# ⌄ Mining Big Data - Assignment 3b

Student Names and Numbers:

Ebubekir Pulat - a1833363

Chunyu Zhang - a1751743

```python
1 import pandas as pd
2 from mlxtend.preprocessing import TransactionEncoder
3 from mlxtend.frequent_patterns import fpgrowth, association_rules
4
5
6 def frequent_patterns(file_path):
7     data = pd.read_csv(file_path)
8     # Group the dataset with member number and date, to get the items in sa
9     transactions = data.groupby(['Member_number', 'Date'])['itemDescription
10    trans_encoder = TransactionEncoder()
11    arr = trans_encoder.fit(transactions).transform(transactions)
12    df = pd.DataFrame(arr, columns = trans_encoder.columns_)
13
14    frequent_itemsets = fpgrowth(df, min_support = 0.001, use_colnames = Tr
15
16    return frequent_itemsets
```

```python
1 # My teammate needs the output includes frequent itemsest, alongside their
2 def get_martrix(frequent_itemsets):
3     output_matrix = [[list(itemsets), support] for itemsets, support in zip
4     # Get the sorted list, so that it will be easy to analyze
5     output_matrix = sorted(output_matrix, key = lambda x: x[1], reverse = 1
6     return output_matrix
```

```python
1 def get_rules(frequent_itemsets):
2     rules = association_rules(frequent_itemsets, metric="confidence", min_t
3     rules['interest'] = rules['confidence'] – rules['consequent support']
4     return rules
```

```
1 def get_conf_supp(frequent_itemsets, rules):
2   results = dict()
3   for index, rule in rules.iterrows():
4     (item1,) = rule['antecedents']
5     (item2,) = rule['consequents']
6     itemsets = item1 + ', ' + item2
7     if (item2 + ', ' + item1) not in results.keys():
8       results[itemsets] = [rule['support'], rule['confidence']]
9     else:
10      if results[((item2 + ', ' + item1))][0] < rule['confidence']:
11        results[itemsets] = [rule['support'], rule['confidence']]
12        del results[((item2 + ', ' + item1))]
13      else:
14        continue
15  results = sorted(results.items(), key=lambda item: item[1][0], reverse =
16  return results
```

```
1 def get_interest(rules):
2     interest = []
3     for index, rule in rules.iterrows():
4         (item1,) = rule['antecedents']
5         (item2,) = rule['consequents']
6         itemsets = [item1, item2]
7         interest.append([itemsets, rule["interest"]])
8     return interest
```

```
1 # Use this function to print the frequent
2 def print_frequent(results):
3     for item in results:
4         print(f'patterns:[{item[0]}], support: {item[1][0]}, confidence: {ite
```

```
1 file_path = 'Groceries data train.csv'
2
3 # Print the results after the fp-growth
4 frequent_itemsets = frequent_patterns(file_path)
5 print(frequent_itemsets)
6
7 # Print the results
8 rules = get_rules(frequent_itemsets)
9 results = get_conf_supp(frequent_itemsets, rules)
10 print("\nThe results: ")
11 print_frequent(results)
12
13 interest = get_interest(rules)
```

```
patterns:[frozen vegetables, rolls/buns], support: 0.0015106826846989425,
patterns:[waffles, whole milk], support: 0.0015106826846989425, confidenc
patterns:[curd, tropical fruit], support: 0.0014387454139989928, confiden
patterns:[curd, rolls/buns], support: 0.0014387454139989928, confidence:
patterns:[chocolate, rolls/buns], support: 0.0014387454139989928, confider
patterns:[hamburger meat, whole milk], support: 0.0014387454139989928, co
patterns:[domestic eggs, sausage], support: 0.0014387454139989928, confide
patterns:[coffee, rolls/buns], support: 0.0014387454139989928, confidence
patterns:[onions, whole milk], support: 0.0014387454139989928, confidence
patterns:[berries, whole milk], support: 0.0014387454139989928, confidence
patterns:[butter, bottled water], support: 0.0013668081432990432, confider
patterns:[coffee, yogurt], support: 0.0013668081432990432, confidence: 0.(
patterns:[oil, whole milk], support: 0.0013668081432990432, confidence: 0.
patterns:[sugar, whole milk], support: 0.0012948708725990935, confidence:
patterns:[dessert, whole milk], support: 0.0012948708725990935, confidence
patterns:[coffee, other vegetables], support: 0.0012948708725990935, confi
patterns:[long life bakery product, whole milk], support: 0.0012948708725(
patterns:[napkins, other vegetables], support: 0.0012948708725990935, con
patterns:[napkins, whole milk], support: 0.0012948708725990935, confidence
patterns:[frozen vegetables, soda], support: 0.0012229336018991439, confio
patterns:[frozen vegetables, sausage], support: 0.0012229336018991439, co
patterns:[UHT-milk, other vegetables], support: 0.0012229336018991439, co
patterns:[candy, whole milk], support: 0.0012229336018991439, confidence:
patterns:[beverages, whole milk], support: 0.0012229336018991439, confider
patterns:[salty snack, rolls/buns], support: 0.0012229336018991439, confio
patterns:[frozen vegetables, root vegetables], support: 0.001150996331199:
patterns:[hamburger meat, soda], support: 0.0011509963311991942, confidenc
patterns:[grapes, soda], support: 0.0011509963311991942, confidence: 0.10!
patterns:[processed cheese, rolls/buns], support: 0.0011509963311991942, (
patterns:[meat, other vegetables], support: 0.0011509963311991942, confide
patterns:[napkins, pastry], support: 0.0011509963311991942, confidence: 0.
patterns:[salty snack, whole milk], support: 0.0011509963311991942, confio
patterns:[waffles, rolls/buns], support: 0.0011509963311991942, confidence
patterns:[white bread, rolls/buns], support: 0.0010790590604992446, confio
patterns:[chicken, other vegetables], support: 0.0010790590604992446, con
patterns:[hamburger meat, other vegetables], support: 0.001079059060499924
patterns:[hamburger meat, rolls/buns], support: 0.0010790590604992446, co
patterns:[specialty bar, whole milk], support: 0.0010790590604992446, con
patterns:[grapes, whole milk], support: 0.0010790590604992446, confidence
patterns:[frozen meals, soda], support: 0.0010790590604992446, confidence
patterns:[napkins, soda], support: 0.0010790590604992446, confidence: 0.0(
patterns:[ham, other vegetables], support: 0.0010790590604992446, confider
patterns:[berries, other vegetables], support: 0.0010790590604992446, con
patterns:[hard cheese, rolls/buns], support: 0.0010790590604992446, confio
patterns:[hygiene articles, whole milk], support: 0.0010071217897992951, (
patterns:[chocolate, sausage], support: 0.0010071217897992951, confidence
patterns:[chocolate, yogurt], support: 0.0010071217897992951, confidence:
patterns:[dessert, rolls/buns], support: 0.0010071217897992951, confidence
patterns:[hamburger meat, yogurt], support: 0.0010071217897992951, confide
patterns:[herbs, whole milk], support: 0.0010071217897992951, confidence:
patterns:[UHT-milk, tropical fruit], support: 0.0010071217897992951, confi
patterns:[grapes, other vegetables], support: 0.0010071217897992951, confi
patterns:[frozen meals, whole milk], support: 0.0010071217897992951, confi
```

```
patterns:[frozen meals, other vegetables], support: 0.00100712178979929510
patterns:[napkins, rolls/buns], support: 0.00100712178979929510, confidenc
patterns:[salty snack, other vegetables], support: 0.00100712178979929510,
patterns:[waffles, other vegetables], support: 0.00100712178979929510, con
```

```python
1  def get_5_output(file_path):
2    frequent_itemsets = frequent_patterns(file_path)
3    rules = get_rules(frequent_itemsets)
4    results = get_conf_supp(frequent_itemsets, rules)
5    print_frequent(results[:5])
6    return results
```

```python
1  # for Trainning dataset
2  print('5 patterns of train dataset:')
3  train_results = get_5_output('Groceries data train.csv')
4
5
6  # For test dataset
7  print('\n5 patterns of test dataset:')
8  test_results = get_5_output('Groceries data test.csv')
```

```
5 patterns of train dataset:
patterns:[rolls/buns, whole milk], support: 0.007913099776994462, confidenc
patterns:[other vegetables, whole milk], support: 0.006977915257895115, con
patterns:[yogurt, whole milk], support: 0.006546291633695417, confidence: 0
patterns:[whole milk, soda], support: 0.00597079346809582, confidence: 0.05
patterns:[other vegetables, soda], support: 0.005467232573196173, confidenc

5 patterns of test dataset:
patterns:[other vegetables, whole milk], support: 0.0033055967172005017, co
patterns:[tropical fruit, whole milk], support: 0.0020517496865382423, conf
patterns:[bottled water, other vegetables], support: 0.0015958053117519663,
patterns:[pastry, whole milk], support: 0.0013678331243588283, confidence:
patterns:[beef, whole milk], support: 0.0011398609369656903, confidence: 0.
```
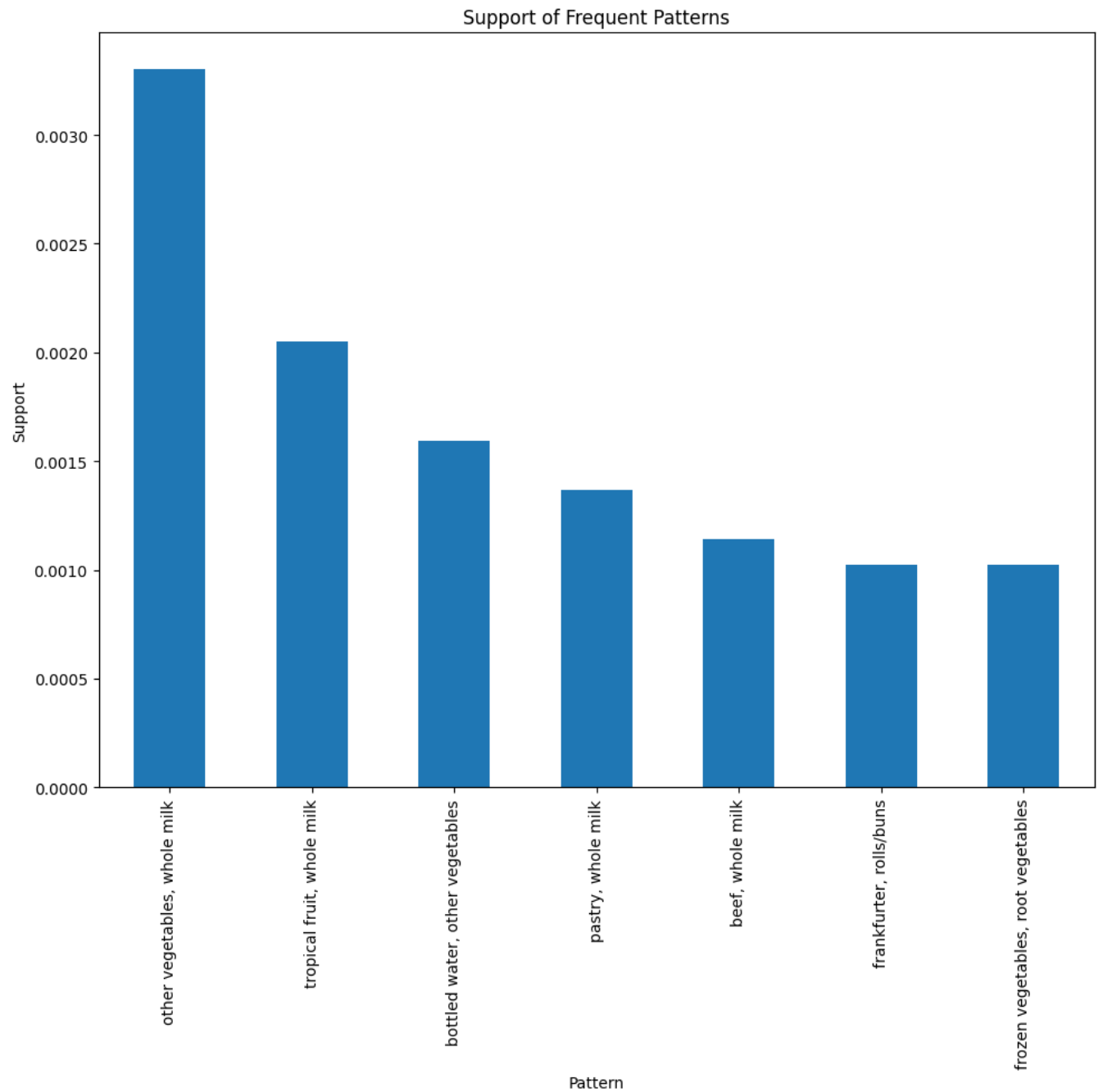
```python
1  df = pd.DataFrame({
2      'Pattern': [],
3      'Support': [],
4      'Confidence': []
5  })
6  for item in test_results:
7    new_row = pd.DataFrame({
8      'Pattern': [item[0]],
9      'Support': [item[1][0]],
10     'Confidence': [item[1][1]]})
11   df = pd.concat([df, new_row], ignore_index=True)
```
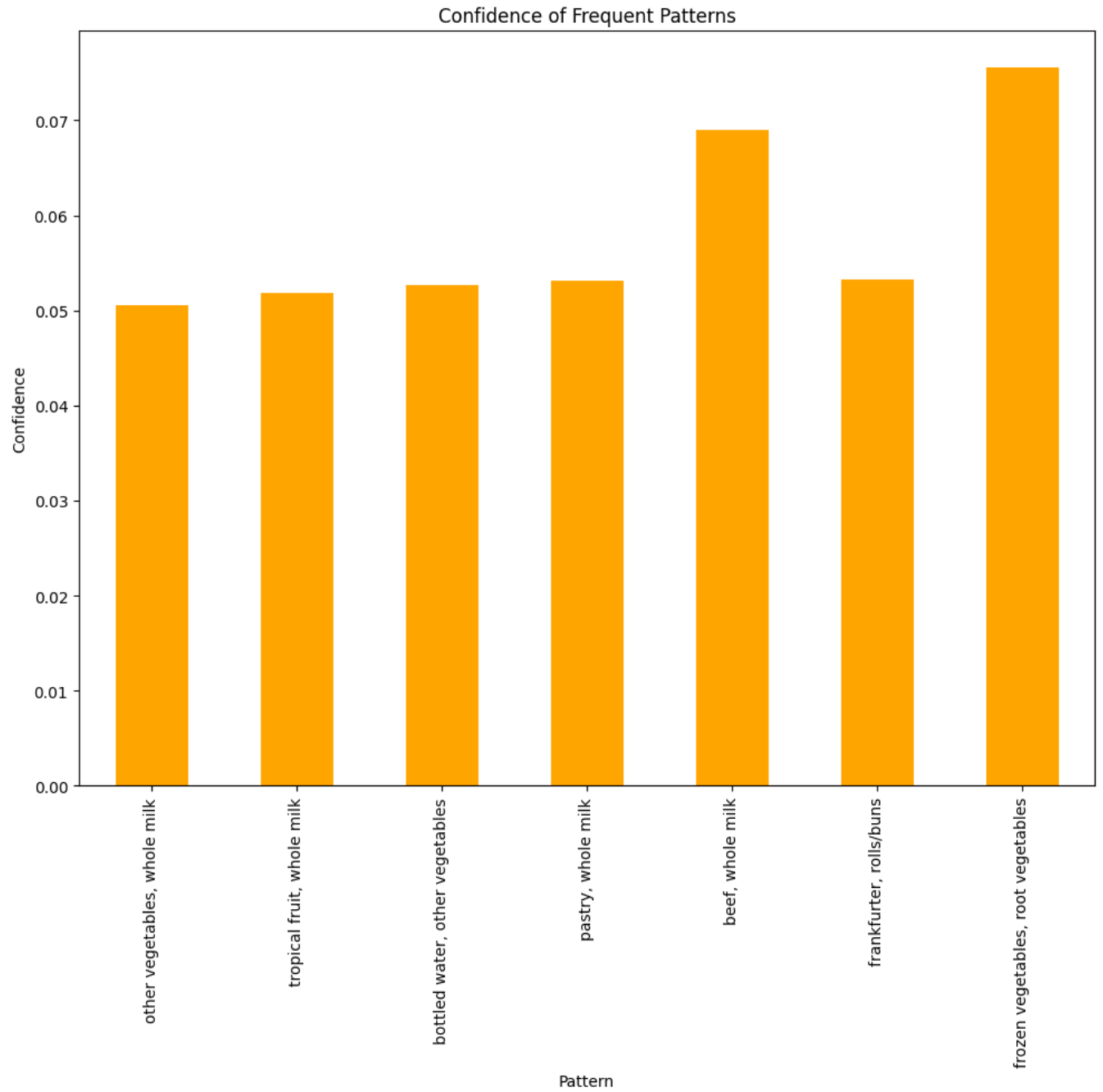
```python
1 import matplotlib.pyplot as plt
2 fig, ax = plt.subplots(1, 1, figsize=(10, 10))
3 df.plot.bar(x='Pattern', y='Support', ax=ax, legend=False)
4 ax.set_ylabel('Support')
5 ax.set_title('Support of Frequent Patterns')
6 plt.tight_layout()
7 plt.show()
```

```python
1 import matplotlib.pyplot as plt
2 fig, ax = plt.subplots(1, 1, figsize=(10, 10))
3 df.plot.bar(x='Pattern', y='Support', ax=ax, legend=False)
```

Support of Frequent Patterns

```
1  fig, ax = plt.subplots(1, 1, figsize=(10, 10))
2
```

```
3 df.plot.bar(x='Pattern', y='Confidence', ax=ax, color='orange', legend=Fals
4 ax.set_ylabel('Confidence')
5 ax.set_title('Confidence of Frequent Patterns')
6 plt.tight_layout()
7 plt.show()
```

Confidence of Frequent Patterns

```
1 # Week 11 Workshop Solutions were used in the development of the code in th
2 # https://myuni.adelaide.edu.au/courses/91968/files/14319268?wrap=1
```

```
 1 import numpy as np
 2
 3 # Reading in Training and Test Datasets
 4 train_set_df = pd.read_csv('Groceries data train.csv')
 5 test_set_df = pd.read_csv('Groceries data test.csv')
 6
 7 train_rows, test_cols = np.shape(train_set_df)
 8 test_rows, test_cols = np.shape(test_set_df)
 9
10 print("Train:", np.shape(train_set_df))
11 print("Test:", np.shape(test_set_df))
```

```
Train: (27000, 7)
Test: (11765, 7)
```

```
1 train_set = np.array(train_set_df)
2 test_set = np.array(test_set_df)
```

```
1 print(train_set[0:3])
```

```
[[3021 '30/01/2015' 'frankfurter' 2015 1 30 4]
 [1292 '24/10/2015' 'pork' 2015 10 24 5]
 [4206 '4/04/2014' 'root vegetables' 2014 4 4 4]]
```

```
1 print(test_set[0:3])
```

```
[[3481 '8/03/2015' 'candy' 2015 3 8 6]
 [1254 '19/04/2015' 'white wine' 2015 4 19 6]
 [2835 '28/01/2014' 'domestic eggs' 2014 1 28 1]]
```

```
1 train_set = np.array([[row[0], row[2]] for row in train_set])
2 test_set = np.array([[row[0], row[2]] for row in test_set])
```

```
1 train_set = train_set[train_set[:, 0].argsort()]        # this part is (bas
2 test_set = test_set[test_set[:, 0].argsort()]           # this part is (bas
```

```python
1  ids = []
2  items = []
3
4  for i in range(train_rows):
5      if train_set[i][0] not in ids:
6          ids.append(train_set[i][0])
7
8      if train_set[i][1] not in items:
9          items.append(train_set[i][1])
10
11 for i in range(test_rows):
12     if test_set[i][0] not in ids:
13         ids.append(test_set[i][0])
14
15     if test_set[i][1] not in items:
16         items.append(test_set[i][1])
17
18 user_count = len(ids)
19 item_count = len(items)
```

```python
1  new_id_dict = {}
2  new_id_reverse_dict = {}
3
4  for i in range(len(ids)):
5      new_id_dict[ids[i]] = i
6      new_id_reverse_dict[i] = ids[i]
```

```python
1  item_id_dict = {}
2  item_id_reverse_dict = {}
3
4  for i in range(item_count):
5      item_id_dict[items[i]] = i
6      item_id_reverse_dict[i] = items[i]
```

```python
1  for i in range(train_rows):
2      train_set[i][0] = new_id_dict[train_set[i][0]]
3      train_set[i][1] = item_id_dict[train_set[i][1]]
4
5  for i in range(test_rows):
6      test_set[i][0] = new_id_dict[test_set[i][0]]
7      test_set[i][1] = item_id_dict[test_set[i][1]]
```

```
1 import copy
2
3 user_item_matrix = [[0 for i in range(item_count)] for j in range(user_coun
4 user_item_matrix = np.array(user_item_matrix)
5
6 test_user_item_matrix = copy.deepcopy(user_item_matrix)                # t
7 collab_filter_user_item_matrix = copy.deepcopy(user_item_matrix)       # t
```

```
1 for i in range(train_rows):
2     user = int(train_set[i][0])
3     item = int(train_set[i][1])
4     user_item_matrix[user][item] += 1
```

```
1 alt_user_item_matrix = copy.deepcopy(user_item_matrix)       # this part is
```

```
1 for i in range(user_count):
2     row_mean = np.mean(user_item_matrix[i])          # this part is (based o
3     for j in range(item_count):
4         alt_user_item_matrix[i][j] -= row_mean
```

```
1 # Method used to calculate cosine similarity in below function based off in
2 def cosine_similarity(vec1, vec2):
3     len1 = len(vec1)
4     len2 = len(vec2)
5     cos_sim = (np.dot(vec1, vec2)) / (len1 * len2)       # this part is (bas
6     return cos_sim
```

```
1 # Following function based on [8] (Week 11 Workshop Solutions)
2 def get_similarity_vec(user, item, user_item_matrix):
3     similarity_vec = []
4     for i in range(user_count):
5         if i != user and user_item_matrix[i][item] != 0:
6             cosine_sim = cosine_similarity(alt_user_item_matrix[user], alt_
7             similarity_vec.append([cosine_sim, i])
8
9     similarity_vec = sorted(similarity_vec, reverse=True)
10    return similarity_vec
```

```python
1  # [8] (Week 11 Workshop Solutions) assisted the development of the followir
2  def weighted_avg(user, item, n, original_u_i_matrix, modified_u_i_matrix):
3      sim_vec = get_similarity_vec(user, item, modified_u_i_matrix)
4      neighbours = []
5      cos_sims = []
6
7      n = min(n, len(sim_vec))
8      for i in range(n):
9          neighbours.append(sim_vec[i][1])
10         cos_sims.append(sim_vec[i][0])
11
12     if n == 0 or max(cos_sims) == 0:
13         return 0
14
15     weighted_sum = 0
16     for i in range(n):
17         original_item_count = original_u_i_matrix[neighbours[i]][item]
18         cos_sim = cos_sims[i]
19         weighted_sum += (cos_sim * original_item_count)
20
21     weight_avg = weighted_sum / sum(cos_sims)
22     return weight_avg
```

```python
1  # Following function based on [8] (Week 11 Workshop Solutions)
2  def recommend(user, original_u_i_matrix, modified_u_i_matrix, top_recs=1, t
3      recommendation_ratings = []
4
5      for i in range(item_count):
6          if original_u_i_matrix[user][i] == 0:
7              recommendation_rating = weighted_avg(user, i, top_users, origir
8              if recommendation_rating == 0:
9                  continue
10             recommendation_ratings.append([recommendation_rating, i])
11
12     if len(recommendation_ratings) == 0:
13         return 0
14
15     recommendation_ratings = sorted(recommendation_ratings, reverse=True)
16     return recommendation_ratings[:top_recs]
```

```python
1  for i in range(test_rows):
2      user = int(test_set[i][0])
3      item = int(test_set[i][1])
4      test_user_item_matrix[user][item] += 1
```

```
1  import math
2
3  # Formula for root mean squared error (RMSE) used to develop function below
4  def RMSE(errors):
5      rmse = 0
6      for error in errors:
7          error_squared = error * error
8          rmse += error_squared
9
10     rmse = rmse / len(errors)
11     return math.sqrt(rmse)        # this part is (based on) from [5]
```

```
1  collab_filter_losses = []
2
3  # [8] (Week 11 Workshop Solutions) assisted the development of the code in
4  for i in range(user_count):
5      recommendations = recommend(i, user_item_matrix, alt_user_item_matrix)
6      if recommendations == 0:
7          continue
8
9      for j in range(len(recommendations)):
10         recommended_item = recommendations[j][1]
11         if user_item_matrix[i][recommended_item] == 0:
12             collab_filter_user_item_matrix[i][recommended_item] = recommend
13             loss = collab_filter_user_item_matrix[i][recommended_item] - te
14             collab_filter_losses.append(loss)
15
16 print("Root Mean Squared Error (RMSE) for Collaborative Filtering:", RMSE(c
```

Root Mean Squared Error (RMSE) for Collaborative Filtering: 2.3536843146554

```
1  freq_itemsets = interest
2
3  for i in range(len(freq_itemsets)):
4      freq_itemset = freq_itemsets[i][0]
5      for j in range(len(freq_itemset)):
6          freq_itemset[j] = item_id_dict[freq_itemset[j]]
```

```
1  def get_itemset_relevance(user, freq_itemset, user_item_matrix):
2      item_relevances = []
3      items = freq_itemset[0]
4      interest = freq_itemset[1]
5
6      for i in range(len(items)):
7          item = items[i]
8          item_relevances.append(user_item_matrix[user][item] * interest)
9
10     return sum(item_relevances)
```

```
1  pattern_recs_true = 0
2  pattern_recs_false = 0
3
4  for i in range(user_count):
5      relevant_itemsets = []
6      for j in range(item_count):
7          for k in range(len(freq_itemsets)):
8              if user_item_matrix[i][j] == 1 and j in freq_itemsets[k][0]:
9                  relevant_itemsets.append(freq_itemsets[k])
10
11     if len(relevant_itemsets) == 0:
12         continue
13
14     itemset_relevances = []
15     for j in range(len(relevant_itemsets)):
16         itemset_relevances.append([get_itemset_relevance(i, relevant_itemse
17
18     itemset_relevances = sorted(itemset_relevances, reverse=True)
19     top_freq_itemset = relevant_itemsets[itemset_relevances[0][1]][0]
20
21     print("\nRecommendations for User", new_id_reverse_dict[i])
22     for j in range(len(top_freq_itemset)):
23         if user_item_matrix[user][top_freq_itemset[j]] == 0 and test_user_i
24             print(item_id_reverse_dict[top_freq_itemset[j]])
25             pattern_recs_true += 1
26         elif user_item_matrix[user][top_freq_itemset[j]] == 0:
27             print(item_id_reverse_dict[top_freq_itemset[j]])
28             pattern_recs_false += 1
29
30 pattern_recs_accuracy = pattern_recs_true / (pattern_recs_true + pattern_re
31 print("\nAccuracy of Recommendations Generated From Frequent Itemsets:", pa
```

⇥▾

```
Recommendations for User 1000
grapes
```

```
Recommendations for User 1001
processed cheese
rolls/buns

Recommendations for User 1002
ham
whole milk

Recommendations for User 1003
processed cheese
rolls/buns

Recommendations for User 1004
processed cheese
rolls/buns

Recommendations for User 1005
margarine

Recommendations for User 1006
processed cheese
rolls/buns

Recommendations for User 1008
UHT—milk
tropical fruit

Recommendations for User 1009
napkins
pastry

Recommendations for User 1010
UHT—milk
tropical fruit

Recommendations for User 1011
processed cheese
rolls/buns

Recommendations for User 1012
processed cheese
rolls/buns

Recommendations for User 1013
ham
whole milk

Recommendations for User 1014
ham
whole milk

Recommendations for User 1015
ham
```

. _ ...

## Reference List

1. https://stackoverflow.com/questions/2828059/sorting-arrays-in-numpy-by-column
2. https://stackoverflow.com/questions/2612802/how-do-i-clone-a-list-so-that-it-doesnt-change-unexpectedly-after-assignment
3. https://numpy.org/doc/stable/reference/generated/numpy.mean.html
4. https://numpy.org/doc/stable/reference/generated/numpy.dot.html
5. https://www.w3schools.com/python/ref_math_sqrt.asp
6. https://myuni.adelaide.edu.au/courses/91968/pages/module-8-online-learning?module_item_id=3274637
7. https://en.wikipedia.org/wiki/Cosine_similarity#:~:text=In%20data%20analysis%2C%20cosine%20similarity,the%20product%20of%20their%20lengths.
8. https://myuni.adelaide.edu.au/courses/91968/files/14319268?wrap=1
9. https://docs.oracle.com/en/cloud/saas/planning-budgeting-cloud/pfusu/insights_metrics_RMSE.html

```python
1  import matplotlib.pyplot as plt
2  import pandas as pd
3  import seaborn as sns
4
5  data = pd.read_csv('Groceries data train.csv')
6  data.head()
```

|   | Member_number | Date | itemDescription | year | month | day | day_of_week |
|---|---|---|---|---|---|---|---|
| **0** | 3021 | 30/01/2015 | frankfurter | 2015 | 1 | 30 | 4 |
| **1** | 1292 | 24/10/2015 | pork | 2015 | 10 | 24 | 5 |
| **2** | 4206 | 4/04/2014 | root vegetables | 2014 | 4 | 4 | 4 |
| **3** | 4369 | 25/08/2015 | onions | 2015 | 8 | 25 | 1 |
| **4** | 1522 | 1/07/2014 | waffles | 2014 | 7 | 1 | 1 |

```
1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 27000 entries, 0 to 26999
Data columns (total 7 columns):
 #   Column          Non-Null Count  Dtype
---  ------          --------------  -----
 0   Member_number   27000 non-null  int64
 1   Date            27000 non-null  object
 2   itemDescription 27000 non-null  object
 3   year            27000 non-null  int64
 4   month           27000 non-null  int64
 5   day             27000 non-null  int64
 6   day_of_week     27000 non-null  int64
dtypes: int64(5), object(2)
memory usage: 1.4+ MB
```

```
1 data['Date'] = pd.to_datetime(data['Date'],dayfirst=True)
2 data['year'] = data['Date'].dt.year
3 data['month'] = data['Date'].dt.month
4 data['day'] = data['Date'].dt.day
5 data['day_of_week'] = data['Date'].dt.dayofweek
6 basket = data.groupby(['Member_number', 'Date'])['itemDescription'].apply(l
7 basket['itemDescription'] = basket['itemDescription'].apply(lambda x: ', '.
8 # basket_sets = basket['itemDescription'].str.get_dummies(sep=', ')
9 basket.head(10)
```

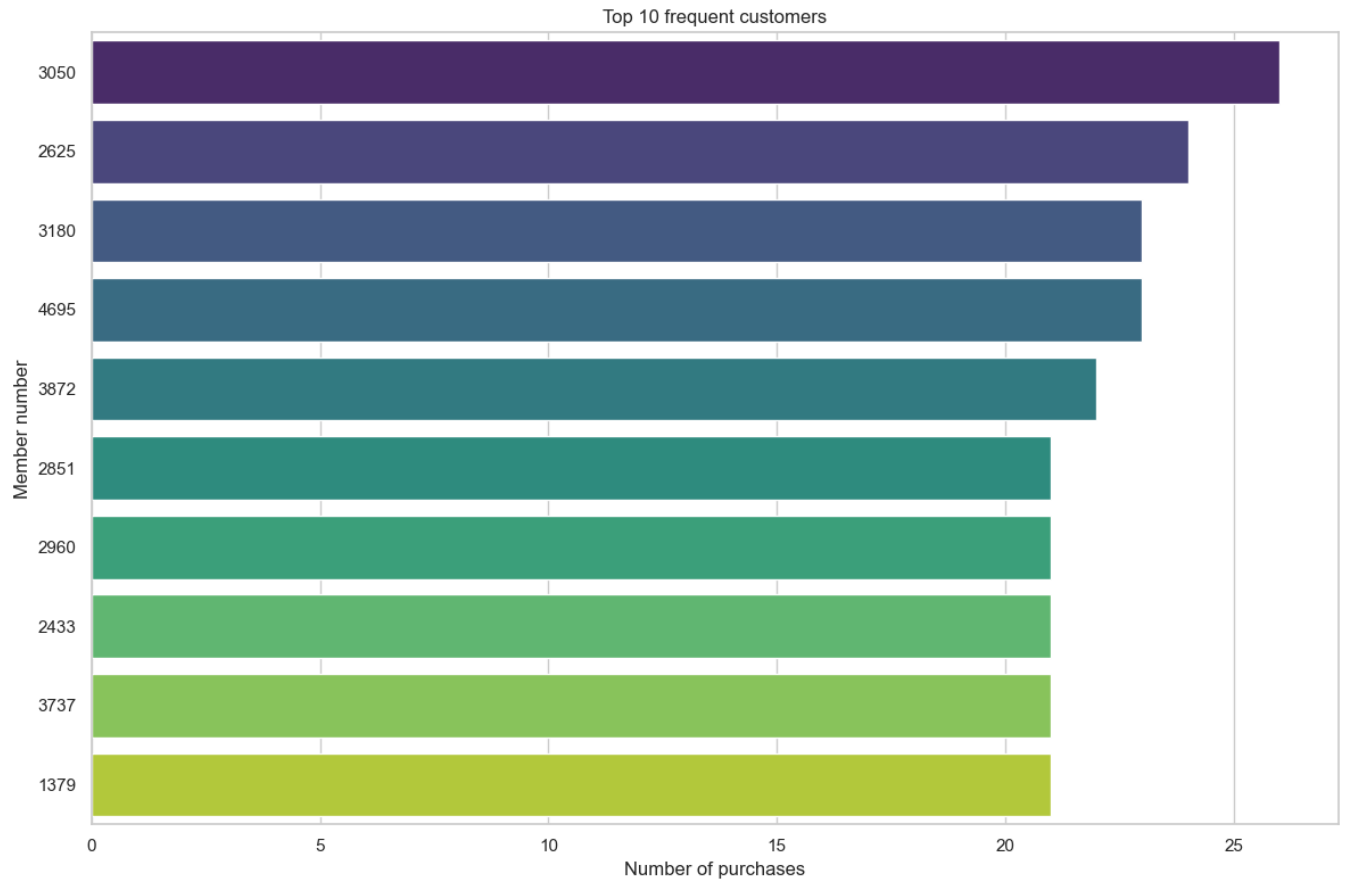|   | Member_number | Date | itemDescription |
|---|---|---|---|
| 0 | 1000 | 2014-06-24 | pastry |
| 1 | 1000 | 2015-03-15 | sausage, yogurt |
| 2 | 1000 | 2015-05-27 | soda, pickled vegetables |
| 3 | 1000 | 2015-07-24 | misc. beverages, canned beer |
| 4 | 1000 | 2015-11-25 | sausage |
| 5 | 1001 | 2014-07-02 | rolls/buns, whole milk, sausage |
| 6 | 1001 | 2014-12-12 | soda, whole milk |
| 7 | 1001 | 2015-01-20 | soda, whipped/sour cream |
| 8 | 1001 | 2015-02-05 | frankfurter, curd |
| 9 | 1001 | 2015-04-14 | white bread, beef |

```
1 sales_count = data['Member_number'].apply(lambda x: str(x)).value_counts().
```

```
 2 sales_count.columns = ['Member_number', 'count']
 3
 4 sns.set(style="whitegrid")
 5
 6
 7 plt.figure(figsize=(12,8))
 8 sns.barplot(y='Member_number', x='count', data=sales_count, palette='viridi
 9 plt.title('Top 10 frequent customers')
10 plt.xlabel('Number of purchases')
11 plt.ylabel('Member number')
12 plt.tight_layout()
13 plt.savefig("top_customer.png")
14 plt.show()
```

C:\Users\ebube\AppData\Local\Temp\ipykernel_8092\419175320.py:8: FutureWarn

Passing `palette` without assigning `hue` is deprecated and will be removed

sns.barplot(y='Member_number', x='count', data=sales_count, palette='viri
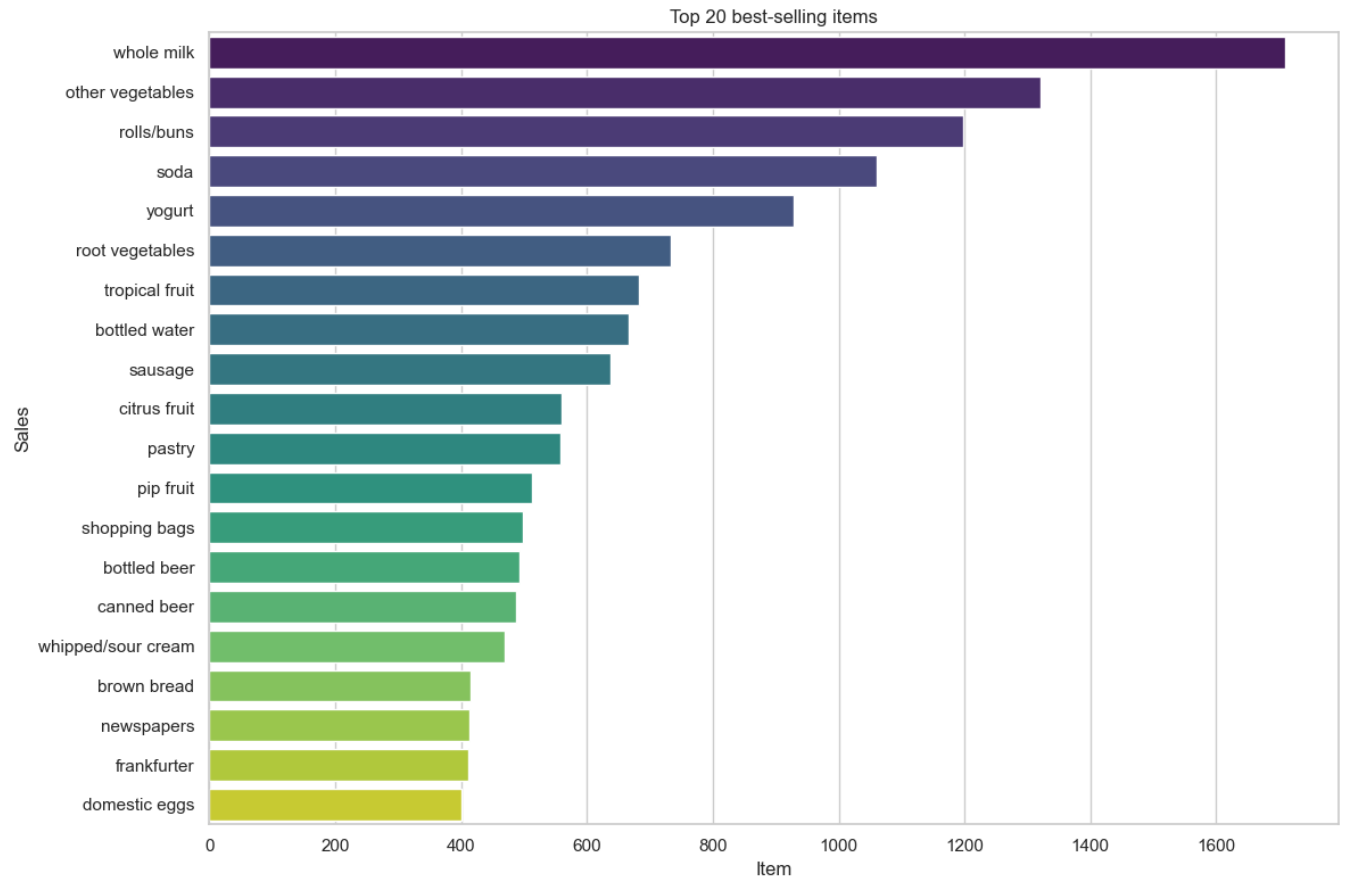


Top 10 frequent customers

```
1 sales_count = data['itemDescription'].value_counts().reset_index()[:20]
```

```
 1  sales_count = data['ItemDescription'].value_counts().reset_index()[:20]
 2  sales_count.columns = ['itemDescription', 'count']
 3
 4
 5  sns.set(style="whitegrid")
 6
 7  # 绘制直方图
 8  plt.figure(figsize=(12,8))
 9  sns.barplot(y='itemDescription', x='count', data=sales_count, palette='viri
10  plt.title('Top 20 best-selling items')
11  plt.xlabel('Item')
12  plt.ylabel('Sales')
13  plt.xticks()
14  plt.tight_layout()
15  plt.savefig('topsales.png', dpi=300)
16
17  plt.show()
```

```
C:\Users\ebube\AppData\Local\Temp\ipykernel_8092\2162360578.py:9: FutureWar
Passing `palette` without assigning `hue` is deprecated and will be removed
sns.barplot(y='itemDescription', x='count', data=sales_count, palette='vi
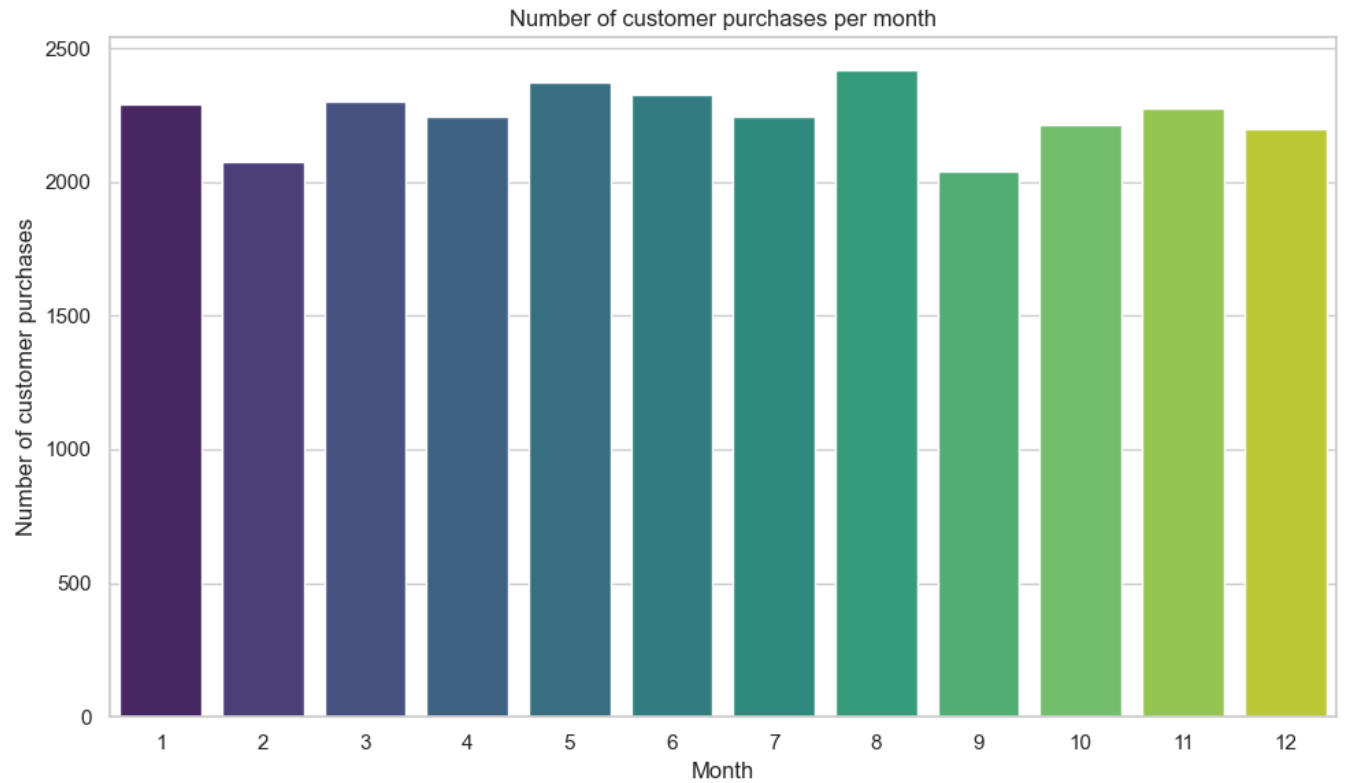```



Top 20 best-selling items

```
1 data['month'] = data['month'].astype(int)
```

```
 1 data['month'] = data['month'].astype(int)
 2
 3
 4 monthly_counts = data['month'].value_counts().sort_index()
 5
 6 monthly_counts_df = monthly_counts.reset_index()
 7 monthly_counts_df.columns = ['month', 'count']
 8
 9 sns.set(style="whitegrid")
10
11 plt.figure(figsize=(10,6))
12 sns.barplot(x='month', y='count', data=monthly_counts_df, palette='viridis'
13 plt.title('Number of customer purchases per month')
14 plt.xlabel('Month')
15 plt.ylabel('Number of customer purchases')
16 plt.xticks(rotation=0)
17 plt.tight_layout()
18 plt.savefig("monthly_counts.jpg",dpi=300)
19 plt.show()
```

⯈ C:\Users\ebube\AppData\Local\Temp\ipykernel_8092\2750937065.py:12: FutureWa

   Passing `palette` without assigning `hue` is deprecated and will be removed

   sns.barplot(x='month', y='count', data=monthly_counts_df, palette='viridi



Number of customer purchases per month

```
1 basket_sets = basket['itemDescription'].str.get_dummies(sep=', ')
```

```
1 basket_sets
```

| | Instant food products | UHT-milk | abrasive cleaner | artif. sweetener | baby cosmetics | bags | baking powder | bathroom cleaner |
|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **1** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **2** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **3** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **4** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| **13896** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **13897** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **13898** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **13899** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **13900** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

```
1 data = pd.read_csv('Groceries data test.csv')
2 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11765 entries, 0 to 11764
Data columns (total 7 columns):
 #   Column          Non-Null Count  Dtype
---  ------          --------------  -----
 0   Member_number   11765 non-null  int64
 1   Date            11765 non-null  object
 2   itemDescription 11765 non-null  object
 3   year            11765 non-null  int64
 4   month           11765 non-null  int64
 5   day             11765 non-null  int64
 6   day_of_week     11765 non-null  int64
dtypes: int64(5), object(2)
memory usage: 643.5+ KB
```