

```
1 !pip install implicit
2 !pip install gradio
```

```
1 import implicit
2 import numpy as np
3 import pandas as pd
4 from scipy.sparse import csr_matrix
5 import tensorflow as tf
6 from tqdm import tqdm
7 from sentence_transformers import SentenceTransformer
8 from tensorflow.keras import layers, Model, losses, optimizers, callbacks
9 from sklearn.metrics.pairwise import cosine_similarity
10 import pickle
11 from tensorflow.keras.models import load_model
12 import gradio as gr
```

```
1 orders = pd.read_csv("/content/drive/MyDrive/Recommender_Project/data/orders_small.csv")
2 order_products = pd.read_csv("/content/drive/MyDrive/Recommender_Project/data/order_products_small.csv")
3 products = pd.read_csv("/content/drive/MyDrive/Recommender_Project/data/products_small.csv")
4 aisles = pd.read_csv("/content/drive/MyDrive/Recommender_Project/data/aisles.csv")
5 departments = pd.read_csv("/content/drive/MyDrive/Recommender_Project/data/departments.csv")
```

```
1 orders.head()
```

	order_id	user_id	eval_set	order_number	order_dow	order_hour_of_day	days_since_prior_order	
0	361493	27	prior	1	3		9	NaN
1	1662354	27	prior	2	2		17	6.0
2	965677	27	prior	3	3		8	1.0
3	2842504	27	prior	4	5		13	2.0
4	1007361	27	prior	5	2		15	4.0

```
1 products.head()
```

	product_id	product_name	aisle_id	department_id	
0	1	Chocolate Sandwich Cookies	61	19	
1	2	All-Seasons Salt	104	13	
2	3	Robust Golden Unsweetened Oolong Tea	94	7	
3	8	Cut Russet Potatoes Steam N' Mash	116	1	
4	10	Sparkling Orange Juice & Prickly Pear Beverage	115	7	

Next steps: [Generate code with products](#) [New interactive sheet](#)

```
1 order_products.head()
```

	order_id	product_id	add_to_cart_order	reordered	
0	4	46842	1	0	
1	4	26434	2	1	
2	4	39758	3	1	
3	4	27761	4	1	
4	4	10054	5	1	

```
1 aisles.head()
```



	aisle_id	aisle	grid
0	1	prepared soups salads	grid
1	2	specialty cheeses	grid
2	3	energy granola bars	grid
3	4	instant foods	grid
4	5	marinades meat preparation	grid

Next steps: [Generate code with aisles](#) [New interactive sheet](#)

1	departments.head()	grid
0	department_id	department
1	1	frozen
2	2	other
3	3	bakery
4	4	produce
5	5	alcohol

Next steps: [Generate code with departments](#) [New interactive sheet](#)

## Collaborative Filtering

```

1 # --- Build unique (user, product) interactions ---
2 # Merge orders and order_products to get user-product interactions
3 user_products = orders[['order_id','user_id']].merge(
4     order_products[['order_id','product_id']], on='order_id'
5 )
6
7 # Map to contiguous indices for embeddings
8 user_ids = user_products['user_id'].unique()
9 product_ids = user_products['product_id'].unique()
10
11 user_map = {uid: i for i, uid in enumerate(user_ids)}
12 product_map = {pid: i for i, pid in enumerate(product_ids)}
13
14 user_products['user_idx'] = user_products['user_id'].map(user_map)
15 user_products['product_idx'] = user_products['product_id'].map(product_map)
16
17 # Get unique positive pairs (user_idx, product_idx)
18 positives = user_products[['user_idx','product_idx']].drop_duplicates().values
19 print("Unique positives:", len(positives))
20
21 # --- Negative sampling ---
22 # For each positive interaction, sample a fixed number of negative interactions (user, unpurchased item)
23 n_neg = 4
24 num_users = len(user_map)
25 num_items = len(product_map)
26
27 pairs = []
28 # Create a dictionary mapping user index to a set of item indices they have purchased
29 user_pos = user_products.groupby('user_idx')['product_idx'].apply(set).to_dict()
30
31 # Iterate through positive interactions and sample negatives
32 for u, i in tqdm(positives, desc="Sampling negatives"):
33     pairs.append((u, i, 1)) # Add the positive interaction with label 1
34     negs = 0
35     while negs < n_neg:
36         # Randomly sample an item index
37         j = np.random.randint(0, num_items)
38         # Check if the sampled item has not been purchased by the user
39         if j not in user_pos[u]:
40             pairs.append((u, j, 0)) # Add the negative interaction with label 0
41             negs += 1
42
43 # Convert the list of pairs to a numpy array
44 pairs = np.array(pairs, dtype=np.int32)

```

```

45 # Separate users, items, and labels
46 users_arr, items_arr, labels_arr = pairs[:,0], pairs[:,1], pairs[:,2].astype(np.float32)
47
48 # Shuffle the data and split into training and validation sets
49 perm = np.random.permutation(len(users_arr))
50 users_arr, items_arr, labels_arr = users_arr[perm], items_arr[perm], labels_arr[perm]
51
52 split = int(0.9 * len(users_arr))
53 users_train, items_train, labels_train = users_arr[:split], items_arr[:split], labels_arr[:split]
54 users_val, items_val, labels_val = users_arr[split:], items_arr[split:], labels_arr[split:]
55
56 print("Train samples:", len(users_train), "Val samples:", len(users_val))
57
58 # Build tf.data.Dataset for efficient training
59 batch_size = 4096
60 train_ds = tf.data.Dataset.from_tensor_slices(((users_train, items_train), labels_train))\
61     .shuffle(100000).batch(batch_size).prefetch(tf.data.AUTOTUNE)
62 val_ds = tf.data.Dataset.from_tensor_slices(((users_val, items_val), labels_val))\
63     .batch(batch_size).prefetch(tf.data.AUTOTUNE)
64
65 # Print the number of unique users and items
66 num_users, num_items

```

Unique positives: 1637937  
Sampling negatives: 100%|██████████| 1637937/1637937 [00:29<00:00, 56283.46it/s]  
Train samples: 7370716 Val samples: 818969  
(10000, 20000)

```

1 # === TensorFlow CF model (MF) + training ===
2
3 # Define hyperparameters
4 embedding_dim = 64      # Dimension of user and item embeddings; try 32, 64, or 128
5 lr = 1e-3                # Learning rate for the optimizer
6 epochs = 6                # Number of training epochs
7 model_save_path = "tf_cf_mf.h5" # Path to save the best model during training
8
9 # --- Model: simple Matrix Factorization (dot product + bias) ---
10 class MatrixFactorization(Model):
11     def __init__(self, num_users, num_items, embed_dim):
12         super().__init__()
13         # Embedding layers for users and items
14         self.user_emb = layers.Embedding(num_users, embed_dim, embeddings_initializer="he_normal")
15         self.item_emb = layers.Embedding(num_items, embed_dim, embeddings_initializer="he_normal")
16         # Bias layers for users and items
17         self.user_bias = layers.Embedding(num_users, 1, embeddings_initializer="zeros")
18         self.item_bias = layers.Embedding(num_items, 1, embeddings_initializer="zeros")
19         # Sigmoid activation to output a probability-like score
20         self.sig = layers.Activation("sigmoid")
21
22     def call(self, inputs):
23         user_idx, item_idx = inputs # Input tensors: user and item indices (batch,)
24         u = self.user_emb(user_idx)    # User embeddings (batch, embed_dim)
25         v = self.item_emb(item_idx)   # Item embeddings (batch, embed_dim)
26         ub = tf.squeeze(self.user_bias(user_idx), axis=-1) # User biases (batch,)
27         ib = tf.squeeze(self.item_bias(item_idx), axis=-1) # Item biases (batch,)
28         dot = tf.reduce_sum(u * v, axis=1)                  # Dot product of embeddings (batch,)
29         x = dot + ub + ib                                # Add biases to the dot product
30         out = self.sig(x)                                 # Apply sigmoid activation
31         return out
32
33 # Instantiate the model
34 model = MatrixFactorization(num_users, num_items, embedding_dim)
35
36 # Compile the model
37 model.compile(
38     optimizer=optimizers.Adam(learning_rate=lr),
39     loss=losses.BinaryCrossentropy(), # Use Binary Crossentropy for implicit feedback (binary classification)
40     metrics=[tf.keras.metrics.BinaryAccuracy(name="acc")]) # Metric to monitor during training
41 )
42
43 # Define callbacks for training
44 es = callbacks.EarlyStopping(monitor="val_loss", patience=2, restore_best_weights=True) # Stop training if validation loss doesn't decrease
45 ck = callbacks.ModelCheckpoint(model_save_path, save_best_only=True, monitor="val_loss") # Save the best model based on validation loss
46
47 # Train the model
48 history = model.fit(train_ds, validation_data=val_ds, epochs=epochs, callbacks=[es, ck])
49

```

```

50 # Save the final model and extract embeddings
51 model.save(model_save_path, include_optimizer=False)
52 user_embeddings = model.user_emb.get_weights()[0] # Extract user embeddings
53 item_embeddings = model.item_emb.get_weights()[0] # Extract item embeddings
54
55 print("Training done. Embedding shapes:", user_embeddings.shape, item_embeddings.shape)
56
57
58 # === Recommend helper: returns product_ids and scores ===
59 # Build product_index_to_id array (col index -> product_id) from product_map
60 product_index_to_id = np.array(list(product_map.keys()))
61
62 # Build user -> set(item_idx) for masking (do this once)
63 user_seen = orders[['order_id','user_id']].merge(order_products[['order_id','product_id']], on='order_id')
64 user_seen['user_idx'] = user_seen['user_id'].map(user_map)
65 user_seen['item_idx'] = user_seen['product_id'].map(product_map)
66 user_seen = user_seen.dropna(subset=['user_idx','item_idx']).astype({'user_idx':int,'item_idx':int})
67 user_seen_map = user_seen.groupby('user_idx')['item_idx'].apply(set).to_dict()
68
69
70

```

```

Epoch 1/6
1800/1800 ━━━━━━━━━━ 0s 28ms/step - acc: 0.8446 - loss: 0.5221WARNING:absl:You are saving your model as an HDF5 file via
1800/1800 ━━━━━━━━━━ 55s 30ms/step - acc: 0.8446 - loss: 0.5220 - val_acc: 0.8616 - val_loss: 0.3261
Epoch 2/6
1797/1800 ━━━━━━ 0s 29ms/step - acc: 0.8646 - loss: 0.3184WARNING:absl:You are saving your model as an HDF5 file via
1800/1800 ━━━━━━ 55s 30ms/step - acc: 0.8646 - loss: 0.3184 - val_acc: 0.8698 - val_loss: 0.3033
Epoch 3/6
1796/1800 ━━━━━━ 0s 29ms/step - acc: 0.8718 - loss: 0.2964WARNING:absl:You are saving your model as an HDF5 file via
1800/1800 ━━━━━━ 54s 30ms/step - acc: 0.8718 - loss: 0.2964 - val_acc: 0.8730 - val_loss: 0.2921
Epoch 4/6
1797/1800 ━━━━━━ 0s 27ms/step - acc: 0.8764 - loss: 0.2833WARNING:absl:You are saving your model as an HDF5 file via
1800/1800 ━━━━━━ 51s 28ms/step - acc: 0.8764 - loss: 0.2833 - val_acc: 0.8775 - val_loss: 0.2817
Epoch 5/6
1797/1800 ━━━━━━ 0s 28ms/step - acc: 0.8831 - loss: 0.2687WARNING:absl:You are saving your model as an HDF5 file via
1800/1800 ━━━━━━ 84s 29ms/step - acc: 0.8831 - loss: 0.2687 - val_acc: 0.8816 - val_loss: 0.2718
Epoch 6/6
1798/1800 ━━━━━━ 0s 29ms/step - acc: 0.8897 - loss: 0.2533WARNING:absl:You are saving your model as an HDF5 file via
1800/1800 ━━━━━━ 54s 30ms/step - acc: 0.8897 - loss: 0.2533 - val_acc: 0.8843 - val_loss: 0.2645
WARNING:absl:You are saving your model as an HDF5 file via `model.save()` or `keras.saving.save_model(model)`. This file format is
Training done. Embedding shapes: (10000, 64) (20000, 64)

```

```

1 def recommend_for_user_CF(user_id, top_k=10, return_scores=True):
2     # Check if the user ID exists in our mapping
3     if user_id not in user_map:
4         raise ValueError(f"User ID {user_id} not found in the user map.")
5
6     u_idx = user_map[user_id]           # Get the internal index for the user
7     u_vec = user_embeddings[u_idx]      # Get the user's embedding vector (embed_dim,)
8
9     # Calculate scores for all items by taking the dot product of user embedding with all item embeddings
10    scores = item_embeddings @ u_vec   # (num_items,)
11
12    # Mask items the user has already seen (purchased)
13    seen = user_seen_map.get(u_idx, set()) # Get the set of item indices the user has seen
14    scores[list(seen)] = -np.inf        # Set scores of seen items to negative infinity to exclude them from top-k
15
16    # Get the indices of the top-k highest scores
17    top_idx = np.argpartition(-scores, range(top_k))[:top_k]
18
19    # Sort the top-k indices by their scores in descending order
20    top_idx = top_idx[np.argsort(-scores[top_idx])]
21
22    # Get the original product IDs for the top-k indices
23    prod_ids = product_index_to_id[top_idx]
24
25    # Return product IDs and scores if requested, otherwise just product IDs
26    if return_scores:
27        return list(prod_ids), list(scores[top_idx])
28    return list(prod_ids)
29

```

```

1 # Example usage: recommend for user 27 (real Instacart ID)
2 example_uid = 27
3 pids, sc = recommend_for_user_CF(example_uid, top_k=10)
4 print("TF-MF recommendations (product_id, score):"

```

```

5 # Iterate through the recommended product IDs and their scores
6 for pid, s in zip(pids, sc):
7     # Find the product name using the product_id from the products DataFrame
8     name = products.loc[products['product_id']==pid,'product_name'].values[0]
9     # Print the product name, ID, and recommendation score
10    print(f"- {name} (id={pid}) score={s:.4f}")

TF-MF recommendations (product_id, score):
- Sparkling Lemon Water (id=21709) score=4.1237
- Clementines, Bag (id=28199) score=3.7631
- Organic Whole String Cheese (id=22035) score=3.7572
- Peach Pear Flavored Sparkling Water (id=26620) score=3.7144
- Vitamin C Super Orange Dietary Supplement (id=2238) score=3.6471
- Kiwi Sandia Sparkling Water (id=12576) score=3.5266
- Dark Chocolate Nuts & Sea Salt Bars (id=45645) score=3.5177
- Earl Grey Tea (id=48946) score=3.4537
- Organic Reduced Fat 2% Milk (id=5785) score=3.4254
- Hass Avocados (id=12341) score=3.4170

```

## Content Based Filtering

```

1 # Merge products with aisle and department names to create a richer product description
2 products_full = products.merge(aisles, on="aisle_id").merge(departments, on="department_id")
3
4 # Create an enriched text column by concatenating product name, aisle, and department
5 # This text will be used to generate content-based embeddings
6 products_full['full_text'] = (
7     products_full['product_name'].fillna("") + # Use fillna("") to handle potential missing values gracefully
8     " [aisle: " + products_full['aisle'].fillna("") + "] " +
9     " [dept: " + products_full['department'].fillna("") + "]"
10 )
11
12 # Display the head of the dataframe to show the new 'full_text' column
13 print(products_full[['product_name','full_text']].head(5))

```

	product_name \
0	Chocolate Sandwich Cookies
1	All-Seasons Salt
2	Robust Golden Unsweetened Oolong Tea
3	Cut Russet Potatoes Steam N' Mash
4	Sparkling Orange Juice & Prickly Pear Beverage

  

	full_text
0	Chocolate Sandwich Cookies [aisle: cookies cak...
1	All-Seasons Salt [aisle: spices seasonings] [d...
2	Robust Golden Unsweetened Oolong Tea [aisle: t...
3	Cut Russet Potatoes Steam N' Mash [aisle: froz...
4	Sparkling Orange Juice & Prickly Pear Beverage...

```

1 # Load a pre-trained sentence transformer model
2 # This model will be used to convert text descriptions into numerical vectors (embeddings)
3 st_model = SentenceTransformer('all-MiniLM-L6-v2')
4
5 # Prepare the product texts and their corresponding product IDs
6 product_texts = products_full['full_text'].tolist()
7 product_ids = products_full['product_id'].tolist()
8
9 # Encode the enriched product texts into embeddings
10 # show_progress_bar=True displays a progress bar during encoding
11 # convert_to_numpy=True ensures the output is a NumPy array
12 # normalize_embeddings=True scales the embeddings to unit length, for cosine similarity
13 item_embeddings_content = st_model.encode(
14     product_texts,
15     show_progress_bar=True,
16     convert_to_numpy=True,
17     normalize_embeddings=True
18 )
19
20 # Print the shape of the generated embeddings
21 print("Embeddings shape:", item_embeddings_content.shape)

```

Batches: 100% 625/625 [00:06<00:00, 126.93it/s]  
 Embeddings shape: (20000, 384)

```

1 # Map product_id -> row index in embedding matrix
2 product_id_to_index = {pid: i for i, pid in enumerate(product_ids)}

1 def recommend_for_user_CBF(user_id, top_k=10):
2     # Find the orders placed by the user
3     user_orders = orders[orders['user_id']==user_id]['order_id']
4
5     # Find all unique products purchased by the user across their orders
6     user_products_hist = order_products[order_products['order_id'].isin(user_orders)]['product_id'].unique()
7
8     # If the user has no purchase history, return an empty list
9     if len(user_products_hist) == 0:
10         return []
11
12     # Get the indices of the purchased items in the content-based embedding matrix
13     # Filter out products that might not be in the products_full dataframe used for embeddings
14     hist_indices = [product_id_to_index[pid] for pid in user_products_hist if pid in product_id_to_index]
15
16     # Calculate the user's profile by averaging the content embeddings of their purchased items
17     user_profile = item_embeddings_content[hist_indices].mean(axis=0).reshape(1, -1)
18
19     # Compute the cosine similarity between the user profile and all item embeddings
20     sims = cosine_similarity(user_profile, item_embeddings_content)[0]
21
22     # Mask already purchased items by setting their similarity scores to a very low value
23     for idx in hist_indices:
24         sims[idx] = -1 # Use -1 as cosine similarity is in [-1, 1]
25
26     # Get the indices of the top-k items with the highest similarity scores
27     top_idx = np.argpartition(-sims, range(top_k))[:top_k]
28
29     # Sort the top-k indices based on their similarity scores in descending order
30     top_idx = top_idx[np.argsort(-sims[top_idx])]
31
32     # Get the original product IDs and their corresponding similarity scores for the top-k indices
33     rec_ids = [product_ids[i] for i in top_idx]
34     rec_scores = [sims[i] for i in top_idx]
35
36     # Return the recommendations as a list of (product_id, score) tuples
37     return list(zip(rec_ids, rec_scores))

```

```

1 # Example: get content-based recommendations for user 27
2 recs = recommend_for_user_CBF(27, top_k=10)
3
4 print("Content-based recommendations for user 27 (with aisle+dept):")
5 # Iterate through the recommendations and print details
6 for pid, score in recs:
7     # Find the product details (name, aisle, department) using the product_id
8     product_info = products_full.loc[products_full['product_id']==pid].iloc[0]
9     name = product_info['product_name']
10    aisle = product_info['aisle']
11    dept = product_info['department']
12    # Print the recommendation with its details and score
13    print(f"- {name} | aisle={aisle}, dept={dept} (sim={score:.3f})")

```

Content-based recommendations for user 27 (with aisle+dept):

- Organic Nuts & Chocolate Beverage | aisle=nuts seeds dried fruit, dept=snacks (sim=0.793)
- Strawberry | aisle=yogurt, dept=dairy eggs (sim=0.784)
- Organic Jelly Beans | aisle=candy chocolate, dept=snacks (sim=0.782)
- Fruit Medley, Organic | aisle=fruit vegetable snacks, dept=snacks (sim=0.780)
- Nutri-Grain Strawberry | aisle=breakfast bars pastries, dept=breakfast (sim=0.776)
- Yogurt, Lowfat, Organic, Vanilla Bean | aisle=yogurt, dept=dairy eggs (sim=0.766)
- Organic Strawberry Energy Chews | aisle=candy chocolate, dept=snacks (sim=0.763)
- Organic Chunky Strawberry Fruit Bars | aisle=ice cream ice, dept=frozen (sim=0.763)
- Fruit Spread, Organic, Strawberry | aisle=spreads, dept=pantry (sim=0.763)
- Organic Mixed Berry Snack Bar | aisle=energy granola bars, dept=snacks (sim=0.760)

```

1 def recommend_hybrid(user_id, top_k=10, alpha=0.6):
2     """
3         Generates hybrid recommendations for a user by combining CF and CBF scores.
4
5     Args:
6         user_id (int): The ID of the user for whom to generate recommendations.
7         top_k (int): The number of top recommendations to return.
8         alpha (float): The weighting factor for CF scores (0.0 to 1.0).
9             A higher alpha gives more weight to collaborative filtering.

```

```

10
11     Returns:
12         list of tuples: A list of (product_id, hybrid_score) for the top_k recommended products.
13             Returns an empty list if the user has no purchase history.
14 """
15 # Check if the user ID exists in our mapping (important for CF part)
16 if user_id not in user_map:
17     # If user not in CF model, fallback to pure CBF or popularity if no history
18     print(f"Warning: User ID {user_id} not found in CF user map. Using only CBF if history exists.")
19     # Proceed to CBF part, which handles users with no history
20     return recommend_for_user_CBF(user_id, top_k=top_k) # CBF returns (pid, sim)
21
22 u_idx = user_map[user_id]
23
24 # --- CF scores ---
25 # Get the user's embedding from the CF model
26 u_vec_cf = user_embeddings[u_idx]
27 # Calculate CF scores for all items using the dot product
28 scores_cf = item_embeddings @ u_vec_cf # (num_items,)
29
30 # --- CBF scores ---
31 # Find the orders placed by the user to get purchase history for CBF
32 user_orders = orders[orders['user_id']==user_id]['order_id']
33 user_products_hist = order_products[order_products['order_id'].isin(user_orders)][['product_id']].unique()
34
35 # We'll return empty if no history for CBF, relying on the initial check for CF.
36 if len(user_products_hist) == 0:
37     print(f"Warning: User ID {user_id} has no purchase history for CBF.")
38     return []
39
40 # Get the indices of the purchased items in the content-based embedding matrix
41 # Filter out products that might not be in the products_full dataframe used for embeddings
42 hist_indices = [product_id_to_index[pid] for pid in user_products_hist if pid in product_id_to_index]
43
44 # Calculate the user's content-based profile by averaging embeddings of purchased items
45 user_profile = item_embeddings_content[hist_indices].mean(axis=0).reshape(1, -1)
46 # Compute cosine similarity between the user profile and all item content embeddings
47 scores_cbf = cosine_similarity(user_profile, item_embeddings_content)[0]
48
49 # --- Combine scores ---
50 # Calculate the final hybrid score as a weighted sum of CF and CBF scores
51 scores = alpha * scores_cf + (1 - alpha) * scores_cbf
52
53 # Mask already purchased items by setting their scores to negative infinity
54 # This prevents recommending items the user has already bought
55 seen = set(hist_indices) # Use indices from CBF history as a proxy for seen items
56 for idx in seen:
57     scores[idx] = -np.inf
58
59 # Get the indices of the top-k items with the highest hybrid scores
60 top_idx = np.argpartition(-scores, range(top_k))[:top_k]
61
62 # Sort the top-k indices based on their hybrid scores in descending order
63 top_idx = top_idx[np.argsort(-scores[top_idx])]

64 # Get the original product IDs and their corresponding hybrid scores for the top-k indices
65 rec_ids = [product_ids[i] for i in top_idx]
66 rec_scores = [scores[i] for i in top_idx]
67
68 # Return the recommendations as a list of (product_id, score) tuples
69 return list(zip(rec_ids, rec_scores))

```

```

1 # Example: get hybrid recommendations for user 27 with alpha=0.6
2 recs = recommend_hybrid(27, top_k=10, alpha=0.6)
3
4 print("Hybrid recommendations for user 27:")
5 # Iterate through the recommendations and print details
6 for pid, score in recs:
7     # Find the product details (name, aisle, department) using the product_id from products_full
8     product_info = products_full.loc[products_full['product_id']==pid,'product_name'].values[0]
9     aisle = products_full.loc[products_full['product_id']==pid,'aisle'].values[0]
10    dept = products_full.loc[products_full['product_id']==pid,'department'].values[0]
11    # Print the recommendation with its details and hybrid score
12    print(f"- {product_info} | aisle={aisle}, dept={dept} (hybrid_score={score:.3f})")

```

```
Hybrid recommendations for user 27:
- Snack Bags | aisle=more household, dept=household (hybrid_score=2.891)
- Minis Original Saltine Crackers | aisle=crackers, dept=snacks (hybrid_score=2.788)
- Peach Sparkling Energy Water | aisle=energy sports drinks, dept=beverages (hybrid_score=2.723)
- Berry Punch Flavored Fruit Drink | aisle=refrigerated, dept=beverages (hybrid_score=2.504)
- Caramel Sauce | aisle=ice cream toppings, dept=snacks (hybrid_score=2.501)
- Organic Rice Vinegar | aisle=asian foods, dept=international (hybrid_score=2.478)
- Bean and Cheese Gorditas | aisle=prepared meals, dept=deli (hybrid_score=2.474)
- Yerba Mate Sparkling Classic Gold | aisle=tea, dept=beverages (hybrid_score=2.464)
- 24/7 Performance Cat Litter | aisle=cat food care, dept=pets (hybrid_score=2.447)
- Organic Yukon Gold Potato Bag | aisle=fresh vegetables, dept=produce (hybrid_score=2.385)
```

```
1 class HybridRecommender:
2     """
3         A class to combine Collaborative Filtering (CF) and Content-Based Filtering (CBF)
4         recommendations.
5     """
6     def __init__(self, cf_model=None, user_embeddings=None, item_embeddings_cf=None,
7                  item_embeddings_cbf=None, product_ids=None,
8                  user_map=None, product_map=None, product_id_to_index=None,
9                  alpha=0.6):
10    """
11        Initializes the HybridRecommender with necessary models, embeddings, and mappings.
12
13    Args:
14        cf_model (tf.keras.Model, optional): The trained Collaborative Filtering model. Defaults to None.
15        user_embeddings (np.ndarray, optional): User embeddings from the CF model. Defaults to None.
16        item_embeddings_cf (np.ndarray, optional): Item embeddings from the CF model. Defaults to None.
17        item_embeddings_cbf (np.ndarray, optional): Item embeddings from the CBF model. Defaults to None.
18        product_ids (np.ndarray, optional): Array of original product IDs. Defaults to None.
19        user_map (dict, optional): Mapping from original user_id to internal index. Defaults to None.
20        product_map (dict, optional): Mapping from original product_id to internal index (for CF). Defaults to None.
21        product_id_to_index (dict, optional): Mapping from original product_id to internal index (for CBF). Defaults to None.
22        alpha (float, optional): Weighting factor for CF vs CBF scores (0.0 to 1.0). Defaults to 0.6.
23    """
24     self.cf_model = cf_model
25     self.user_embeddings = user_embeddings
26     self.item_embeddings_cf = item_embeddings_cf
27     self.item_embeddings_cbf = item_embeddings_cbf
28     self.product_ids = product_ids
29     self.user_map = user_map
30     self.product_map = product_map # Keep if needed for other CF related lookups
31     self.product_id_to_index = product_id_to_index # For CBF index lookups
32     self.alpha = alpha
33
34     # -----
35     # Save and Load
36     # -----
37     def save(self, path="recommender"):
38         """
39             Saves the recommender components to disk.
40
41         Args:
42             path (str, optional): Base path for saving files. Defaults to "recommender".
43         """
44             # Save embeddings
45             np.save(f"{path}_user_embeddings.npy", self.user_embeddings)
46             np.save(f"{path}_item_embeddings_cf.npy", self.item_embeddings_cf)
47             np.save(f"{path}_item_embeddings_cbf.npy", self.item_embeddings_cbf)
48             np.save(f"{path}_product_ids.npy", self.product_ids)
49
50             # Save mappings and alpha
51             with open(f"{path}_mappings.pkl", "wb") as f:
52                 pickle.dump({
53                     "user_map": self.user_map,
54                     "product_map": self.product_map,
55                     "product_id_to_index": self.product_id_to_index,
56                     "alpha": self.alpha
57                 }, f)
58
59             # Save CF model (optional, depends on if you need to load the full model later)
60             if self.cf_model is not None:
61                 try:
62                     self.cf_model.save(f"{path}_cf_model.h5", include_optimizer=False)
63                 except Exception as e:
64                     print(f"Warning: Could not save TF model: {e}")
65
```

```

66
67     @classmethod
68     def load(cls, path="recommender"):
69         """
70             Loads the recommender components from disk.
71
72         Args:
73             path (str, optional): Base path where files were saved. Defaults to "recommender".
74
75         Returns:
76             HybridRecommender: An instance of the HybridRecommender.
77         """
78
79         # Load embeddings
80         user_embeddings = np.load(f"{path}_user_embeddings.npy")
81         item_embeddings_cf = np.load(f"{path}_item_embeddings_cf.npy")
82         item_embeddings_cbf = np.load(f"{path}_item_embeddings_cbf.npy")
83         product_ids = np.load(f"{path}_product_ids.npy")
84
85         # Load mappings
86         with open(f'{path}_mappings.pkl', "rb") as f:
87             maps = pickle.load(f)
88
89         # Load CF model
90         cf_model = None
91         try:
92             # compile=False to avoid needing the original training setup
93             cf_model = load_model(f"{path}_cf_model.h5", compile=False)
94         except:
95             # Handle cases where the model file might not exist or is corrupted
96             print(f"Warning: Could not load TF model from {path}_cf_model.h5")
97             pass
98
99         return cls(cf_model=cf_model,
100                 user_embeddings=user_embeddings,
101                 item_embeddings_cf=item_embeddings_cf,
102                 item_embeddings_cbf=item_embeddings_cbf,
103                 product_ids=product_ids,
104                 user_map=maps["user_map"],
105                 product_map=maps.get("product_map"), # Use .get for backward compatibility if product_map wasn't saved
106                 product_id_to_index=maps["product_id_to_index"],
107                 alpha=maps.get("alpha", 0.6)) # Use .get for backward compatibility
108
109     # -----
110     # Recommend
111     # -----
112     def recommend(self, user_id, top_k=10, orders=None, order_products=None):
113         """
114             Generates hybrid recommendations for a user.
115
116         Args:
117             user_id (int): The ID of the user for whom to generate recommendations.
118             top_k (int, optional): The number of top recommendations to return. Defaults to 10.
119             orders (pd.DataFrame, optional): DataFrame containing order information. Required for CBF history lookup.
120             order_products (pd.DataFrame, optional): DataFrame containing order product information. Required for CBF history
121
122         Returns:
123             list of tuples: A list of (product_id, hybrid_score) for the top_k recommended products.
124                     Returns an empty list if the user is not in the user map or has no purchase history for CBF.
125         Raises:
126             ValueError: If orders or order_products DataFrames are not provided for CBF calculation.
127         """
128         # Ensure user exists in the CF map (needed for user embedding)
129         if user_id not in self.user_map:
130             print(f"User ID {user_id} not found in user map.")
131         return []
132
133         u_idx = self.user_map[user_id]
134
135         # --- CF scores ---
136         # Calculate CF scores for all items using the dot product of user embedding with all item embeddings
137         u_vec_cf = self.user_embeddings[u_idx]
138         scores_cf = self.item_embeddings_cf @ u_vec_cf # (num_items,)
139
140         # --- CBF scores ---
141         # CBF requires user purchase history, which is in the provided DataFrames
142         if orders is None or order_products is None:

```

```

143         raise ValueError("Need orders and order_products DataFrames to compute CBF")
144
145     # Find the orders placed by the user
146     user_orders = orders[orders['user_id']==user_id]['order_id']
147
148     # Find all unique products purchased by the user across their orders
149     user_products_hist = order_products[order_products['order_id'].isin(user_orders)]['product_id'].unique()
150
151     if len(user_products_hist) == 0:
152         print(f"User ID {user_id} has no purchase history for CBF.")
153         return []
154
155     # Get the indices of the purchased items in the content-based embedding matrix
156     # Filter out products that might not be in the products_full dataframe used for embeddings
157     hist_indices = [self.product_id_to_index[pid] for pid in user_products_hist if pid in self.product_id_to_index]
158
159     # Calculate the user's content-based profile by averaging embeddings of purchased items
160     user_profile = self.item_embeddings_cbf[hist_indices].mean(axis=0).reshape(1, -1)
161     # Compute cosine similarity between the user profile and all item content embeddings
162     scores_cbf = cosine_similarity(user_profile, self.item_embeddings_cbf)[0]
163
164     # --- Combine scores ---
165     # Calculate the final hybrid score as a weighted sum of CF and CBF scores
166     scores = self.alpha * scores_cf + (1 - self.alpha) * scores_cbf
167
168     # Mask already purchased items by setting their scores to negative infinity
169     # This prevents recommending items the user has already bought
170     # Use indices from CBF history as a proxy for seen items
171     seen_indices = set(hist_indices)
172     scores[list(seen_indices)] = -np.inf
173
174     # --- Get Top-K recommendations ---
175     # Get the indices of the top-k items with the highest hybrid scores
176     top_idx = np.argpartition(-scores, range(top_k))[:top_k]
177
178     # Sort the top-k indices based on their hybrid scores in descending order
179     top_idx = top_idx[np.argsort(-scores[top_idx])]
180
181     # Get the original product IDs and their corresponding hybrid scores for the top-k indices
182     rec_ids = [self.product_ids[i] for i in top_idx]
183     rec_scores = [scores[i] for i in top_idx]
184
185     return list(zip(rec_ids, rec_scores))

```

```

1 # Create and save recommender
2 recommender = HybridRecommender(
3     cf_model=model,
4     user_embeddings=user_embeddings,
5     item_embeddings_cf=item_embeddings,
6     item_embeddings_cbf=item_embeddings_content,
7     product_ids=np.array(product_ids),
8     user_map=user_map,
9     product_map=product_map,
10    product_id_to_index=product_id_to_index,
11    alpha=0.6
12 )
13
14 recommender.save("/content/drive/MyDrive/Recommender_Project/models/my_recommender")

```

WARNING:absl:You are saving your model as an HDF5 file via `model.save()` or `keras.saving.save\_model(model)`. This file format is

```

1 # Load recommender
2 recommender_loaded = HybridRecommender.load(
3     "/content/drive/MyDrive/Recommender_Project/models/my_recommender"
4 )
5
6 # --- Precompute most popular products (for fallback) ---
7 popular_items = order_products["product_id"].value_counts().index.tolist()
8
9 def recommend_with_history(user_id, top_k=10):
10    try:
11        user_id = int(user_id)
12
13        # --- Get recommendations ---
14        recs = recommender_loaded.recommend(

```

```

15     user_id, top_k=top_k, orders=orders, order_products=order_products
16 )
17
18 # --- Past purchases (last order) ---
19 last_order = orders[orders["user_id"] == user_id]["order_id"].max()
20 if pd.isna(last_order):
21     past_items = []
22 else:
23     past_items = order_products[order_products["order_id"] == last_order]["product_id"].tolist()
24
25 # --- Build past purchases table ---
26 past_data = []
27 for pid in past_items[:top_k]:
28     row = products_full.loc[products_full["product_id"] == pid].iloc[0]
29     past_data.append([row["product_name"], row["aisle"], row["department"]])
30 past_df = pd.DataFrame(past_data, columns=["Product", "Aisle", "Department"])
31 if past_df.empty:
32     past_df = pd.DataFrame([["No past purchases found.", "-", "-"]],
33                             columns=["Product", "Aisle", "Department"])
34
35 # --- Build recommendations table (with fallback) ---
36 rec_data = []
37 if recs:
38     for pid, score in recs:
39         row = products_full.loc[products_full["product_id"] == pid].iloc[0]
40         rec_data.append([row["product_name"], row["aisle"], row["department"], f"{score:.3f}"])
41 else:
42     rec_data.append(["⚠️ No personalized recs", "-", "-"])
43     rec_data.append(["💡 Popular items:", "-", "-"])
44     for pid in popular_items[:top_k]:
45         row = products_full.loc[products_full["product_id"] == pid].iloc[0]
46         rec_data.append([row["product_name"], row["aisle"], row["department"], "-"])
47 rec_df = pd.DataFrame(rec_data, columns=["Product", "Aisle", "Department", "Score"])
48
49 return past_df, rec_df
50
51 except Exception as e:
52     return pd.DataFrame([f"🔴 Error: {str(e)}", "-", "-"]),
53                             columns=["Product", "Aisle", "Department"]), pd.DataFrame()

```

Warning: Could not load TF model from /content/drive/MyDrive/Recommender\_Project/models/my\_recommender\_cf\_model.h5

```

1 # --- Gradio UI ---
2 with gr.Blocks() as demo:
3     gr.Markdown("# 🛒 Instacart Personalized Recommender")
4     gr.Markdown("Compare a user's **past purchases** with their **new recommendations** (with fallback).")
5
6     with gr.Row():
7         user_id_in = gr.Number(value=27, label="User ID")
8         topk_in = gr.Slider(1, 20, value=10, step=1, label="Top-K Recommendations")
9
10    with gr.Row():
11        past_out = gr.Dataframe(headers=["Product", "Aisle", "Department"], label="Past Purchases", interactive=False)
12        rec_out = gr.Dataframe(headers=["Product", "Aisle", "Department", "Score"], label="Recommended Products", interactive=
13
14    btn = gr.Button("Get Recommendations")
15    btn.click(fn=recommend_with_history, inputs=[user_id_in, topk_in], outputs=[past_out, rec_out])
16
17 demo.launch(share=True) # share=True for public link

```

Colab notebook detected. To show errors in colab notebook, set debug=True in launch()  
\* Running on public URL: <https://98ee3c93e67de01756.gradio.live>

This share link expires in 1 week. For free permanent hosting and GPU upgrades, run `gradio deploy` from the terminal in the working directory.

27 1 20

Past Purchases

Product	Aisle
Natural Artisan Water	water seltzer
Natural Artesian Water	water seltzer
Pure Coconut Water	juice nectars
Sparkling Water Grapefruit	water seltzer

Recommended Products

Aisle	Department	Score
more household	household	2.891
crackers	snacks	2.788
energy sports drink	beverages	2.723
refrigerated	beverages	2.504