

Mining temporal data – Frequent sequences

Algorithmic Data Analysis – Coding Assignment 2

Giulia Ortolani

March 25, 2024

Resources: No resources

Collaborations: No collaborations

The Generalized Sequential Pattern Mining (GSP) algorithm is an apriori-like algorithm for solving sequential pattern mining, here I try to implement it in Python and use it for mining frequent patterns from discrete sequences. The algorithm is then applied to the CRSW dataset, which contains two months worth of weather data in Kuopio (January–February 2019).

GSP algorithm implementation idea

The implemented algorithm takes a dataset (list of itemsets) and a minimum support threshold. It first identifies frequent items (items with support greater than or equal to the minimum support threshold). Then, starting from $k = 1$, it iteratively generates candidate sequences of length $k + 1$ by joining pairs of sequences from the frequent patterns of length k . The support of these candidates is calculated and only those with support greater than or equal to the minimum support threshold are retained. This process continues until no more frequent patterns can be found. The output of the code is a dictionary where the keys are the frequent patterns and the values are their corresponding support counts.

Application to the CRSW dataset

The weather time series analysis using the Generalized Sequential Pattern (GSP) mining algorithm has revealed interesting patterns in the occurrences of clouds (C), precipitation (R), sunshine (S), and wind (W).

Setting different minimum support threshold

If we set the minimum support threshold to **30**, we get three frequent patterns:

- 'CRW', with a support of 225. This is cloudy, rainy and windy weather.
- 'CR', with a support of 99. This is cloudy and rainy weather.
- 'CW', with a support of 237. This is cloudy and windy weather.

Setting different values for max gap and max span

If we change the values for max gap and max span we can obtain these different results:

Min supp thresh	Max span	Max gap	Frequent patterns
10	5	10	[('CRW', 225), ('CR', 99), ('CW', 237)]
1	1	1	[('CRW', 75), ('CR', 99), ('CW', 237)]
200	20	20	[('CW', 237)]
100	10	10	[('CRW', 225), ('CW', 237)]

Table 1: Frequent patterns using different parameters

A. Code

```

1 import copy
2 ## GSP
3 # implementation of the generalized sequential pattern mining (GSP) algorithm for
  mining frequent patterns from discrete sequences
4
5 # def get_sequences_from_file(file_path):
6 #     sequences = []
7 #     with open(file_path, 'r') as file:
8 #         for line in file:
9 #             timestamp, sequence = line.strip().split()[0], line.strip().split()[1:]
10 #             sequences.append(sequence)
11 #     return sequences
12
13 def get_sequences_from_file(file_path):
14     sequences = []
15     with open(file_path, 'r') as file:
16         for line in file:
17             timestamp, sequence = line.strip().split()[0], line.strip().split()[1:]
18             split_sequence = [list(itemset) for itemset in sequence]
19             sequences.append(split_sequence)
20     return sequences
21
22
23
24 def join_sequences(frequent_items, maxspan, maxgap):
25     # candidate (k + 1)-sequences are generated by combining pairs of frequent k-
    sequences
26     candidates = []
27
28     # Fix a canonical order on the sequences and ensure that each candidate (k + 1)-
    subsequence is generated by only one pair of k-sequences ??
29
30
31     for i in range(len(frequent_items)):
32         Sa = list(frequent_items[i])
33
34         for j in range(len(frequent_items)):
35
36             Sb = list(frequent_items[j])
37
38             Sc = []
39
40             if check_candidate_gen(Sa, Sb) != 0:
41                 # If the first itemset of Sa is a singleton, Sc is generated by
    prepending it to Sb
42                 if len(list(Sa)[0]) == 1: # here I'm checking if the list
43                     Sa = "".join(Sa)
44                     Sb = "".join(Sb)
45                     Sc.insert(0, Sa)
46                     Sc.append(Sb)
47                     Sc = "".join(Sc)
48                     candidates.append(Sc)
49

```

```

50         # If the first itemset of Sa is not a singleton, Sc is generated by
replacing the first itemset of Sb with it
51         elif len(list(Sa)[0])>1:
52             Sa = "".join(Sa)
53             Sb = "".join(Sb)
54             Sc.insert(0,Sa)
55             Sc.append(Sb[1:])
56             Sc = "".join(Sc)
57             candidates.append(Sc)
58
59         # If the last itemset of Sb is a singleton, Sc is generated by
appending it to Sa
60         elif len(list(Sb)[-1])==1:
61             Sa = "".join(Sa)
62             Sb = "".join(Sb)
63             Sc.append(Sa)
64             Sc.append(Sb[-1])
65             Sc = "".join(Sc)
66             candidates.append(Sc)
67
68         # If the last itemset of Sb is not a singleton, Sc is generated by
replacing the last itemset of Sa with it
69         elif len(list(Sb)[-1])>1:
70             Sa = "".join(Sa)
71             Sb = "".join(Sb)
72             Sc.append(Sa[0:(len(Sa)-2)])
73             Sc.append(Sb[-1])
74             Sc = "".join(Sc)
75             candidates.append(Sc)
76
77     return candidates
78
79 def check_candidate_gen(Sa, Sb):
80     # removing an item from the first itemset of Sa and removing an item from the
81     # last itemset of Sb should result in the same sequence So
82
83     Sa1 = [list(item) for item in Sa]
84     Sb2 = [list(item) for item in Sb]
85
86     Sa_red = copy.deepcopy(Sa1)
87     Sb_red = copy.deepcopy(Sb2)
88
89     # first itemset of Sa: Sa_red[0]
90     # last itemset of Sb: Sb_red[-1]
91
92     k = 0
93
94     for i in range(len(Sa_red[0])):
95         # remove an item from the first itemset of Sa
96         if i == len(Sa_red[0]):
97             break
98         Sa_red[0].pop(i)
99         Sa_red = [item for item in Sa_red if item != []]
100
101     for j in range(len(Sb_red[-1])):
102         # remove an item from the last itemset of Sb
103         if j == len(Sb_red[-1]):
104             break
105         Sb_red[-1].pop(j)
106         Sb_red = [item for item in Sb_red if item != []]
107
108     k += (Sa_red == Sb_red)
109
110     Sb_red = copy.deepcopy(Sa1)
111     Sa_red = copy.deepcopy(Sb2)
112

```

```

113     #print(k)
114     return k
115
116 def is_subsequence(candidate, sequence, max_span, max_gap):
117
118     # CONSTRAINTS
119     # 1. maximum span limit the time difference between the first and last itemsets
120     # of a subsequence
121     # 2. maximum gap limit the number of gaps between successive itemsets of a
122     # subsequence
123
124     # Initialize indices for candidate and sequence
125     i, j = 0, 0
126     span_start = None
127
128     # Iterate through both lists
129     while i < len(candidate) and j < len(sequence[0]):
130         # If the elements match, move to the next element in both lists
131         if candidate[i] == sequence[0][j] and abs(i-j) < max_gap:
132             if span_start is None:
133                 span_start = j # Start of the span
134             i += 1
135             # Move to the next element in the larger list
136             j += 1
137
138     # If all elements of the candidate list are found in the sequence in the same
139     # order, check constraints
140     if i == len(candidate):
141         span_end = j - 1 # End of the span
142         span_length = span_end - span_start
143
144         # Check constraints
145         if span_length <= max_span:
146             return True
147
148     return False
149
150 def count_support(dataset, candidates, maxspan, maxgap):
151
152     support = {}
153     # the support of subsequence S in sequence D is the number of occurrences of S in
154     # D
155
156     for candidate in candidates:
157         candidate = list(candidate)
158         for itemset in dataset:
159             if is_subsequence(candidate, itemset, maxspan, maxgap):
160                 itemset = "".join(itemset[0])
161                 if itemset in support:
162                     support[itemset] += 1
163                 else:
164                     support[itemset] = 1
165
166     return support
167
168 def gsp(dataset, minsup, maxspan, maxgap):
169     k = 1
170     item_counts = {}
171
172     # Fk: all frequent items
173
174     # count support of items.... or itemsets?
175     for sequence in dataset:
176         for itemset in sequence:
177             for i in itemset:
178                 if i in item_counts:

```

```

175         item_counts[i] += 1
176     else:
177         item_counts[i] = 1
178
179     # itemsets:
180     # for sequence in dataset:
181     #     for itemset in sequence:
182     #         key = ''.join(itemset)
183     #         if key in item_counts:
184     #             item_counts[key] += 1
185     #         else:
186     #             item_counts[key] = 1
187
188     # Fk:
189     frequent_items = {item: count for item, count in item_counts.items() if count >=
190     minsup}
191
192     result = {}
193
194     # while Fk != 0
195     while len(frequent_items)!=0:
196
197         # generate C_k+1 by joining pairs of sequences from F_k
198         candidates = join_sequences(list(frequent_items.keys()), maxspan, maxgap)
199         candidate_support = count_support(dataset, candidates, maxspan, maxgap)
200
201         # S in C_k+1
202         frequent_items = {pattern: support for pattern, support in candidate_support.
203         items() if support >= minsup} # add to result if the support is > minsup
204         result.update(frequent_items) # F_k+1
205
206         k += 1
207
208     return result # dict object made of 'pattern' and 'support'
209
210 ## Apply the algorithm on the CRSW dataset:
211 # '2019-01-03_itemsets-CRSW.txt' Two months worth of weather in Kuopio (
212 # JanuaryFebruary 2019), obtained from https://en.ilmatieteenlaitos.fi/download-
213 # observations#!/
214 # Each line represents weather events during one hour as an itemset. The first column
215 # contains the contextual attribute, i.e. indication of time in Year-Month-
216 # Day_Hour format.
217 # The second column contains the itemset. C, R S and W stand respectively for clouds,
218 # precipitation, sunshine and wind.
219
220 file_path = "2019-01-03_itemsets-CRSW.txt"
221 minsup = 100
222 maxspan = 10
223 maxgap = 10
224 dataset = get_sequences_from_file(file_path)
225 result = gsp(dataset, minsup, maxspan, maxgap)
226
227 # Print the frequent patterns
228 print(f"Minimum support threshold = {minsup}")
229 print(f"Maximum span = {maxspan}")
230 print(f"Maximum gap = {maxgap}")
231 print("Frequent patterns:")
232 print(list(result.items()))
233
234 # Try different values for the minimum support threshold, different constraints
235 # such as max gap and max span,
236 # and considering the data either as one long sequence or with each day as a separate
237 # short sequence. You might
238 # also try replacing repeated occurrences of the same itemset by a single copy, i.e.
239 # represent constant weather
240 # during successive hours with only one itemset.

```

231 # Report on the type and number of patterns obtained under different conditions.