

```

import pandas as pd

import numpy as np

df = pd.read_csv('./sales_data_sample.csv', encoding='unicode_escape')

df.head

df.info

#Columns to Remove

to_drop = ['ADDRESSLINE1', 'ADDRESSLINE2', 'STATE', 'POSTALCODE', 'PHONE']

df = df.drop(to_drop, axis=1)

#Check for null values

df.isnull().sum()

#

df.dtypes

#ORDERDATE Should be in date time

df['ORDERDATE'] = pd.to_datetime(df['ORDERDATE'])


#MonetaryValue : Revenue generated by the customers

import datetime as dt

snapshot_date = df['ORDERDATE'].max() + dt.timedelta(days = 1)

df_RFM = df.groupby(['CUSTOMERNAME']).agg({

    'ORDERDATE' : lambda x : (snapshot_date - x.max()).days,

    'ORDERNUMBER' : 'count',

    'SALES' : 'sum'

})


#Rename the columns

df_RFM.rename(columns = {

    'ORDERDATE' : 'Recency',

    'ORDERNUMBER' : 'Frequency',

    'SALES' : 'MonetaryValue'

}, inplace=True)

df_RFM.head()

```

```

# Divide into segments

# We create 4 quartile ranges

df_RFM['M'] = pd.qcut(df_RFM['MonetaryValue'], q = 4, labels = range(1,5))
df_RFM['R'] = pd.qcut(df_RFM['Recency'], q = 4, labels = list(range(4,0,-1)))
df_RFM['F'] = pd.qcut(df_RFM['Frequency'], q = 4, labels = range(1,5))

df_RFM.head()

#Create another column for RFM score
df_RFM['RFM_Score'] = df_RFM[['R', 'M', 'F']].sum(axis=1)
df_RFM.head()

#
def rfm_level(df):
    if bool(df['RFM_Score'] >= 10):
        return 'High Value Customer'

    elif bool(df['RFM_Score'] < 10) and bool(df['RFM_Score'] >= 6):
        return 'Mid Value Customer'

    else:
        return 'Low Value Customer'

df_RFM['RFM_Level'] = df_RFM.apply(rfm_level, axis = 1)
df_RFM.head()

#
# Time to perform KMeans
data = df_RFM[['Recency', 'Frequency', 'MonetaryValue']]
data.head()

# Our data is skewed we must remove it by performing log transformation
data_log = np.log(data)
data_log.head()

```

#Standardization

```
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()

scaler.fit(data_log)

data_normalized = scaler.transform(data_log)

data_normalized = pd.DataFrame(data_normalized, index = data_log.index,
columns=data_log.columns)

data_normalized.describe().round(2)
```

#Fit KMeans and use elbow method to choose the number of clusters

```
import matplotlib.pyplot as plt

import seaborn as sns

from sklearn.cluster import KMeans
```

```
sse = {}
```

```
for k in range(1, 21):
```

```
    kmeans = KMeans(n_clusters = k, random_state = 1)
```

```
    kmeans.fit(data_normalized)
```

```
    sse[k] = kmeans.inertia_
```

```
#
```

```
plt.figure(figsize=(10,6))
```

```
plt.title('The Elbow Method')
```

```
plt.xlabel('K')
```

```
plt.ylabel('SSE')
```

```
plt.style.use('ggplot')
```

```
sns.pointplot(x=list(sse.keys()), y = list(sse.values()))
```

```
plt.text(4.5, 60, "Largest Angle", bbox = dict(facecolor = 'lightgreen', alpha = 0.5))
```

```
plt.show()
```

5 number of clusters seems good

```
kmeans = KMeans(n_clusters=5, random_state=1)
```

```
kmeans.fit(data_normalized)
```

```
cluster_labels = kmeans.labels_
```

```
data_rfm = data.assign(Cluster = cluster_labels)
```

```
data_rfm.head()
```