In [102 original_df= pd_read_esy('online_retail_ff.csy') df= original_df: csy(') df= original_df: csy(
4 489434 21732 STRAWBERRY CERAMIC TRINKET BOX 24 2009-12-01 07.4500 1.25 1388.50 United Kingdom 1067366 581587 22899 CHILDRENS APRON DOLLY GIRL 6 2011-12-09 12:50.00 2.10 12680.0 France 1067367 581587 23254 CHILDRENS CUTLERY DOLLY GIRL 4 2011-12-09 12:50.00 4.15 12680.0 France 1067368 581587 23255 CHILDRENS CUTLERY CIRCUS PARADE 4 2011-12-09 12:50.00 4.95 12680.0 France 1067370 581587 POST POSTAGE 1 2011-12-09 12:50.00 18.00 12680.0 France 1067371 rows × 8 columns Data Cleaning	
Out[103 (1067371, 8) In [104 df.info()	
3 Quantity 1067371 non-null int64 4 InvoiceDate 1067371 non-null object 5 Price 1067371 non-null float64 6 Customer ID 824364 non-null float64 7 Country 1067371 non-null object dtypes: float64(2), int64(1), object(5) memory usage: 65.1+ MB In [105 df['Price'].describe() Out[105 count 1.067371e+06 mean 4.649388e+00 std 1.235531e+02 min -5.359436e+04	
25% 1.250000e+00 50% 2.10000e+00 75% 4.150000e+00 max 3.897000e+04 Name: Price, dtype: float64 In [106 df.isnull().sum() Out[106 Invoice 0 StockCode 0 Description 4382 Quantity 0 InvoiceDate 0 InvoiceDate 0	
Price 0 Customer ID 243007 Country 0 dtype: int64 In [107 df.dropna(subset=['Customer ID'], inplace=True) df.shape Out[107 (824364, 8) In [108 for col in df.columns:	
Invoice 576339 542 579196 533 580727 529 578270 442 573576 435 502263 1 CS14128 1 S46102 1 CS46106 1	
C49427 1 Name: count, Length: 44876, dtype: int64	
72751C 1 9011E 1 Name: count, length: 4646, dtype: int64	
RUSSIAN FOLKART STACKING TINS 1 BLUE BARQUE FLOCK CANDLE HOLDER 1 Name: count, Length: \$299, dtype: int64Value counts for column: Quantity Quantity 1 153683 12 119842 2 117920 6 30249 4 62209 1968 1 1960 1	
9360 1 6336 1 -80995 1 Name: count, Length: 643, dtype: int64Value counts for column: InvoiceDate 2011-11-14 15:27:00 543 2011-12-05 17:17:00 534 2011-12-05 17:17:00 530 2011-11-23 13:39:00 444 2011-03-11 14:09:00 436 2011-02-11 11:44:00 1	
2010-10-04 11:52:00 1 2010-10-04 11:08:00 1 2010-10-05-04 16:28:00 1 Name: count, length: 41439, dtype: int64	
75.19 1 280.91 1 129.95 1 1.77 1 224.69 1 Name: count, Length: 1022, dtype: int64	
14580.0 1 16154.0 1 16443.0 1 15233.0 1 17948.0 1 Name: count, Length: 5942, dtype: int64	
Netherlands 5140 Spain 3811 Belgium 3123 Switzerland 3064 Portugal 2504 Australia 1913 Channel Islands 1664 Italy 1634 Norway 1455 Sweden 1345 Cyprus 1176 Finland 1049 Austria 938 Demmark 817 Greece 663	
Japan 592 USA 535 Poland 535 Unspecified 524 United Arab Emirates 386 Singapore 346 Israel 324 Malta 299 Iceland 253 Canada 258 Lithuania 189 RSA 123 Brazil 94 Thálland 76	
Rorea 63	
Out[109 Invoice 576339 542 579196 533 580727 529 578270 442 573576 435 502263 1 514528 1 546108 1 494927 1 Name: count, Length: 44876, dtype: int64	
<pre>In (110.</pre>	
327 80995 1 Name: count, Length: 438, dtype: int64 In [111. df.dtypes Out[11. Invoice int32 StockCode object Description object Quantity int64 InvoiceDate object Price float64 Country object Country object dtype: object	
<pre>In [112_ df['InvoiceDate'] = pd.to_datetime(df['InvoiceDate']) df['Customer ID'] = df['Customer ID'].astype(int) In [113_ df.dtypes Out[113_ Invoice</pre>	
<pre>country dtype: object In [114_ df.columns = df.columns.str.strip().str.lower().str.replace(' ', '_') In [115_ df.columns Out[115_ Index(['invoice', 'stockcode', 'description', 'quantity', 'invoicedate',</pre>	
Out [116 Invoice stockcode quantity invoicedate price customer_id country total_price 0 489434 85048 12 2009-12-01 07:45:00 6.75 13085 United Kingdom 81.00 1 489434 79323W 12 2009-12-01 07:45:00 6.75 13085 United Kingdom 81.00 2 489434 79323W 12 2009-12-01 07:45:00 6.75 13085 United Kingdom 81.00 3 489434 22041 48 2009-12-01 07:45:00 2.10 13085 United Kingdom 100.80 4 489434 2132 24 2009-12-01 07:45:00 1.25 13085 United Kingdom 30.00	
1067366 581587 22899 6 2011-12-09 12:50:00 2.10 12680 France 12.60 1067367 581587 23254 4 2011-12-09 12:50:00 4.15 12680 France 16.60 1067368 581587 23255 4 2011-12-09 12:50:00 4.15 12680 France 16.60 1067369 581587 22138 3 2011-12-09 12:50:00 4.95 12680 France 14.85 1067370 581587 POST 1 2011-12-09 12:50:00 18.00 12680 France 18.00 805620 rows × 8 columns	
RFM Feature Engineering In [117_ snapshot_date = df['invoicedate'].max() + pd.Timedelta(days=1) In [118_ rfm_data = df.groupby('customer_id').agg({ 'invoicedate': lambda date: (snapshot_date - date.max()).days,	
'invoice': 'frequency',	
12349 19 4 4428.69 12350 310 1 334.40 18283 4 22 2736.65 18284 432 1 461.68 18285 661 1 427.00 18286 477 2 1296.43	
18287 43 7 4182.99 5881 rows × 3 columns EDA In [120 rfm_cols= ['recency', 'frequency', 'monetary'] for col in rfm_cols: plt.figure(figu	
2000 - 1750 - 1500 - 1250 - 1000 - 750 -	
500 - 250 -	
1000 - 600 -	
1000	
200	
customer_id 12346 5.789960 2.564949 11.258774 12347 1.098612 2.197225 8.636632 12348 4.330733 1.791759 7.611051 12349 2.995732 1.609438 8.396085 12350 5.739793 0.693147 5.815324	
1828 1.609438 3.135494 7.914855 1828 6.070738 0.693147 6.137036 1828 6.495266 0.693147 6.059123 1828 6.169611 1.098612 7.168141 1828 3.784190 2.079442 8.339021 5881 rows × 3 columns	
plt.figure(figsize=(6,4)) sns.histplot(x= log_transformation(col), kde=True) 600- 500- 400- 300-	
200 - 100 -	
1200 - 1000 -	
350 - 300 - 250 -	
In [123 merged_df= pd.merge(df, log_transformation, on='customer_id') merged_df	
Invoice Stockcode quantity Invoicedate price customer_jid country total_price recency requency requency monetary	
805615 581587 2289 6 2011-12-09 12:50:00 2.10 12680 France 12.60 0.693147 1.609438 6.781977 805616 581587 23254 4 2011-12-09 12:50:00 4.15 12680 France 16.60 0.693147 1.609438 6.781977 805617 581587 23255 4 2011-12-09 12:50:00 4.15 12680 France 16.60 0.693147 1.609438 6.781977 805618 581587 22138 3 2011-12-09 12:50:00 4.95 12680 France 14.85 0.693147 1.609438 6.781977 805619 581587 POST 1 2011-12-09 12:50:00 18.00 12680 France 18.00 0.693147 1.609438 6.781977 805620 rows × 11 columns In [124. sns.heatmap (merged_df.corr (numeric_only= True), annot=True)	
Out [124 <axes:> invoice - 1</axes:>	
total_price = 0.0003	
Customer Segmentation with K-Means Clustering In [125 scaler= StandardScaler() df_scaled= scaler.fit_transform(log_transformation) df_scaled Out [125 array([[0.85632068,	
[1.3085854 , -1.05756818, -0.55461809], [1.3085854 , -0.5566389 , 0.24403565], [-0.42984452 , 0.65512034 , 1.07880297]]) In [126_ inertia= [] k_range= range(1,11) for k in k_range: kmeans= RMeans(n_clusters=k) kmeans, fit(df_scaled) inertia.append(kmeans.inertia_) inertia	
Out[126 [17643.0, 8614.946158654508, 6382.42305711746, 4349.228804709554, 4129.989857438568, 3586.990413231882, 320.8055188684775, 2320.167401845461, 2688.12466404995, 2498.7204000535157] In [127 plt.plot(k_range, inertia, marker='o') Out[127 [<matplotlib.lines.line2d 0xle3356a1580="" at="">]</matplotlib.lines.line2d>	
18000 - 14000 - 12000 - 10000 - 8000 - 6000 - 4000 -	
In [128 kmeans_final= KMeans (n_clusters=4) kmeans_final Out [128 v KMeans 0 0	
Out[129_ array([1, 1, 0,, 3, 0, 1]) In [130_ rfm_data('cluster')= cluster_labels rfm_data('cluster').value_counts() Out[130_ cluster	
In [131. # Calculate the mean RFM values and count for each cluster cluster_summary = rfm_data_groupby('cluster').agg({ ''recquency': 'mean', 'requency': 'mean', ''anotary': 'mean', 'cluster': 'size' # Get the number of customers in each cluster }).rename(columns={'cluster': 'Size'}) cluster_summary.sort_values(by='monetary', ascending=False, inplace=True) cluster_summary Cut[131.	
1 27.505042 19.329412 11001 208703 1190 0 228.783913 5.075665 1995.093748 1467 2 28.378997 3.045527 888.233196 1252 3 386.987830 1.376268 323.610716 1972 1n 132.	
rfm_data.to_csv('rfm_data.csv') cut(132. recency fequency monetary cluster Segment customer_id 12346 326 12 77556.46 1 Lost Customers 12347 2 8 5633.32 1 Lost Customers 12348 75 5 2019.40 0 Champions 12349 19 4 4428.69 2 At-Risk Loyalists	
<pre>5881 rows × 5 columns In [133 plt.figure(figsize=(8, 6)) sns.scatterplot(data=rfm_data, x='recency', y= frequency', hue='cluster', size='monetary', size='monetary', palette='viridid', palette='viridid',</pre>	

In [101... import pandas as pd

import numpy as np
import seaborn as sns

import warnings

Loading Data

warnings.filterwarnings('ignore')

import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans