

```
In [1]: import pandas as pd

In [3]: customers = pd.read_csv(r"C:\Users\yashu\Downloads\Customers.csv")
products = pd.read_csv(r"C:\Users\yashu\Downloads\Products.csv")
transactions = pd.read_csv(r"C:\Users\yashu\Downloads\Transactions.csv")

In [5]: customers

Out[5]:
```

	CustomerID	CustomerName	Region	SignupDate
0	C0001	Lawrence Carroll	South America	2022-07-10
1	C0002	Elizabeth Lutz	Asia	2022-02-13
2	C0003	Michael Rivera	South America	2024-03-07
3	C0004	Kathleen Rodriguez	South America	2022-10-09
4	C0005	Laura Weber	Asia	2022-08-15
...
195	C0196	Laura Watts	Europe	2022-06-07
196	C0197	Christina Harvey	Europe	2023-03-21
197	C0198	Rebecca Ray	Europe	2022-02-27
198	C0199	Andrea Jenkins	Europe	2022-12-03
199	C0200	Kelly Cross	Asia	2023-06-11

200 rows × 4 columns

```
In [7]: products

Out[7]:
```

	ProductID	ProductName	Category	Price
0	P001	ActiveWear Biography	Books	169.30
1	P002	ActiveWear Smartwatch	Electronics	346.30
2	P003	ComfortLiving Biography	Books	44.12
3	P004	BookWorld Rug	Home Decor	95.69
4	P005	TechPro T-Shirt	Clothing	429.31
...
95	P096	SoundWave Headphones	Electronics	307.47
96	P097	BookWorld Cookbook	Books	319.34
97	P098	SoundWave Laptop	Electronics	299.93
98	P099	SoundWave Mystery Book	Books	354.29
99	P100	HomeSense Sweater	Clothing	126.34

100 rows × 4 columns

```
In [9]: transactions

Out[9]:
```

	TransactionID	CustomerID	ProductID	TransactionDate	Quantity	TotalValue	Price
0	T00001	C0199	P067	2024-08-25 12:38:23	1	300.68	300.68
1	T00112	C0146	P067	2024-05-27 22:23:54	1	300.68	300.68
2	T00166	C0127	P067	2024-04-25 07:38:55	1	300.68	300.68
3	T00272	C0087	P067	2024-03-26 22:55:37	2	601.36	300.68
4	T00363	C0070	P067	2024-03-21 15:10:10	3	902.04	300.68
...
995	T00496	C0118	P037	2024-10-24 08:30:27	1	459.86	459.86
996	T00759	C0059	P037	2024-06-04 02:15:24	3	1379.58	459.86
997	T00922	C0018	P037	2024-04-05 13:05:32	4	1839.44	459.86
998	T00959	C0115	P037	2024-09-29 10:16:02	2	919.72	459.86
999	T00992	C0024	P037	2024-04-21 10:52:24	1	459.86	459.86

1000 rows × 7 columns

```
In [11]: from sklearn.preprocessing import OneHotEncoder
import numpy as np

In [13]: merged_data = transactions.merge(customers, on="CustomerID", how="left").merge(products, on="ProductID", how="left")

In [15]: merged_data

Out[15]:
```

	TransactionID	CustomerID	ProductID	TransactionDate	Quantity	TotalValue	Price_x	CustomerName	Region	SignupDate	ProductName	Category	Price_y
0	T00001	C0199	P067	2024-08-25 12:38:23	1	300.68	300.68	Andrea Jenkins	Europe	2022-12-03	ComfortLiving Bluetooth Speaker	Electronics	300.68
1	T00112	C0146	P067	2024-05-27 22:23:54	1	300.68	300.68	Brittany Harvey	Asia	2024-09-04	ComfortLiving Bluetooth Speaker	Electronics	300.68
2	T00166	C0127	P067	2024-04-25 07:38:55	1	300.68	300.68	Kathryn Stevens	Europe	2024-04-04	ComfortLiving Bluetooth Speaker	Electronics	300.68
3	T00272	C0087	P067	2024-03-26 22:55:37	2	601.36	300.68	Travis Campbell	South America	2024-04-11	ComfortLiving Bluetooth Speaker	Electronics	300.68
4	T00363	C0070	P067	2024-03-21 15:10:10	3	902.04	300.68	Timothy Perez	Europe	2022-03-15	ComfortLiving Bluetooth Speaker	Electronics	300.68
...
995	T00496	C0118	P037	2024-10-24 08:30:27	1	459.86	459.86	Jacob Holt	South America	2022-01-22	SoundWave Smartwatch	Electronics	459.86
996	T00759	C0059	P037	2024-06-04 02:15:24	3	1379.58	459.86	Mrs. Kimberly Wright	North America	2024-04-07	SoundWave Smartwatch	Electronics	459.86
997	T00922	C0018	P037	2024-04-05 13:05:32	4	1839.44	459.86	Tyler Haynes	North America	2024-09-21	SoundWave Smartwatch	Electronics	459.86
998	T00959	C0115	P037	2024-09-29 10:16:02	2	919.72	459.86	Joshua Hamilton	Asia	2024-11-11	SoundWave Smartwatch	Electronics	459.86
999	T00992	C0024	P037	2024-04-21 10:52:24	1	459.86	459.86	Michele Cooley	North America	2024-02-05	SoundWave Smartwatch	Electronics	459.86

1000 rows × 13 columns

```
In [17]: customer_spending = merged_data.groupby('CustomerID')['TotalValue'].sum()

In [19]: customer_spending
```

```
Out[19]: CustomerID
C0001    3354.52
C0002    1862.74
C0003    2725.38
C0004    5354.88
C0005    2034.24
...
C0196    4982.88
C0197    1928.65
C0198     931.83
C0199    1979.28
C0200    4758.60
Name: TotalValue, Length: 199, dtype: float64

In [21]: avg_transaction_value = merged_data.groupby('CustomerID')['TotalValue'].mean()

In [23]: avg_transaction_value
```

```
Out[23]: CustomerID
C0001    670.904000
C0002    465.685000
C0003    681.345000
C0004    669.360000
C0005    678.080000
...
C0196    1245.720000
C0197    642.883333
C0198    465.915000
C0199    494.820000
C0200    951.720000
Name: TotalValue, Length: 199, dtype: float64

In [25]: category_purchases = merged_data.groupby(['CustomerID', 'Category']).size().unstack(fill_value=0)

In [27]: category_purchases

Out[27]:
```

	Category	Books	Clothing	Electronics	Home Decor
CustomerID					
	C0001	1	0	3	1
	C0002	0	2	0	2
	C0003	0	1	1	2
	C0004	3	0	2	3
	C0005	0	0	2	1

	C0196	1	1	0	2
	C0197	0	0	2	1
	C0198	0	1	1	0
	C0199	0	0	2	2
	C0200	1	2	1	1

199 rows × 4 columns

```
In [45]: encoder = OneHotEncoder(sparse_output=False, handle_unknown='ignore')
region_encoded = pd.DataFrame(encoder.fit_transform(customers[['Region']]), columns=encoder.get_feature_names_out(['Region']))
region_encoded.index = customers['CustomerID']

In [47]: customer_features = pd.concat([
    customer_spending.rename("TotalSpending"),
    avg_transaction_value.rename("AvgTransactionValue"),
    category_purchases,
    region_encoded
], axis=1).fillna(0)

In [49]: customer_features

Out[49]:
```

	TotalSpending	AvgTransactionValue	Books	Clothing	Electronics	Home Decor	Region_Asia	Region_Europe	Region_North America	Region_South America
CustomerID										
	C0001	3354.52	670.904000	1.0	0.0	3.0	1.0	0.0	0.0	1.0
	C0002	1862.74	465.685000	0.0	2.0	0.0	2.0	1.0	0.0	0.0
	C0003	2725.38	681.345000	0.0	1.0	1.0	2.0	0.0	0.0	1.0
	C0004	5354.88	669.360000	3.0	0.0	2.0	3.0	0.0	0.0	1.0
	C0005	2034.24	678.080000	0.0	0.0	2.0	1.0	1.0	0.0	0.0

	C0197	1928.65	642.883333	0.0	0.0	2.0	1.0	0.0	1.0	0.0
	C0198	931.83	465.915000	0.0	1.0	1.0	0.0	0.0	1.0	0.0
	C0199	1979.28	494.820000	0.0	0.0	2.0	2.0	0.0	1.0	0.0
	C0200	4758.60	951.720000	1.0	2.0	1.0	1.0	0.0	0.0	0.0
	C0180	0.00	0.000000	0.0	0.0	0.0	0.0	1.0	0.0	0.0

200 rows × 10 columns

```
In [53]: from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
scaled_features = scaler.fit_transform(customer_features)
scaled_features_df = pd.DataFrame(scaled_features, index=customer_features.index, columns=customer_features.columns)
scaled_features_df.head()

Out[53]:
```

	TotalSpending	AvgTransactionValue	Books	Clothing	Electronics	Home Decor	Region_Asia	Region_Europe	Region_North America	Region_South America
CustomerID										
	C0001	-0.051884	-0.054781	-0.314627	-1.036192	1.555406	-0.215318	-0.538816	-0.57735	-0.546536
	C0002	-0.862714	-0.903985	-1.213560	0.781689	-1.141830	0.681841	1.855921	-0.57735	-0.546536
	C0003	-0.393842	-0.011575	-1.213560	-0.127252	-0.242751	0.681841	-0.538816	-0.57735	-0.546536
	C0004	1.035375	-0.061170	1.483240	-1.036192	0.656327	1.578999	-0.538816	-0.57735	-0.546536
	C0005	-0.769499	-0.025086	-1.213560	-1.036192	0.656327	-0.215318	1.855921	-0.57735	-0.546536

```
In [63]: import pandas as pd
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.metrics.pairwise import cosine_similarity

In [65]: scaler = StandardScaler()
normalized_features = scaler.fit_transform(customer_features)

In [67]: similarity_matrix = cosine_similarity(normalized_features)

In [71]: lookalike_results = {}
for i, customer in enumerate(customers):
    similarity_scores = list(enumerate(similarity_matrix[i]))
    similarity_scores = sorted(similarity_scores, key=lambda x: x[1], reverse=True)
    top_3 = [(customer_features.index[j], score) for j, score in similarity_scores[1:4]]
    lookalike_results[customer] = top_3

In [73]: lookalike_df = pd.DataFrame.from_dict(lookalike_results, orient='index', columns=['Lookalike1', 'Lookalike2', 'Lookalike3'])
lookalike_df.to_csv('Lookalike.csv', index_label='CustomerID')

In [75]: lookalike_df

Out[75]:
```

	Lookalike1	Lookalike2	Lookalike3
CustomerID	(C0120, 0.8841868127875581)	(C0091, 0.870710992240339)	(C0190, 0.8694368338254506)
CustomerName	(C0134, 0.3490436864136428)	(C0106, 0.9268340500713907)	(C0158, 0.8098204349873207)
Region	(C0031, 0.344486708707058)	(C0129, 0.9069501773231374)	(C0158, 0.8093560554513439)

