# Developing a Content-based Filtering approach for item recommender

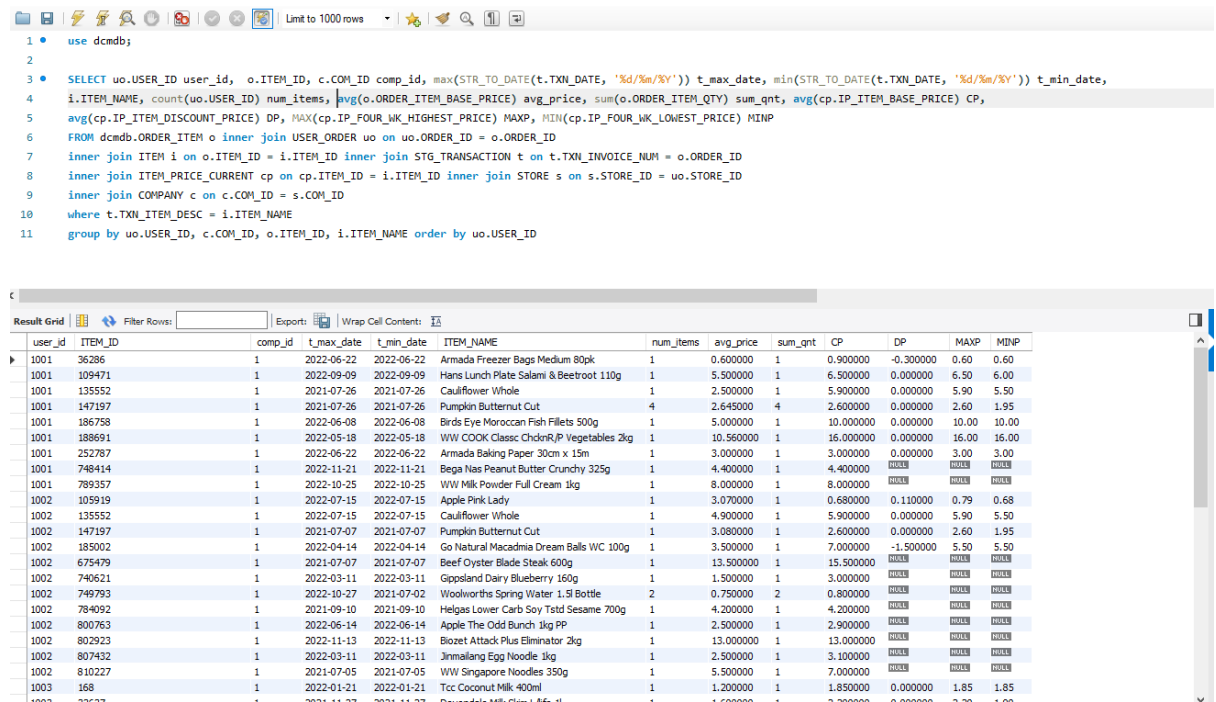- **Advantages of Content-based Filtering approach**
  - Content-based filtering uses item features to recommend other items like what the user likes, based on their previous actions or explicit feedback.
  - Content-based filtering does not require other users' data during recommendations to one user.
  - The model can capture the specific interests of a user, and can recommend niche items that very few other users are interested in.

- **Main components of Content-based Filtering approach**

Recommenders mostly have 3 components:

  - **Candidate Generations**: generating smaller subsets of candidates to recommend to a user, given a huge pool of thousands of items.
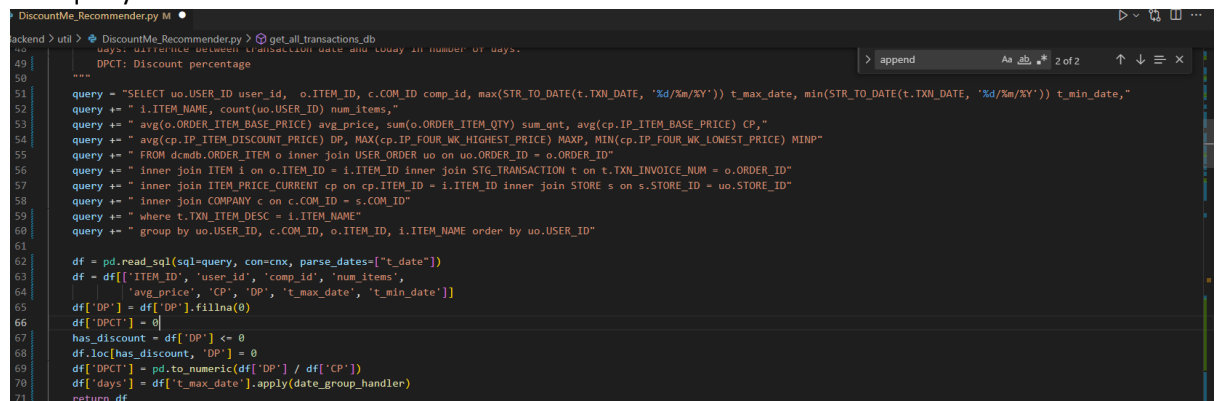
    For this step, I collected user data by fetching the database using the following query.



This query used in code to user data.

- **Scoring Systems**: this step involves scoring to each of the items in the subsets. This is done by the Scoring system.

  One problem to build the scoring system is the unavailability of item rating data which should be in order_item table.

  To solve the problem, I proposed the following algorithm to compute the score for each item in respect to user behaviour:

  1. Rank user items (current ones) based on number of purchases (quantity).

```python
def rank_data(df):
    """ Return the rank for each item per user.
        Args:
            df: all user items
        Returns:
            df_ranked: ranked user items based on number of transactions
    """
    df['rank'] = df.groupby("user_id")["num_items"].rank(
        "dense", ascending=True)
    df_ranked = df[['user_id', 'ITEM_ID', 'comp_id', 'num_items', 'rank']]
    return df_ranked
```

  2. For each current item, compute the similarity with all other items and only conder similarity score > 0.

```python
def compute_similarity_items(all_items, item_id, item_name):
    """ Return the cosine similarity between an item and all items.
        Args:
            all_items: all items on the database
            item_id: target item id
            item_name: target item name
        Returns:
            df_all_items_similarity_scores: data frame store cosine_similarity with all items
    """
    df_all_items_similarity_scores = pd.DataFrame(
        columns=['ITEM_ID', 'ITEM_OTHER', 'SIM'])

    for idx1 in all_items.index:
        sim = _compute_similarity_two_items(
            all_items['ITEM_NAME'][idx1], item_name)
        #add only
        if sim > 0 and all_items['ITEM_ID'][idx1] != item_id:
            data_row = [item_id, all_items['ITEM_ID'][idx1], sim]
            df_all_items_similarity_scores.loc[df_all_items_similarity_scores.size] = data_row


    return df_all_items_similarity_scores
```

  3. Compute the score as a product of rank of related item and similarity score.

```python
# compute score based on rank and sim
df_merged['SCORE'] = df_merged['rank'] * df_merged['SIM']
df_user_item_scores = df_merged[[
    'user_id', 'comp_id', 'ITEM_OTHER', 'SCORE']]
df_user_item_scores.columns = [
    'USER_ID', 'COM_ID', 'ITEM_ID', 'SCORE']
count = min(count, len(df_user_item_scores))

#get top scores
df_user_item_top_scores = df_user_item_scores.sort_values(
    'SCORE', ascending=False).head(count)
```

  4. Select items with top scores.

- **Sample user current items**



- **Recommended items for a user with scoring results.**