

Phân khúc khách hàng bằng RFM và Kmean

Từ những khách hàng tiêu dùng lớn cho đến những khách hàng rời bỏ doanh nghiệp, tất cả những khách hàng đều có nhu cầu và mong muốn đa dạng. Doanh nghiệp muốn khách hàng chi tiêu nhiều hơn từ những chiến dịch tiếp thị chương trình, sản phẩm mới tới khách hàng theo những cách khác nhau. Tuy nhiên, câu hỏi đặt ra là làm thế nào để đưa ra được các chiến dịch tiếp thị phù hợp với những nhóm khách hàng đang có nhu cầu để từ đó tăng tỷ lệ phản hồi từ khách hàng và từ đó tăng doanh số bán hàng. Bài toán đặt ra là làm thế nào để có thể phân khúc khách hàng một cách tương đối chính xác dựa trên “hành vi giao dịch lịch sử” của khách hàng, thuật toán RFM sẽ giúp chúng ta giải quyết vấn đề này một cách nhanh chóng và hiệu quả.

1. RFM (Recently – Frequently – Monterey Value)

Để phân khúc khách hàng, có một vài metrics thường sử dụng như là:

- + Khách hàng mua hàng lần cuối khi nào? (R)
- + Tần suất khách hàng mua hàng nhiều hay ít (F)
- + Khách hàng chi bao nhiêu tiền để mua sản phẩm của công ty (M)

Recently: Để tính ngày mua hàng gần nhất: ta lấy ngày hiện tại trừ đi ngày mua hàng gần nhất của khách hàng, ta sẽ được một giá trị gọi là khoảng cách từ ngày mua hàng gần nhất đến hiện tại của khách hàng. Số này càng nhỏ thì khách hàng càng mua hàng gần hơn

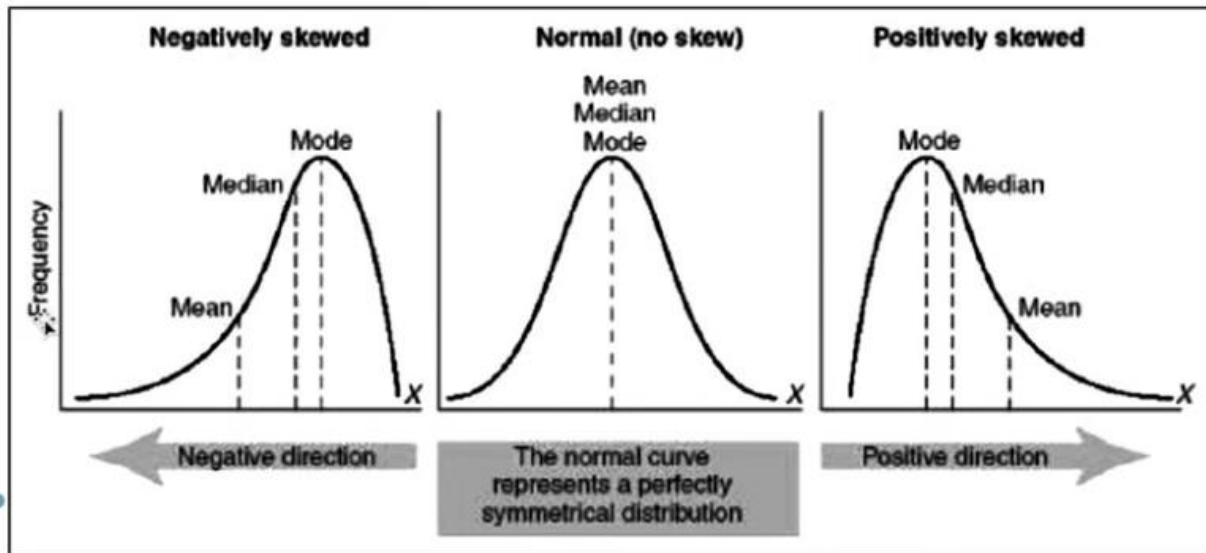
Frequently: Đếm số lượt mua hàng của khách hàng trong khoảng thời gian ta trích xuất dữ liệu.

Montery Value: tính tổng tiền tất cả mỗi khách hàng. (nhân đơn giá với số lượng và tính tổng)

Từ bảng hóa đơn lịch sử mua sắm khách hàng, ta tính các thông số R, F, M này cho mỗi khách hàng. Coi đó là các Feature để phân nhóm khách hàng.

2. Chỉ số Skewness

Chỉ số Skewness dùng để đánh giá dữ liệu có phân phối chuẩn hay không hay bị lệch trái, lệch phải.



Lý tưởng: chỉ số Skewness = 0. Dùng các biện pháp transform để đưa chỉ số Skewness về càng gần 0 càng tốt.

Một số phương pháp:

- + log transformation
- + square root transformation
- + box-cox transformation
- + cube root transformation (căn bậc 3)

3. Lập trình

1. Import một số thư viện cần thiết

```
[2] import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

2. Đọc dữ liệu file excel:

```
[3] df = pd.read_excel("/content/gdrive/MyDrive/Machine_Learning/Customer_Segmentation/data.xlsx")
```

```
[4] df.head(10)
```

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	536365	85123A	WHITE HANGING HEART T-LIGHT HOLDER	6	2010-12-01 08:26:00	2.55	17850.0	United Kingdom
1	536365	71053	WHITE METAL LANTERN	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8	2010-12-01 08:26:00	2.75	17850.0	United Kingdom
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom
5	536365	22752	SET 7 BABUSHKA NESTING BOXES	2	2010-12-01 08:26:00	7.65	17850.0	United Kingdom
6	536365	21730	GLASS STAR FROSTED T-LIGHT HOLDER	6	2010-12-01 08:26:00	4.25	17850.0	United Kingdom
7	536366	22633	HAND WARMER UNION JACK	6	2010-12-01 08:28:00	1.85	17850.0	United Kingdom
8	536366	22632	HAND WARMER RED POLKA DOT	6	2010-12-01 08:28:00	1.85	17850.0	United Kingdom
9	536367	84879	ASSORTED COLOUR BIRD ORNAMENT	32	2010-12-01 08:34:00	1.69	13047.0	United Kingdom

```
[5] # Đếm thử số lượng dòng
df.count()
```

InvoiceNo	541909
StockCode	541909
Description	540455
Quantity	541909
InvoiceDate	541909
UnitPrice	541909
CustomerID	406829
Country	541909
dtype:	int64

Làm bài toán Sample với khoảng 10000 dòng:

0 giây

0 giây

Xem mô tả của cột InvoiceDate, ta thấy kiểu dữ liệu của nó đang ở dạng Object không phải dạng date.

```
0 giây # Tính toán RFM

##### Tính R = Recently: Khách hàng nào mua hàng gần nhất thì R nhỏ, mua hàng càng lâu thì R lớn

# Chuyển từ string sang date
df_not_nan['InvoiceDate'] = pd.to_datetime(df_not_nan['InvoiceDate'], format='%Y-%m-%d %H:%M:%S')

# Lấy ngày lớn nhất trong InvoiceDate + 1

import datetime
current_date = max(df_not_nan['InvoiceDate']) + datetime.timedelta(days=1)

##### Tính M = Moneytary Value
df_not_nan['TotalPay'] = df_not_nan['Quantity'] * df_not_nan['UnitPrice']

# Group by CustomerID để tính RFM
df_customers = df_not_nan.groupby(['CustomerID']).agg(
    {'InvoiceDate': lambda x: (current_date - x.max()).days,
     'InvoiceNo' : 'count',
     'TotalPay' : 'sum'
    }
)
```

Bây giờ ta có một data frame customer với cột InvoiceDate thể hiện khách hàng mua hàng xa hay gần với thời điểm hiện tại, InvoiceNo thể hiện số lượng đơn hàng khách hàng đã mua, TotalPay là tổng tiền khách hàng đã chi để mua hàng.

```
df_customers.head()
```

CustomerID	InvoiceDate	InvoiceNo	TotalPay
12347.0	40	5	133.20
12348.0	249	2	120.88
12349.0	19	2	312.75
12352.0	73	5	80.85
12354.0	233	2	33.30

Rename các cột lại cho đúng ý nghĩa:

```
0 giây [14] df_customers.rename(columns={'InvoiceDate': 'Recently', 'InvoiceNo': 'Frequently', 'TotalPay': 'MonetaryValue'}, inplace = True)
```

0 giây

```
df_customers.rename(columns={'InvoiceDate':'Recently', 'InvoiceNo':'Frequently', 'TotalPay':'MonetaryValue'})
```

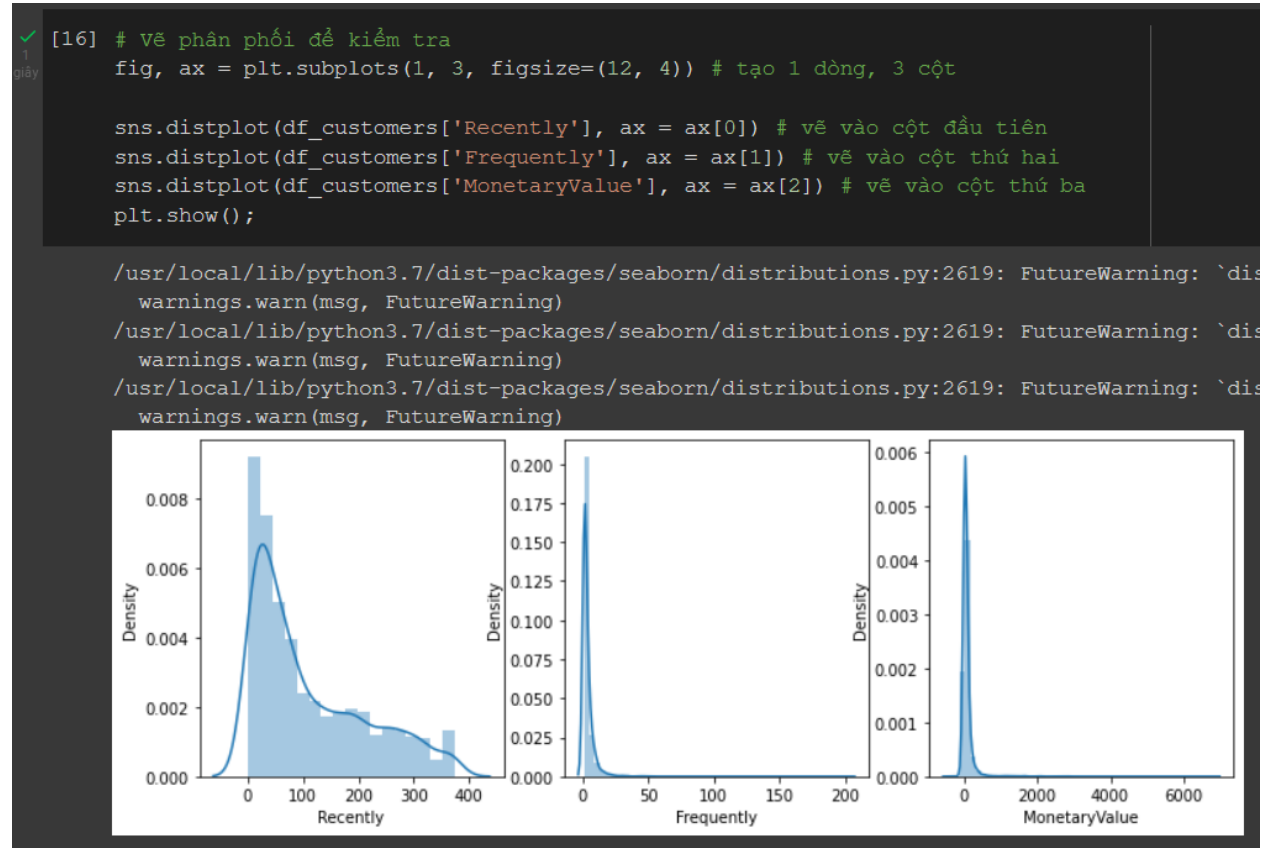
CustomerID	Recently	Frequently	MonetaryValue
12347.0	40	5	133.20
12348.0	249	2	120.88
12349.0	19	2	312.75
12352.0	73	5	80.85
12354.0	233	2	33.30
...
18265.0	73	2	39.60
18272.0	3	11	206.17
18274.0	18	2	-4.65
18283.0	4	21	78.03
18287.0	43	4	68.28

2690 rows x 3 columns

Bây giờ ta có thể dùng để đưa vào model KNN để phân cụm khách hàng, tuy nhiên hiệu quả sẽ không cao vì:

- + Các dữ liệu đang chưa được scale về một khoảng giống nhau.
- + Chưa biết phân phối dữ liệu.

4. Vẽ phân phối để kiểm tra dữ liệu và áp dụng các biện pháp Transform



Ta thấy dữ liệu chưa phân bố chuẩn. Ta phải áp dụng các biện pháp Transform.

Hàm kiểm tra xem áp dụng phương pháp nào để ra được phân phối chuẩn:

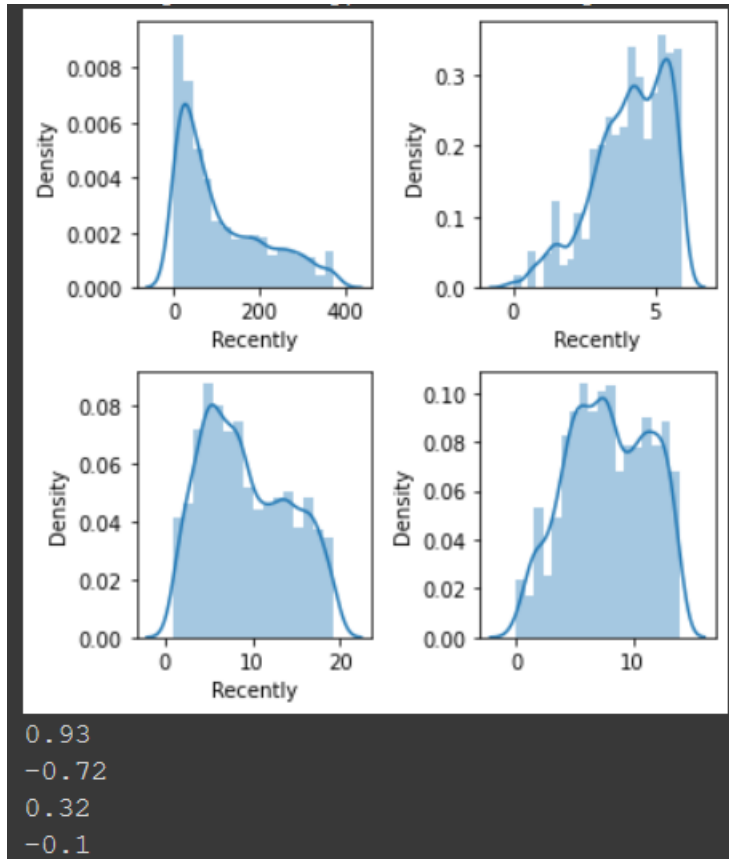
```
from scipy import stats
def analyze_skewness(x):
    fig, ax = plt.subplots(2, 2, figsize=(5,5))
    sns.distplot(df_customers[x], ax=ax[0,0])
    sns.distplot(np.log(df_customers[x]), ax=ax[0,1])
    sns.distplot(np.sqrt(df_customers[x]), ax=ax[1,0])
    sns.distplot(stats.boxcox(df_customers[x])[0], ax=ax[1,1])
    plt.tight_layout()
    plt.show()

    print(df_customers[x].skew().round(2))
    print(np.log(df_customers[x]).skew().round(2))
    print(np.sqrt(df_customers[x]).skew().round(2))
    print(pd.Series(stats.boxcox(df_customers[x])[0]).skew().round(2))
```

Với cột Recently: Xem các phân bố khi áp dụng các biện pháp Transformation bằng lệnh:

```
analyze_skewness('Recently')
```

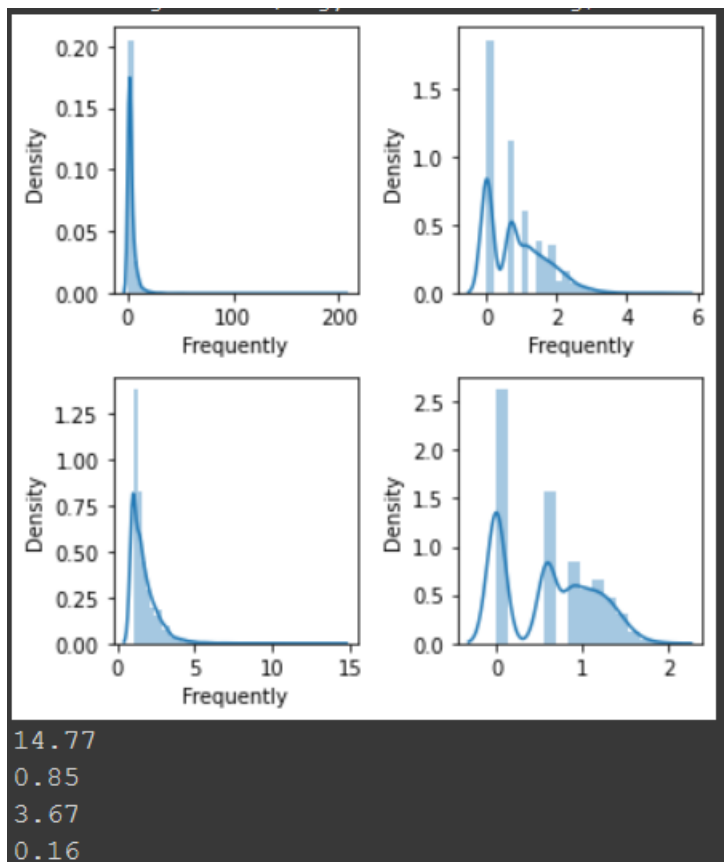
Kết quả: ta thấy box-cox transformation có chỉ số skewness là -0,1 gần với 0 nhất nên nó là tốt nhất.



Với cột Frequently:

```
analyze_skewness('Frequently')
```

Kết quả: Ta thấy box-cox Transformation vẫn cho kết quả tốt nhất (skewness = 0,16)



Với cột MonetaryValue: vì có giá trị âm nên ta không dùng được các biện pháp Transform như box-cox, sqr root hay log mà ta phải dùng cube root (căn bậc 3).

Tiến hành Transform dữ liệu:

```
[25] # Tiến hành Transform dữ liệu

# Tạo 1 dữ liệu riêng để Transform mà không làm thay đổi dữ liệu gốc
df_customers_T = pd.DataFrame()

df_customers_T['Recently'] = stats.boxcox(df_customers['Recently'])[0]
df_customers_T['Frequently'] = stats.boxcox(df_customers['Frequently'])[0]
df_customers_T['MonetaryValue'] = pd.Series(np.cbrt(df_customers['Recently])).values
```

✓ 0 giây

df_customers_T.head()

	Recently	Frequently	MonetaryValue
0	6.152222	1.127547	3.419952
1	12.180045	0.591193	6.291195
2	4.402867	0.591193	2.668402
3	7.832068	1.127547	4.179339
4	11.907953	0.591193	6.153449

5. Scale dữ liệu

✓ 0 giây

```
[28] # Tiến hành scale dữ liệu
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaler.fit(df_customers_T)

df_customers_T = scaler.transform(df_customers_T)
```

✓ 0 giây

pd.DataFrame(df_customers_T).head()

	0	1	2
0	-0.493794	1.012426	-0.533456
1	1.232949	-0.017412	1.262754
2	-0.994917	-0.017412	-1.003616
3	-0.012582	1.012426	-0.058394
4	1.155005	-0.017412	1.176583

6. Đưa dữ liệu vào K-means để phân khúc khách hàng

→ Tiến hành chọn số cụm bằng phương pháp Elbow:

```
[31] from sklearn.cluster import KMeans

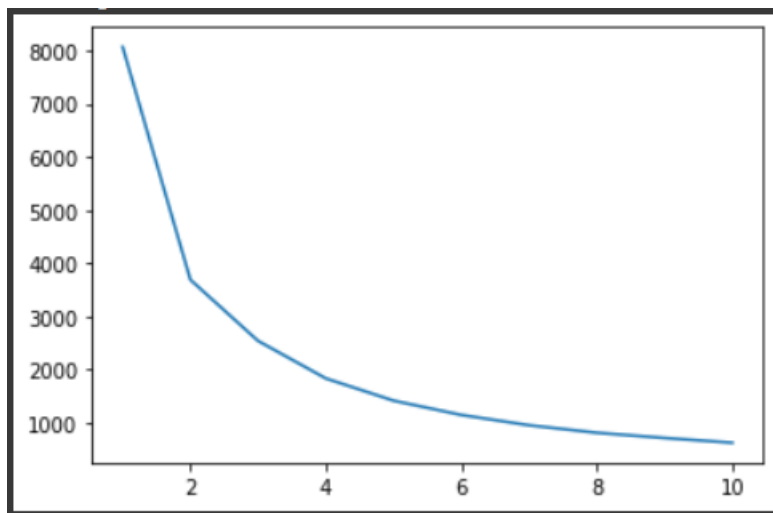
# Chọn K cluster bằng thuật toán Elbow

k_rng = range(1, 11) #10 giá trị k
sse = []

for k in k_rng:
    km = KMeans(n_clusters=k)
    km.fit(df_customers_T)
    sse.append(km.inertia_)

print(sse)
plt.plot(k_rng, sse)
```

Kết quả:



Chọn k sao cho khi k tăng lên 1 thì giá trị SSE không giảm đi quá nhiều. Vì vậy dựa vào đồ thị trên ta chọn $k = 3$ là phù hợp.

Áp dụng thuật toán phân cụm K-means:

```
[42] # Áp dụng model K-Means vào bài toán phân cụm khách hàng
km = KMeans(n_clusters=3)
y_predicted = km.fit_predict(df_customers_T)
y_predicted
```

```
array([0, 1, 0, ..., 0, 0, 0], dtype=int32)
```

```
df_customers['Cluster'] = y_predicted
df_customers.head()
```

	Recently	Frequently	MonetaryValue	Cluster
CustomerID				
12347.0	40	5	133.20	0
12348.0	249	2	120.88	1
12349.0	19	2	312.75	0
12352.0	73	5	80.85	0
12354.0	233	2	33.30	1

Đặc trưng của từng nhóm khách hàng:

```
[52] # Đặc trưng của các cụm khách hàng
df_customers.groupby('Cluster').agg(
    {
        'Recently': 'mean',
        'Frequently': 'mean',
        'MonetaryValue': 'mean'
    }
).round(2)
```

	Recently	Frequently	MonetaryValue
Cluster			
0	29.88	7.40	145.45
1	223.12	1.99	40.10
2	56.46	1.40	27.99

Nhóm khách hàng 0: Mua gần đây, số lần mua khá nhiều và tổng chi phí họ mua rất lớn
→ Một số khách hàng tiềm năng gần đây, cần có những biện pháp tri ân khách hàng.

Nhóm khách hàng 1: Mua hàng rất lâu rồi, số lần mua khá ít và tổng chi phí họ mua ở mức trung bình → đây có thể là nhóm khách hàng mới sử dụng dịch vụ nên e ngại mua số lượng lớn và họ chỉ dừng ở mức mua để dùng thử.

Nhóm khách hàng 2: Mua hàng khá gần, số lần mua ít và chi ít tiền.