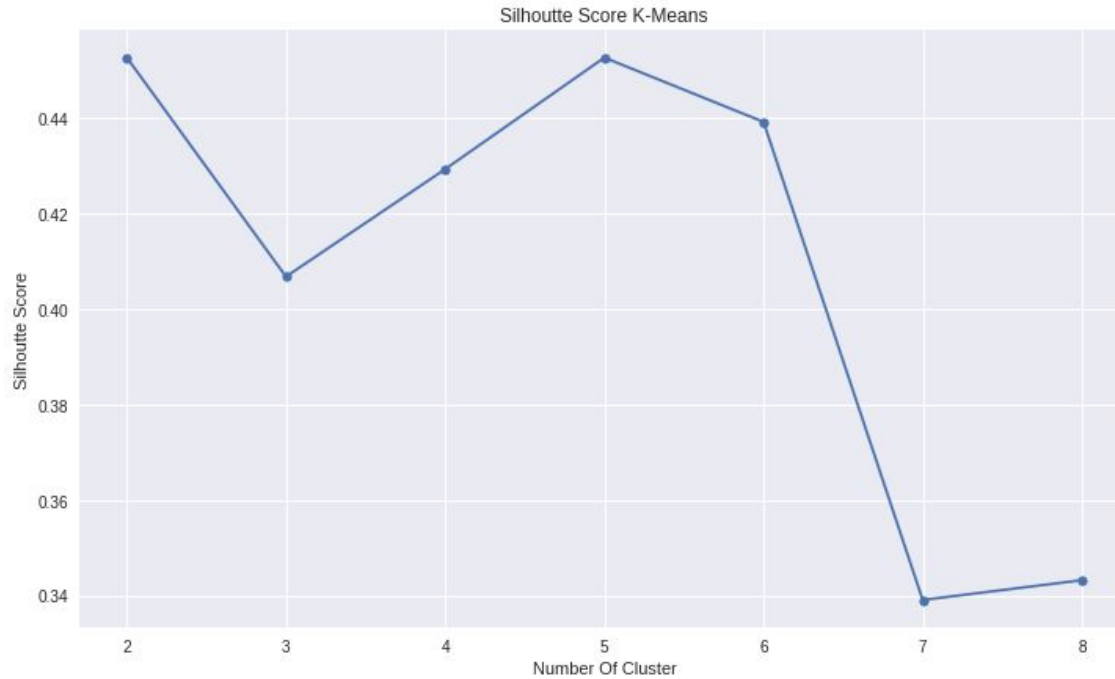
# Insights



Customer Segment by Recency & Frequency

Customer Segment by Recency & Monetary

Customer Segment by Frequency & Monetary

- **Best** : Customer yang melakukan transaksi baru-baru ini, sering melakukan transaksi, dan mempunyai total transaksi yang paling tinggi.
- **Loyal** : Customer yang sudah melakukan transaksi lebih dari 4 kali.
- **Big Spender** : Customer yang melakukan transaksi dengan total transaksi paling tinggi.
- **New** : Customer yang melakukan transaksi baru-baru ini dan baru bertransaksi sebanyak 1 kali.
- **Promising** : Customers yang baru-baru ini melakukan transaksi, serta frekuensi dan total transaksinya diatas rata-rata customers lain.
- **Lost Potential** : Customers yang sudah lama tidak melakukan transaksi, tetapi frekuensi dan total transaksinya diatas rata-rata customers lain.
- **Lost** : Customers yang sudah lama tidak melakukan transaksi, hanya melakukan satu kali transaksi, dan total transaksi sedikit.
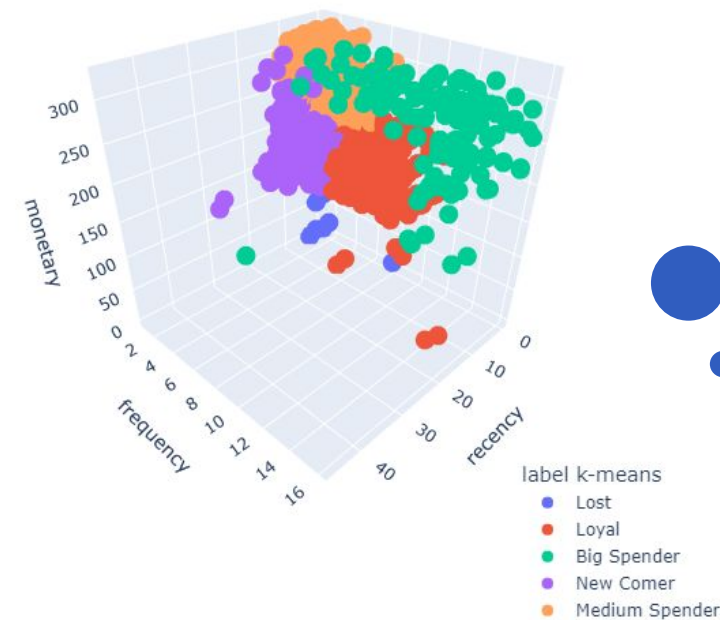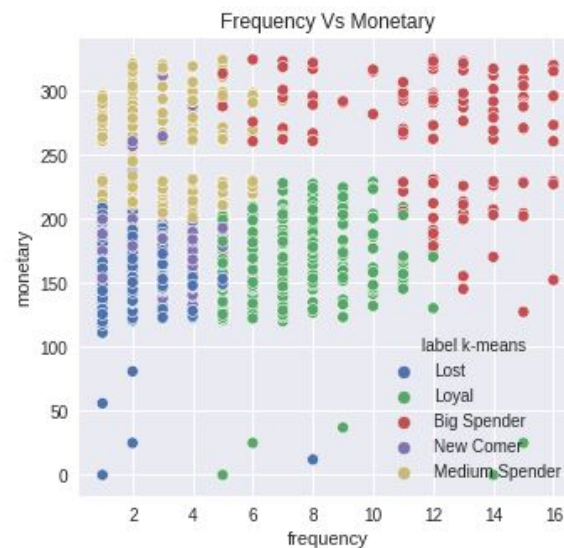
# K-Means



- *Silhoutte Score* terbaik didapatkan pada cluster 2
- Kami memutuskan untuk tidak menggunakan 2 *cluster* karena *cluster* yang terbentuk kemungkinan besar hanya *customer* dengan *frequency* 1 kali dengan *monetary* yang rendah dan *customer* diluar *cluster* tersebut
- Kami akan menggunakan 5 *cluster* karena 5 *cluster* memiliki nilai *silhoutte score* tertinggi setelah 2

# Insights

K-Means

Model 2D Plot



Recency Vs Frequency

Recency Vs Monetary

Frequency Vs Monetary
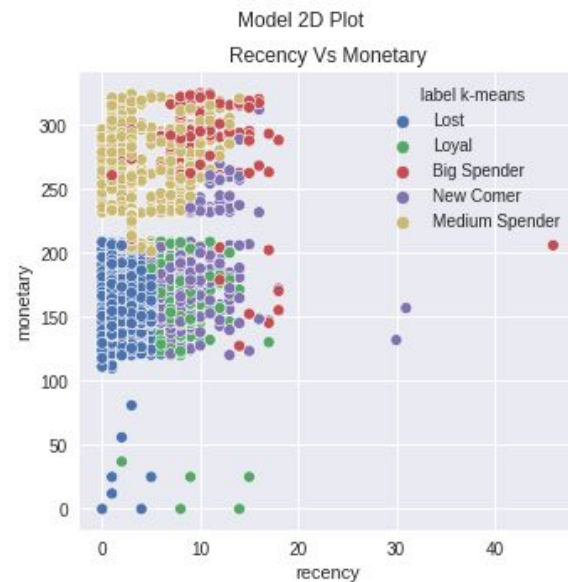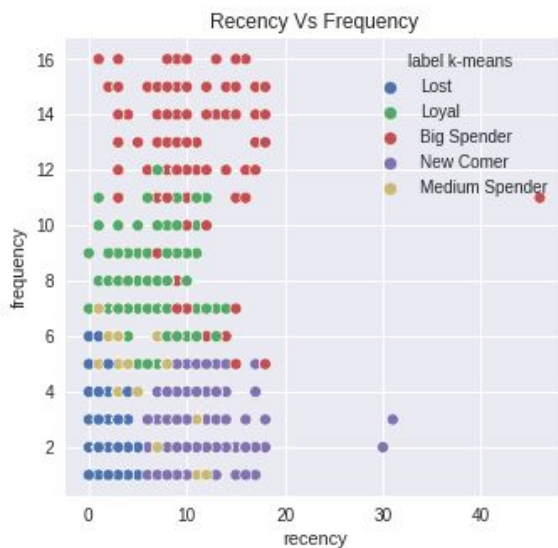
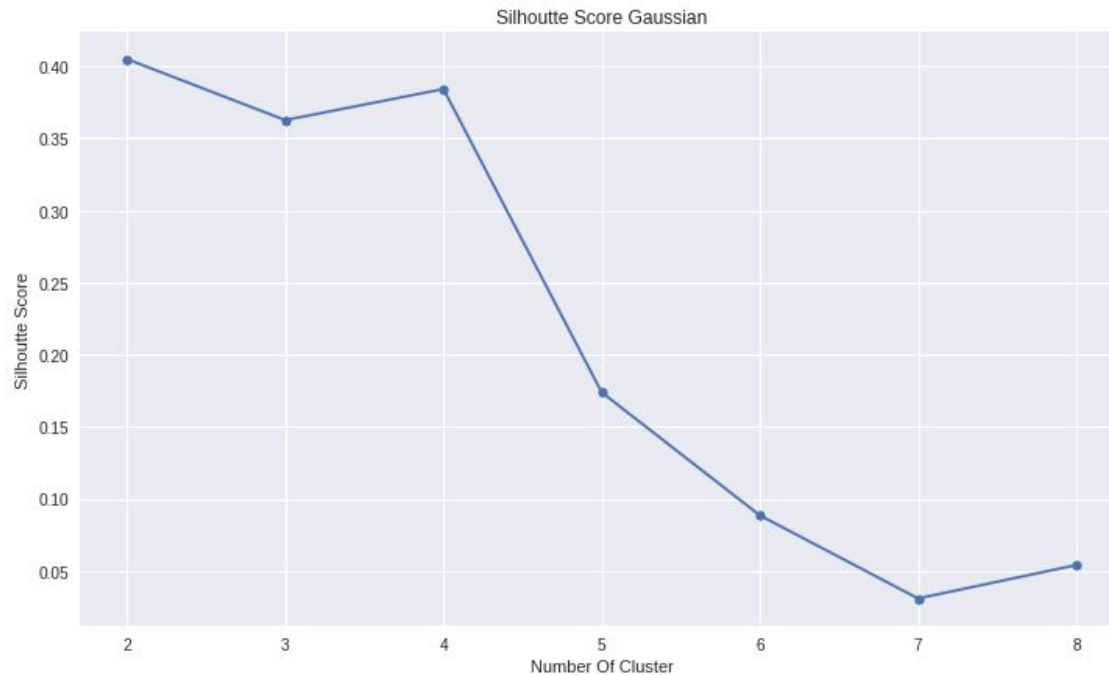# Gaussian



Silhoutte Score Gaussian



Customer Segment Gaussian

- *Silhoutte Score* terbaik didapatkan pada cluster 2
- Kami memutuskan untuk tidak menggunakan 2 *cluster* karena *cluster* yang terbentuk kemungkinan besar hanya *customer* dengan *frequency* 1 kali dengan *monetary* yang rendah dan *customer* diluar *cluster* tersebut
- Kami akan menggunakan 4 *cluster* karena 4 *cluster* memiliki nilai *silhoutte score* tertinggi setelah 2
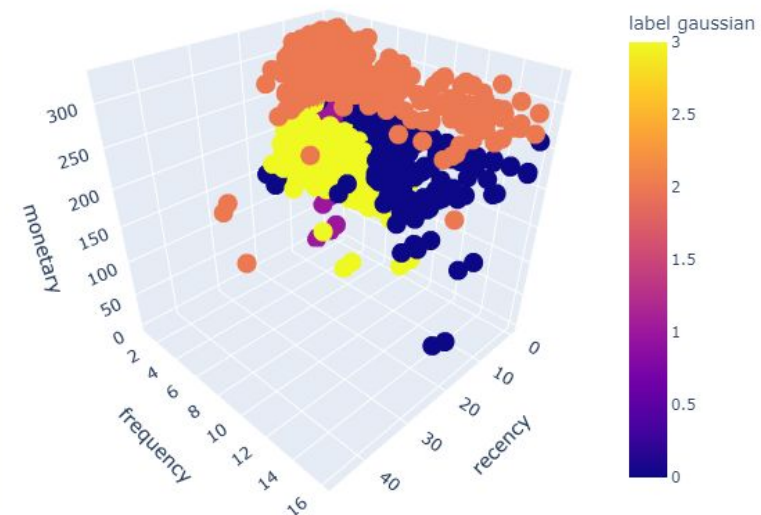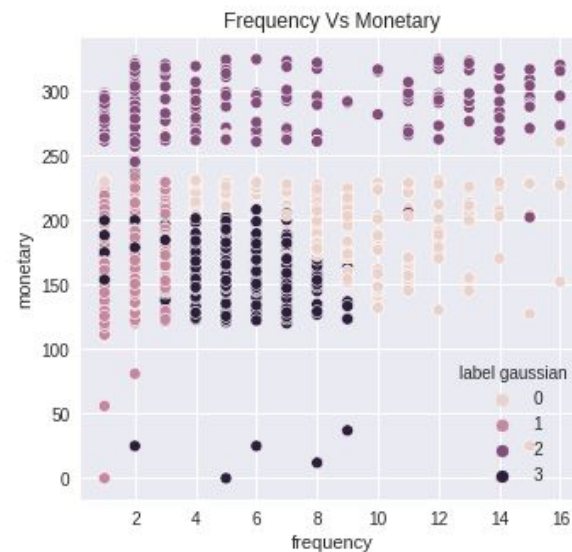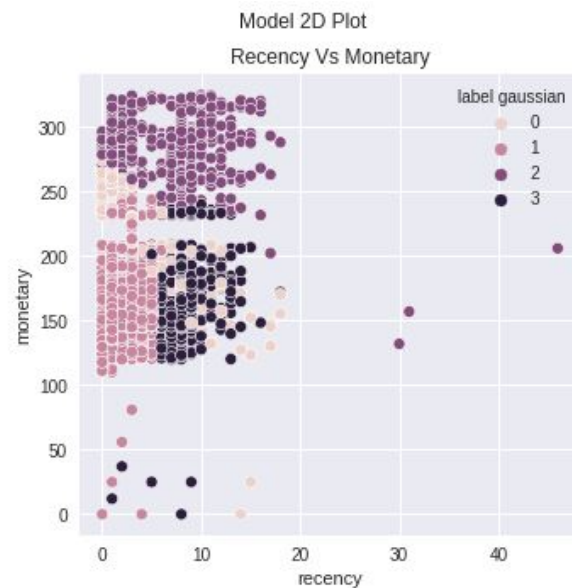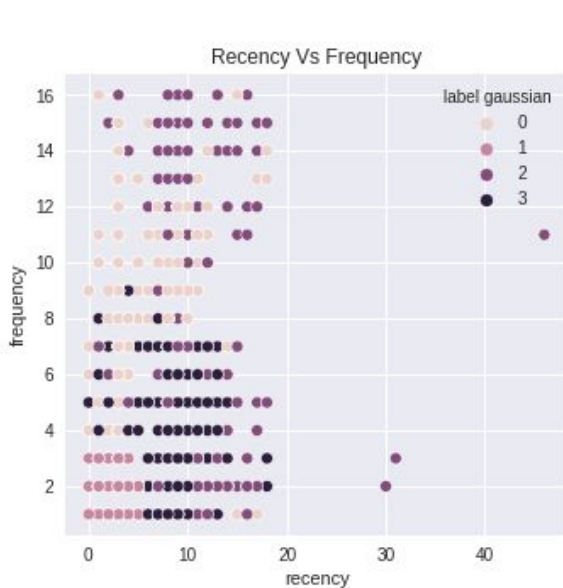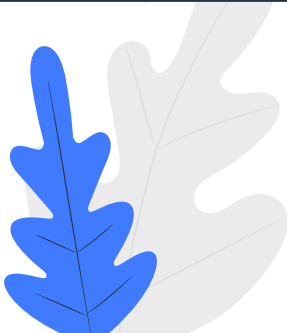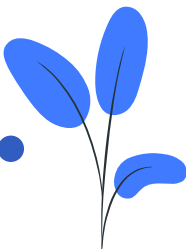
# Insights

## Gaussian



Model 2D Plot

# Summary RFM Segmentation

## Insights

Model RFM Segmentation merupakan model yang memiliki interpretasi paling tinggi dibandingkan model lain & model ini dibuat dengan *domain knowledge* yang kami punya

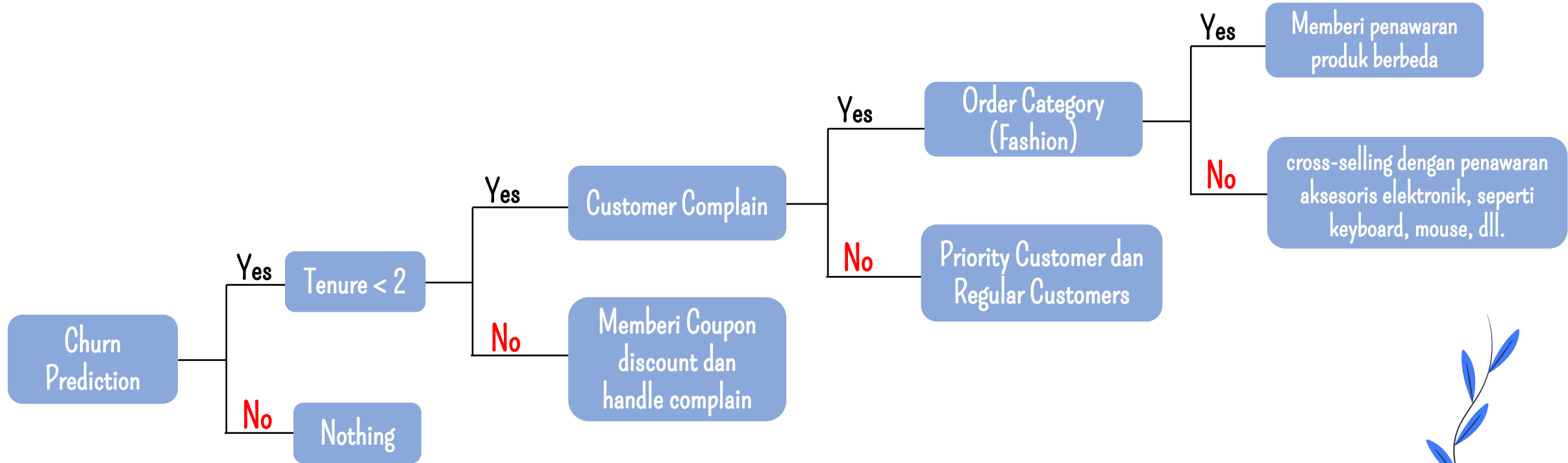| RFM Segment | RFM Segment Score | n customer | mean recency | min recency | max rencency | mean freq | min freq | max freq | mean monetary | min monetary | max monetary | most payment type | avg review score | most product buy |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Best | 7 | 176 | 2.625000 | 0.0 | 3.0 | 8.357955 | 4.0 | 16.0 | 230.968920 | 200.96 | 324.43 | Debit Card | 3.051136 | Fashion |
| Loyal | 6 | 501 | 4.846307 | 3.0 | 7.0 | 4.842315 | 3.0 | 12.0 | 158.280918 | 120.11 | 196.19 | Debit Card | 2.972056 | Laptop & Acc |
| Big Spender | 5 | 712 | 3.200843 | 0.0 | 7.0 | 2.567416 | 1.0 | 15.0 | 244.787219 | 196.67 | 324.26 | Debit Card | 3.005618 | Fashion |
| New | 4 | 888 | 1.010135 | 0.0 | 2.0 | 1.000000 | 1.0 | 1.0 | 138.116137 | 0.00 | 196.10 | Debit Card | 3.087838 | Mobile Phone |
| Promising | 3 | 1336 | 2.079341 | 0.0 | 3.0 | 2.006737 | 1.0 | 9.0 | 153.928451 | 12.00 | 196.37 | Debit Card | 3.058383 | Mobile Phone |
| Lost Potential | 2 | 1671 | 8.461999 | 4.0 | 46.0 | 4.210054 | 1.0 | 16.0 | 195.301556 | 0.00 | 324.99 | Debit Card | 3.115500 | Laptop & Acc |
| Lost | 1 | 346 | 6.132948 | 4.0 | 17.0 | 1.000000 | 1.0 | 1.0 | 141.281647 | 0.00 | 163.22 | Credit Card | 3.080925 | Laptop & Acc |

# Priority Customer Treatment

| RFM Segment | RFM Segment Score | n cus | mean recency | min recency | max rencency | mean freq | min freq | max freq | mean monetary | min monetary | max monetary | most payment type | avg review score | most product buy | sum Exp Loss | sum Grocer Uplift | sum Credit Card Uplift | sum Debit Card Uplift |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Loyal** | 6 | 5 | 3.00000 | 3.0 | 3.0 | 3.000 | 3.0 | 3.0 | 153.3800 | 145.7 | 172.36 | Cash on Delivery | 3.80000 | Mobile Phone | 9740.67 | 600.41 | 1221.61 | 463.76 |
| **New** | 4 | 19 | 1.10526 | 0.0 | 2.0 | 1.000 | 1.0 | 1.0 | 125.2715 | 112.0 | 134.47 | Credit Card | 3.89473 | Mobile Phone | 32903.0 | 1372.46 | 3987.19 | 3389.94 |
| **Promising** | 3 | 23 | 2.30434 | 1.0 | 3.0 | 2.347 | 1.0 | 7.0 | 141.9330 | 120.7 | 159.47 | Cash on Delivery | 3.69565 | Mobile Phone | 48200.0 | 2271.93 | 3944.68 | 3975.56 |
| **Lost Potential** | 2 | 2 | 8.50000 | 8.0 | 9.0 | 5.500 | 5.0 | 6.0 | 12.50000 | 0.0 | 25.00 | E wallet | 2.00000 | Mobile Phone | 225.00 | 0.00 | 25.00 | 25.00 |

| RFM Segment | Strategi |
|---|---|
| Loyal | Loyalty program/reward point dan penawaran barang eksklusif (Cross / Up Selling Strategy) |
| New | Welcome e-mail untuk membangun reletionship, penawaran loyalty program/reward point, dan voucher diskon (Cross / Up Selling Strategy) |
| Promising | Penawaran terbatas secara rutin, voucher diskon dan cashback via e-mail (Retention Strategy) |
| Lost Potential | Penawaran terbatas secara rutin, voucher diskon dan cashback via e-mail (Retention & Reactivate Stretegies) |

Kesimpulan
- Total Expected Loss sebesar $ 910,687
- Estimated Revenue Uplift
  - Order category grocery $42,448
  - Payment Credit Card $ 91,785
  - Payment Debit Card $ 78,543

# CUSTOMER CHURN TREATMENT

# Summary & Recommendations

Dari data visualisasi diperoleh churn ratio memiliki korelasi tenure, complain, cashback Amount, & preferedordercat

Hasil predict churn sangat dipengaruhi oleh tinggi rendahnya Tenure, Complain, Number of Address dan cashback Amount

Hasil Survival Analysis, customer memiliki survival chance terbesar pada No Complain, Marital Status Married, Payment Mode Credit Card, Order Category Grocery

Hasil RFM Segmentation menunjukkan priority customer treatment pada segment Loyal, New, Promising, dan Lost Potential

- Total Expected Loss sebesar $ 910,687
- Estimated Revenue Uplift
  - Order category grocery $42,448
  - Payment Credit Card $ 91,785
  - Payment Debit Card $ 78,543

# Daftar Pustaka

Agrawal, A., Gans, J., & Goldfarb, A. (2018). *Prediction machines: the simple economics of artificial intelligence*. Harvard Business Press.

Al-Sahaf, H., Bi, Y., Chen, Q., Lensen, A., Mei, Y., Sun, Y., Tran, B., Xue, B., & Zhang, M. (2019). A survey on evolutionary machine learning. *Journal of the Royal Society of New Zealand*, *49*(2), 205–228. https://doi.org/10.1080/03036758.2019.1609052

Apte, C. (2010). Invited Applications Paper: The Role of Machine Learning in Business Optimization. *Proceedings of the 27th International Conference on Machine Learning, Haifa, Israel, 2010.*

Brynjolfsson, E., & McAfee, A. (2014). *The second machine age: Work, progress, and prosperity in a time of brilliant technologies*. WW Norton & Company.

Masarifoglu, M., & Buyuklu, A. H. (2019). *Applying Survival Analysis to Telecom Churn Data* (pp. 261–275). American Journal of Theoretical and Applied Statistics.

Pinem, R. J., Afrizal, T., & Saputra, J. (2020). The Relationship of Cashback, Discount, and Voucher toward Decision to Use Digital Payment in Indonesia. *Talent Development & Excellent*, *12*(3s), 2766–2774.

Wu, X., & Meng, S. (2016). E-commerce Customer Churn Prediction Based on. *2016 13th International Conference on Service Systems and Service Management (ICSSSM)*, 1–5.

https://www.kaggle.com/ankitverma2010/ecommerce-customer-churn-analysis-and-prediction