**Task 2: Building a Lookalike Model**

To create a lookalike model, we will:

1. Preprocess and combine datasets.
2. Engineer relevant features using customer profiles and transaction history.
3. Use a similarity measure to find the top 3 similar customers for each target customer.
4. Save the results in the specified format: Map<cust_id, List<cust_id, score>>.

we will use 2 approaches to figure out and save it as lookalike models

1. **K-Nearest Neighbors (KNN) Approach**

We will now use customer and transaction features together for similarity calculation.

**Steps:**

1. **Combine Data:**
   o Merge Transactions with Customers and Products to form a comprehensive dataset.
   o Aggregate transaction data (e.g., total spending, total quantity purchased) for each customer.

2. **Feature Engineering:**
   o Encode categorical features like Region and Category.
   o Calculate summary statistics for transactions (e.g., total spend, average transaction value).

3. **Scale Features:** Normalize numerical columns for similarity calculation.
4. **Fit KNN:** Use the combined feature space for similarity.

```
import pandas as pd
```
5. from sklearn.preprocessing import StandardScaler, LabelEncoder
6. from sklearn.neighbors import NearestNeighbors
7.
8. # Load datasets
9. customers = pd.read_csv('Customers.csv')
10. products = pd.read_csv('Products.csv')
11. transactions = pd.read_csv('Transactions.csv')
12.
13. # Merge datasets
14. merged_data = transactions.merge(customers, on='CustomerID').merge(products, on='ProductID')
15.
16. # Aggregate transaction data for each customer
17. customer_agg = merged_data.groupby('CustomerID').agg({
18.   'TotalValue': ['sum', 'mean'],  # Total and average spending
19.   'Quantity': ['sum'],          # Total quantity purchased
20.   'Category': lambda x: x.nunique(),  # Number of unique categories bought
21.   'Region': 'first',           # Region of the customer
22.   'SignupDate': 'first'        # Signup date
23. }).reset_index()

24. customer_agg.columns = ['CustomerID', 'TotalSpend', 'AvgSpend', 'TotalQuantity', 'UniqueCategories', 'Region', 'SignupDate']

25.

26. # Encode categorical variables

27. le_region = LabelEncoder()

28. customer_agg['RegionEncoded'] = le_region.fit_transform(customer_agg['Region'])

29.

30. # Convert SignupDate to numerical (e.g., days since signup)

31. customer_agg['SignupDate'] = pd.to_datetime(customer_agg['SignupDate'])

32. customer_agg['DaysSinceSignup'] = (pd.Timestamp.now() - customer_agg['SignupDate']).dt.days

33. customer_agg = customer_agg.drop(['Region', 'SignupDate'], axis=1)

34.

35. # Scale features

36. scaler = StandardScaler()

37. scaled_features = scaler.fit_transform(customer_agg.drop('CustomerID', axis=1))

38.

39. # Fit KNN model

40. knn_model = NearestNeighbors(n_neighbors=4, metric='euclidean')

41. knn_model.fit(scaled_features)

42.

43. # Find top 3 similar customers for first 20 customers

44. lookalike_results_knn = {}

45. for idx, customer_id in enumerate(customer_agg['CustomerID'][:20]):

46.     distances, indices = knn_model.kneighbors([scaled_features[idx]], n_neighbors=4)  # Include self

47.     lookalikes = [(customer_agg['CustomerID'][i], 1 / (1 + distances[0][j])) for j, i in enumerate(indices[0][1:])]  # Exclude self

48.     lookalike_results_knn[customer_id] = lookalikes

49.

50. # Save results to CSV

51. lookalike_df_knn = pd.DataFrame({

52.     "CustomerID": list(lookalike_results_knn.keys()),

53.     "Lookalikes": [str(v) for v in lookalike_results_knn.values()]

54. })

55. lookalike_df_knn.to_csv("Lookalike_KNN.csv", index=False)


**2. Collaborative Filtering (Matrix Factorization) Approach**
Use CustomerID and ProductID to create a sparse user-product matrix and apply matrix factorization.
**Steps:**
1. **Prepare User-Product Matrix:**
   o Create a pivot table with CustomerID as rows, ProductID as columns, and TotalValue as values.
   o Fill missing values with 0.

2. **Apply SVD:**
   - o Factorize the matrix using TruncatedSVD to reduce dimensions.
3. **Find Similar Customers:**
   - o Compute cosine similarity between reduced vectors.

```python
from sklearn.decomposition import TruncatedSVD
from sklearn.metrics.pairwise import cosine_similarity

# Create user-product matrix
user_product_matrix = merged_data.pivot_table(index='CustomerID', columns='ProductID',
values='TotalValue', fill_value=0)

# Apply SVD
svd = TruncatedSVD(n_components=10, random_state=42)
reduced_features = svd.fit_transform(user_product_matrix)

# Compute cosine similarity
similarity_matrix = cosine_similarity(reduced_features)
similarity_df = pd.DataFrame(similarity_matrix, index=user_product_matrix.index,
columns=user_product_matrix.index)

# Get top 3 similar customers for first 20 customers
lookalike_results_svd = {}
for customer_id in similarity_df.index[:20]:
    similar_customers = similarity_df[customer_id].sort_values(ascending=False).iloc[1:4]  # Exclude
self
    lookalike_results_svd[customer_id] = list(zip(similar_customers.index, similar_customers.values))

# Save results to CSV
lookalike_df_svd = pd.DataFrame({
    "CustomerID": list(lookalike_results_svd.keys()),
    "Lookalikes": [str(v) for v in lookalike_results_svd.values()]
})
lookalike_df_svd.to_csv("Lookalike_SVD.csv", index=False)
```