

# Technical Report for KDD CUP TASK3

## Main idea

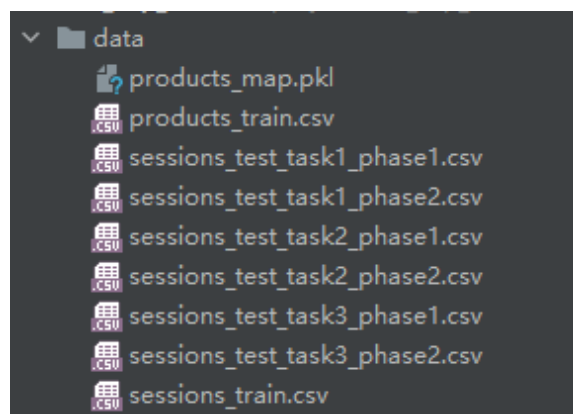
We recommend the next item based on a simple co-visiting graph (the number of times two items co-occur in the same session). But the difference in task3 is that if the item prediction is wrong, it does not mean that the BLEU score is low, because the title may still be very similar compared to ground true. In order to avoid the risk of dissimilarity between the title and the ground true due to incorrect product predictions, we do not directly take the title of the predicted item, but take the intersection of its token and the title of the last item (the titles of two adjacent item are often similar, and we have tried that if the title of the last item is used directly as the predicted title, the BLEU score is also considerable).

## Code run

1. Download data

from <https://drive.google.com/drive/folders/1oKf7yaeo3AxtxWayaRoHla5ESmrlWLed?usp=sharing>. Note that we are not using any external data.

All data looks like the following,



Run our code:

```
python prediction.py
```

## Details

1. Read all items, training set, test set, and the mapping dictionary from item id to title.

```
29 products = read_product_data() # load products
30 hist_data = read_train_data() # load training data
31 test_sessions = read_test_data(task) # load test data
32 test_locale_names = test_sessions['locale'].unique()
33 with open("data/products_map.pkl", "rb") as tf: # load id:title dictionary
34     products_map = pickle.load(tf)
```

2. Construct co-visiting diagram. Here we are based on two considerations. First, only the most recent five items are used for construction per session. Because earlier interactions may have lost their timeliness. Second, among the five items, not any two co-visiting item have the same weight. If the distance between two items (judged by the order of interactions) is greater, the weight will be smaller, so that we can focus more on adjacent interacting items.

In addition, in addition to the construction based on the training set, we can also use the test set of the phase 1 of task 3, and the test sets of task1 and task2 for constructing. For the training set, we use the 4 items of history and the next item. For the other supplementary data, we select five historical commodities due to missing labels.

```

36 # construct co-visiting graph
37 all_id = products['id'].unique()
38 id_to_idx = {id: i for i, id in enumerate(all_id)}
39 idx_to_id = {i: id for i, id in enumerate(all_id)}
40 graph = {id: {} for i, id in enumerate(all_id)}
41 phase1 = pd.read_csv('data/sessions_test_task3_phase1.csv')
42 phase2 = pd.read_csv('data/append_data.csv')
43 for _, row in hist_data.iterrows(): # use training data to construct co-visiting graph
44     items = ([s.strip("\n") for s in row['prev_items'][1:-1].split(" ")] + [row['next_item'].strip()])[:5]
45
46     for i in range(0, len(items)):
47         for j in range(i + 1, len(items)):
48             if (items[j] != items[i]):
49                 try:
50                     graph[items[i]][items[j]] += 1.1 - abs(j - i) * 0.2
51                 except:
52                     graph[items[i]][items[j]] = 1.1 - abs(j - i) * 0.2
53                 try:
54                     graph[items[j]][items[i]] += 1.1 - abs(j - i) * 0.2
55                 except:
56                     graph[items[j]][items[i]] = 1.1 - abs(j - i) * 0.2
57 for _, row in phase1.iterrows(): # use task3 phase1 data to construct co-visiting graph
58     items = ([s.strip("\n") for s in row['prev_items'][1:-1].split(" ")] + [row['next_item'].strip()])[:5]
59     for i in range(0, len(items)):
60         for j in range(i + 1, len(items)):
61             if (items[j] != items[i]):
62                 try:
63                     graph[items[i]][items[j]] += 1.1 - abs(j - i) * 0.2
64                 except:
65                     graph[items[i]][items[j]] = 1.1 - abs(j - i) * 0.2
66                 try:
67                     graph[items[j]][items[i]] += 1.1 - abs(j - i) * 0.2
68                 except:
69                     graph[items[j]][items[i]] = 1.1 - abs(j - i) * 0.2
70 with open("data/graph_session_last_copy.pkl", "wb") as tf:
71     pickle.dump(graph, tf)
72 for _, row in phase2.iterrows(): # use task1& task2 data to construct co-visiting graph
73     items = ([s.strip("\n") for s in row['prev_items'][1:-1].split(" ")] + [row['next_item'].strip()])[:5]
74     for i in range(0, len(items)):
75         for j in range(i + 1, len(items)):
76             if (items[j] != items[i]):
77                 try:
78                     graph[items[i]][items[j]] += 1.1 - abs(j - i) * 0.2
79                 except:
80                     graph[items[i]][items[j]] = 1.1 - abs(j - i) * 0.2
81                 try:
82                     graph[items[j]][items[i]] += 1.1 - abs(j - i) * 0.2
83                 except:
84                     graph[items[j]][items[i]] = 1.1 - abs(j - i) * 0.2
85 with open("data/graph_session_last.pkl", "wb") as tf:
86     pickle.dump(graph, tf)
87 with open("data/graph_session_last_copy.pkl", "rb") as tf:
88     graph = pickle.load(tf)
89 with open("data/graph_session_last.pkl", "rb") as tf:
90     graph1 = pickle.load(tf)

```

3. For each session, we choose the last item in history (Line 127-135 in utils.py). Based on the co-visiting graph, we can get the candidate items that appear together with the item. If the frequency of a certain product is particularly prominent (here we are based on the 6-sigma principle, Line 145), we select it as the final candidate item, otherwise, there is no candidate.

```

127 # obtain the latest record of the corresponding language in the history record, and we predict the next title based on this record
128 for i in range(len(his_list)):
129     try:
130         title = products[locale][his_list[i]]
131         if (isinstance(title, str)):
132             his_titles.append(title)
133             his_titles_id.append(his_list[i])
134     except:
135         pass
136
137 # we only use the last interacted item
138 his_titles = his_titles[:1]
139
140 # get the last item's co-visiting items
141 titles_map = graph(his_titles_id[0])
142 titles_map = sorted(titles_map.items(), key=lambda x: x[1], reverse=True)
143 titles_map = {v[0]: v[1] for v in titles_map}
144 titles_array = np.array(list(titles_map.values()))
145 threshold = np.mean(titles_array) + 6. * np.std(titles_array)
146
147 candidate = []
148 for key, value in titles_map.items():
149     if (
150         value > threshold): # if the maximum item's weight in co-visiting items is larger than a certain threshold (here is 6-sigma as shown in line 145)
151         try:
152             candidate.append(" ".join(products[locale][key].split()).replace(" ", " "))
153         except:
154             pass
155     else:
156         break

```

4. Since the recommendation strategy is simple, if there is no candidate item that is particularly prominent, to avoid recommendation errors, we directly select the last interacted item's title.

```

181 if (len(candidate) == 0):
182     answer = his_titles[0]

```

5. If there is a final candidate, we still will not choose it directly, because we cannot bear the risk of it being wrong and causing a sharp drop in BLEU. Therefore, we take the intersection of the candidate title and the title of the last interacted item. If the number of tokens they overlap is greater than xx% tokens of the candidate, we will choose this intersection as the final predicted title. If there are too few intersections, in order to avoid risks, we still choose the title of the last interacted product as the predicted title.

```

185     else:
186         list1 = split_notation(his_titles[0].lower().split())
187         set1 = set(list1)
188         list2 = split_notation(candidate[0].lower().split())
189         list22 = split_notation(candidate[0].split())
190         set2 = set(list2)
191         set3 = set1 & set2
192         list3 = []
193         list4 = []
194         for iii in range(len(list2)):
195             if (list2[iii] in set3):
196                 list3.append(list22[iii])
197                 if (iii > 0):
198                     if ("%s %s" % (list22[iii - 1], list22[iii]) in candidate[0]):
199                         if (len(list4) > 0 and list4[-1] != " "):
200                             list4.append(" ")
201                             list4.append(list22[iii])
202                         elif ("%s%s" % (list22[iii - 1], list22[iii]) in candidate[0]):
203                             list4.append(list22[iii])
204                     else:
205                         list4.append(list22[iii])
206             else:
207                 if (iii > 0):
208                     if ("%s %s" % (list22[iii - 1], list22[iii]) in candidate[0] and len(list4) > 0 and
209                         list4[-1] != " "):
210                         list4.append(" ")
211                 if (locale == 'JP' and len("".join(list3)) > len("".join(list2)) * 0.8):
212                     list4 = finetune(list4)
213                     answer = "".join(list4)
214                     candidate_count += 1
215                 elif (locale != 'JP' and len(list3) > len(list2) * 0.8):
216                     list4 = finetune(list4)
217                     answer = "".join(list4)
218                     candidate_count += 1
219             else:
220                 answer = his_titles[0]
221

```

6. The final evaluation score of the method on the test set is 0.27130.

233280	ustc-gobble	graded	0.27130	0.99993	Graded successfully!	Sun, 4 Jun 2023 14:54:51	<a href="#">View</a>
--------	-------------	--------	---------	---------	----------------------	-----------------------------	----------------------