Mutual Fund Investment Startup ~Python for Data Analytics~

• • •

Fathan Abdallah





List of Contents:

- 1. Milestone 1 (Data Cleaning and EDA)
- 2. Milestone 2 (Segmentation and Clustering Analysis)
- 3. Milestone 3 (Python Logistic Regression Model)

Python Milestone 1

(Data Cleaning and EDA)



Analytical Objective

- 1. There are 4 types of investing such as saham, pasar uang, pendapatan tetap and campuran so customers can choose among those 4 options
- 2. Because of it is investing company so we gain profit from customers who doing transactIts investing dataset (mutual fund), so we can analyze according to the dataset and macro/micro economic
- 3. Profit comes from transactions fee (buy and sell fee) or profit sharing if gaining capital gain

Data Cleaning Users

Copy

Checking Data Types

Treat Missing / Irrelevant Values

Check for Values & Typo

- Make a copy as df_users_dc
- Melihan ukuran row x kolom dengan .shape()
- Melihat informasi keseluruhan dengan .info()
- Mencari nilai blank
- Convert kolom data (object) menjadi datetime
- Mengubah user_id menjadi string karena tidak diperlukan dalam proses operasi
- Drop user_occupation and user_income_source becaus irrelevant
- In referral_code_used replace blank values with unknown
- Check values each columns

Data Cleaning Users

Data Manipulation

Checking Duplicates

- Clustering for each income range into 3 categories to make easier to identify
- No duplicates found

Data Cleaning Daily Transaction

Copy

Checking Data Types

Treat Missing / Irrelevant Values

- Make a copy as df_DT_dc
- Melihat ukuran row x kolom dengan .shape()
- Melihat informasi keseluruhan dengan .info()
- Mencari nilai blank
- Convert kolom data (object) menjadi datetime
- Mengubah user_id menjadi string karena tidak diperlukan dalam proses operasi
- There are many missing values in the dataset. But it doesnt mean we can just remove it. 0 or blank values means there is no transaction not there is no data. So, yes I consider it as data
- The information we can take from zero/blank data is the frequency of the transaction, so 0/blank values will be treated in data manipulation part

Data Cleaning Daily Transaction

Data Manipulation

- From daily transaction dataset we can take some information such number of frequencies an user balanced
- In this part I split "Daily transaction" into 2 separated dataset
- Those are frequency dataset and last balanced of user_id
 Find the Frequencies
- Make a copy from "df_DT_dc" as "df_DT_freq"
- Find the number of frequencies by using .count() function
- add frequencies between buy and sell, also drop the columns amount because we just got the information we needed
 Find User Balance
- Make a copy from "df_DT_dc" as "df_DT_balance"
- Sort by user_id and date to find out the last balance
- Rename column to make it compact
- Drop duplicates and keep the last one (after sorted) to get the last balance in the last date
- Merge "Frequencies Table" and "Balance Table" as "df_daily_tran"
- Merge "df_users" and "df_daily_tran" as "df_merged"

- 1. Make a copy from df.merged as df_eda
- 2. Analyze the characteristic of the table with .info() function
- 3. Descriptive analysis about numeric information

4. Descriptive analysis about string variable

```
objects = ['user id',
 'user gender',
 'user income range',
 'referral code used']
df eda[objects].describe()
         user id user gender user income range referral code used
            8277
                         8277
                                             8277
                                                                 8277
 count
unique
            8277
                                                3
                                                                    2
  top
         3816789
                         Male
                                            silver
                                                              unknown
  freq
                         5176
                                             6233
                                                                 5322
```

5. Descriptive analysis about date type variable

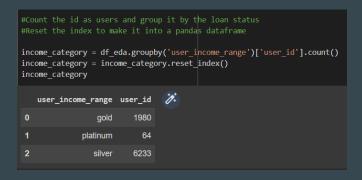
```
df_eda['registration_import_datetime'].describe()
```

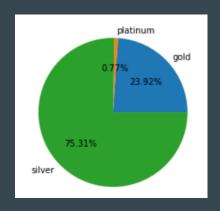
6. Number of Customer

```
[90] #Number of customer
    df_eda['user_id'].count()

8277
```

7. User_income_range





8. Transaction Frequency

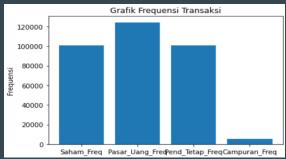
```
Saham_Freq = df_eda['Saham_Freq'].sum()
Pasar_Uang_Freq = df_eda['Pasar_Uang_Freq'].sum()
Pend_Tetap_Freq = df_eda['Pend_Tetap_Freq'].sum()
Campuran_Freq = df_eda['Campuran_Freq'].sum()

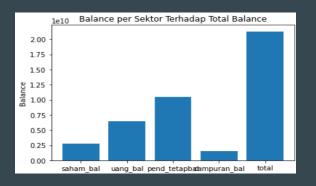
print(f"saham freq: {Saham_Freq}")
print(f"uang freq: {Pasar_Uang_Freq}")
print(f"pendapatan tetap freq: {Pend_Tetap_Freq}")
print(f"campuran freq: {Campuran_Freq}")
```

9. Sector Balance to Total Balance

```
saham_bal = df_eda['saham_bal'].sum()
uang_bal = df_eda['uang_bal'].sum()
pend_tetapbal = df_eda['pend_tetapbal'].sum()
campuran_bal = df_eda['campuran_bal'].sum()
total = df_eda['total_balance'].sum()

print(f"saham balance: {saham_bal}")
print(f"uang balance: {uang_bal}")
print(f"pendapatan tetap balance: {pend_tetapbal}")
print(f"campuran balance: {campuran_bal}")
print(f"total: {total}")
```



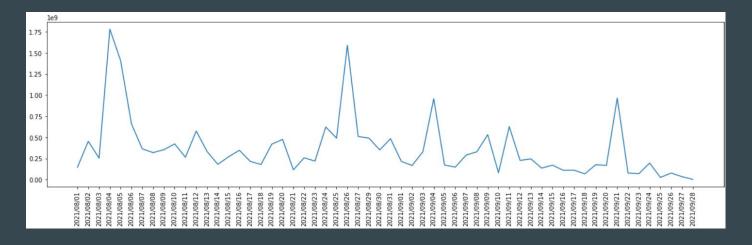


8. Trend

```
plt.figure(figsize =(20, 5))

plt.xticks(rotation = 90)
plt.plot(transaction_trend['trans_date'], transaction_trend['total_balance'])

plt.show()
```



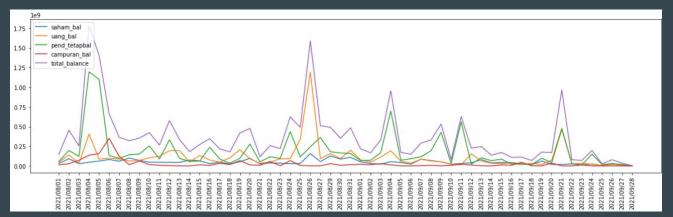
9. Trend Each Product

```
#Create line chart
plt.figure(figsize = (20, 5))
axis = transaction_trend2.columns.tolist()
#axis.head()

for x in axis[1:]: #
   plt.plot(transaction_trend2['trans_date'], transaction_trend2[x])

plt.xticks(rotation = 90)
plt.legend(transaction_trend2.iloc[:,1:],loc = 2)

plt.show()
```



Python Milestone 2

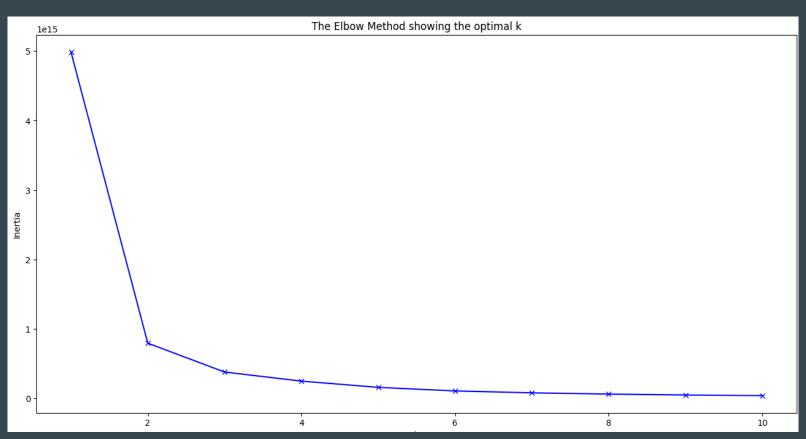
(Segmentation and Clustering Analysis)

000

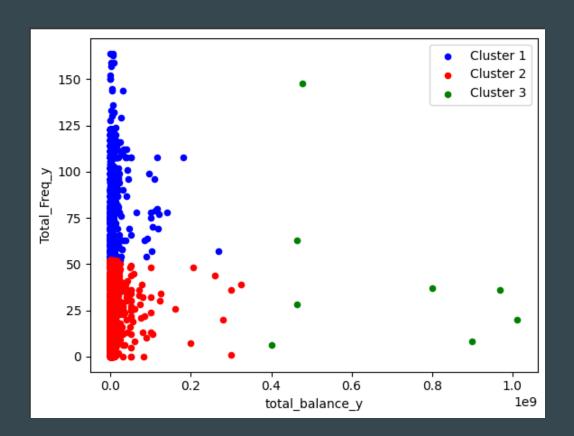
Related Variable

- In this case I want to make cluster according Total_Freq and total_balance
- 1. Total_Freq is choosen because as fund managers, we gain profit from buy/sell transaction fee. So more often costumers doing transactions, we gain more profit
- 2. total_balance reflect the buying power of each customers.
 - Some points we have to notice are it is normal when customers dont have either balance or frequencies in their account. Because investors divided into some classes
- 1. Short term investors = Usually high frequencies, prefer cash in the end of month/quarter (0 balance at the end of the month)
- 2. Long term investors = Low frequencies, making deep analysis of the market before enter positions. As soon as they enter the market, they keep their positions for years at least 1 year

Elbow Method



Clustering



Segmentation

Х	0	1	2
Gender	Male	Male	Female
Age	26	27	39
Income Range	Silver	Silver	Gold
Total Frequency	20	83	43
Total Balance	1.716.676,29	2.349.707,50	685.778.338,25
Total Transaction	916.760,04	2.348.775,22	455.228.150,50
Saham Balance	444.725,37	587.351,74	6.500.000
Pasar Uang Balance	726.924,01	643.763,58	435.568.988,75
Pendapatan Tetap Balance	1.377.780,66	1.069.620,35	527.564.393
Campuran Balance	9.171.888,25	2.890.038,29	12.666.666,60













Learner investors,

Contribute 0,1% of total transaction

Experience investors,

high frequency of transaction. Contribute 0,5% of total transaction

Big money, long term investors, give 99% contribution of total transaction

Cluster 0 (Learner Investors)

X	0
Gender	Male
Age	26
Income Range	Silver
Total Frequency	20
Total Balance	1.716.676,29
Total Transaction	916.760,04
Saham Balance	444.725,37
Pasar Uang Balance	726.924,01
Pendapatan Tetap Balance	1.377.780,66
Campuran Balance	9.171.888,25

Recommendations:

- 1. Give more promotions to gain more transactions or top up on their accounts
- 2. Give education about investing through short videos or articles
- 3. Give promotions with referral because this cluster is huge in numbers

Cluster 1 (Experience Investors)

Х	1
Gender	Male
Age	27
Income Range	Silver
Total Frequency	83
Total Balance	2.349.707,50
Total Transaction	2.348.775,22
Saham Balance	587.351,74
Pasar Uang Balance	643.763,58
Pendapatan Tetap Balance	1.069.620,35
Campuran Balance	2.890.038,29

Recommendations:

- I. Give advance education about retirement income from investing
- 2. Give promotion or coupon from referral because this cluster is huge in numbers

Cluster 2 (Big Money)

х	2
Gender	Female
Age	39
Income Range	Gold
Total Frequency	43
Total Balance	685.778.338,25
Total Transaction	455.228.150,50
Saham Balance	6.500.000
Pasar Uang Balance	435.568.988,75
Pendapatan Tetap Balance	527.564.393
Campuran Balance	12.666.666,60
Campuran Balance	12.666.666,60

Recommendations:

- l. Give transaction fee discount
- 2. Low profit sharing to fund managers

Python Milestone 3

(Python Logistic Regression Model)

000

Understanding Business Problem

- 1. Mutual Investment Company has some investment products: stocks, bond, money market, and mixed investment mutual funds
- 2. Company has some issues such as churn rate, and limitation of campaign budget (30%) and it give impact to the company profit
- 3. According to previous data, cost per campaign is Rp1.000 and the transaction fee (buy and sell) is 0,15%



Steps





Code Preparation

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn import cluster
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix
```

Raw Data Preparation

3651808 2021-09-12 13:05:04 Male 52 silver no referral 0 11 11 3934978 2021-08-30 14:26:12 Male 21 silver no referral 238604 20 20 3944498 2021-08-31 20:37:09 Male 17 silver no referral 160000 19 19 3980156 2021-09-05 07:08:38 Male 24 silver used referral 10000 16 16	user_id	registration_import_datetime	user_gender	user_age	user_income_range	referral_code_used	Total_Transaction	Saham_Freq	Pasar_Uang_Freq	Pend_Tetap_Freq
3944498 2021-08-31 20:37:09 Male 17 silver no referral 160000 19 19	3651808	2021-09-12 13:05:04	Male	52	silver	no referral	0	11	11	11
	3934978	2021-08-30 14:26:12	Male	21	silver	no referral	238604	20	20	20
3980156 2021-09-05 07:08:38 Male 24 silver used referral 10000 16 16	3944498	2021-08-31 20:37:09	Male	17	silver	no referral	160000	19	19	19
	3980156	2021-09-05 07:08:38	Male	24	silver	used referral	10000	16	16	16

Campuran_Freq	Total_Freq	date	saham_bal	uang_bal	pend_tetapbal	campuran_bal	total_balance	churn
0	33	2021- 09-30	30000.000	10000.000	60000.000	0.000	100000	1
0	60	2021- 09-30	800.000	87804.000	110000.000	0.000	198604	1
0	57	2021- 09-30	0.000	0.000	40000.000	0.000	40000	0
0	48	2021- 09-30	0.000	0.000	0.000	0.000	0	0

Category of each variable



EDA

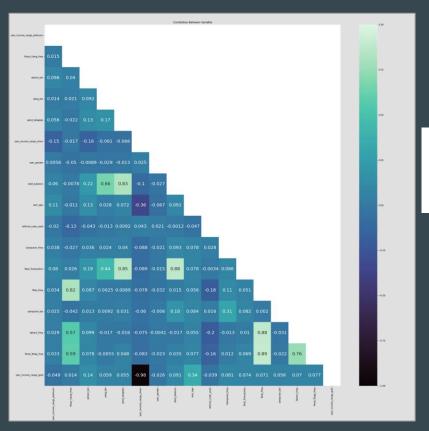
```
user_id
churn
          4809
0
           3468
<AxesSubplot:xlabel='churn', ylabel='user_id'>
   5000
   4000
   3000
   2000
   1000
                   0
                             churn
```

churn rate

3468/(4809+3468) = 0,418



Uji Korelasi dan Contoh





Drop High Correlation Var

```
# Create correlation matrix
corr_matrix = df_churn_new.corr().abs()

# Select upper triangle of correlation matrix
upper = corr_matrix.where(np.triu(np.ones(corr_matrix.shape), k=1).astype(np.bool))

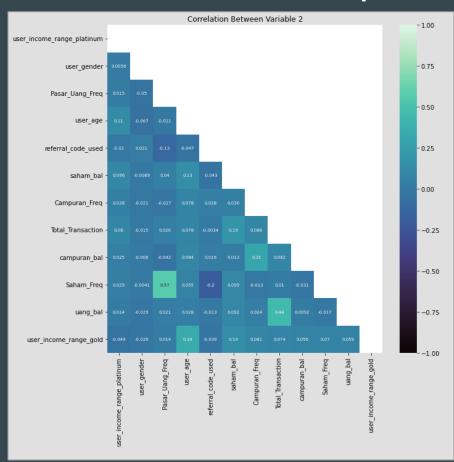
# Find features with correlation greater than 0.7, and add into to_drop list
to_drop = [column for column in upper.columns if any(upper[column] > 0.7)]
to_drop
```

```
['Pend_Tetap_Freq',
 'Total_Freq',
 'pend_tetapbal',
 'total_balance',
 'user_income_range_silver']
```



Dropped columns

New Correlation After Drop Column





Fit Logistic Regression - Accuracy Rate

```
model = LogisticRegression(class_weight='balanced',max_iter=500)
model.fit(x_training, y_training)
```

```
▼ LogisticRegression
LogisticRegression(class_weight='balanced', max_iter=500)
```

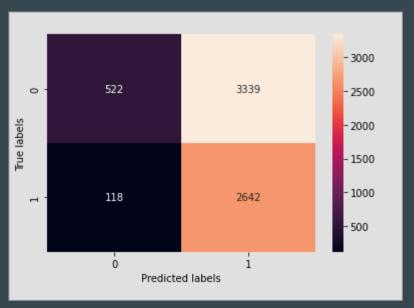
```
# Accuracy dr prediksi model dengan data training model.score(x_training, y_training)
```

0.4778734330161607

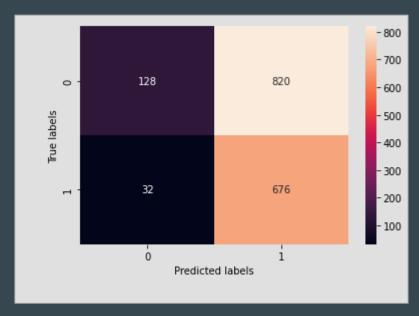




Confusion Matrix & Performance Test



	precision	recall	f1-score	support
Not Churn	0.82	0.14	0.23	3861
Churn	0.44	0.96	0.60	2760
accuracy			0.48	6621
macro avg	0.63	0.55	0.42	6621
weighted avg	0.66	0.48	0.39	6621



	precision	recall	f1-score	support
Not Churn Churn	0.80 0.45	0.14 0.95	0.23 0.61	948 708
accuracy macro avg weighted avg	0.63 0.65	0.54 0.49	0.49 0.42 0.39	1656 1656 1656

Benefit Cost Analysis (Transaction Fee)

prft70 customer.sum()

```
user_id Total_Transaction fee_trans
                                                fee_trans = 0.15% x Total_Transaction
      3816789
                           100000
                                      150,000
                                                           = Rp 22.001.245,962
      3802293
                          8500000
                                   12750.000
      3049927
                           149000
                                     223.500
                                                 user id
                                                                  3802293378330237226103787202383709937402903768...
                                                                                                14667407209
  3
      3836491
                                        0.000
                                                 fee trans
      3783302
                          2889569
                                    4334.354
                                                 dtype: object
# Get 70% threshold of profit
prft70 threshold = df tot trans['fee trans'].quantile(0.7)
                                                                                        fee trans
# Get top 70% customers
prft70_customer = df_tot_trans[df_tot_trans['fee_trans'] > prft70_threshold]
                                                                                        Quantile 70%
# Merge with actual churn outcome
prft70_customer = prft70_customer.merge(y, left_index = True, right_index = True)
```

Benefit Cost Analysis (Churn Customers)

	predicted_score	churn	Flag
0	0.532	0	Medium Churn
1	0.983	0	Mostly Churn
5	0.527	1	Medium Churn
8	0.519	1	Medium Churn
10	0.670	1	Mostly Churn

```
pct70_customer['churn'].value_counts()

1 1603
0 880
Name: cnurn, dtype: int64
```

```
def threshold_category(pct70_customer):
    if (pct70_customer['predicted_score'] <= 0.3 ):
        return 0
    elif (pct70_customer['predicted_score'] > 0.3) and (pct70_customer['predicted_score'] <= 0.6):
        return 'Medium Churn'
    else:
        return 'Mostly Churn'

pct70_customer['Flag'] = pct70_customer.apply(threshold_category, axis = 1)
pct70_customer</pre>
```

Potential Benefit per Campaign

If case transaction fee = 0,15%

Cost campaign = Rp 1000

Churn users = 1603

Potential Benefit = 22.001.245,962 - (Rp1000 x 1603)

= Rp 20.398.245,962

THANKS FOR YOUR ATTENTION