

ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN
KHOA TOÁN - TIN HỌC



BTTH TUẦN 5: KHAI THÁC DỮ LIỆU

Sinh viên thực hiện: Nguyễn Công Hoài Nam
Mã số sinh viên: 21280099

Ngày 11 tháng 5 năm 2024

MỤC LỤC

1	Sử dụng thư viện	1
2	Cài đặt lại thuật toán Vertical Apriori	3
2.1	Xử lý dữ liệu	3
2.2	Cài đặt thuật toán	4
2.3	So sánh kết quả	6

I. Sử dụng thư viện

- Cài đặt thư viện pyECLAT (đã cài đặt)
- Sử dụng lại file data.csv

```

1 import numpy as np
2 import pandas as pd
3 from pyECLAT import ECLAT
4
5
6 dataframe = pd.read_csv('data.csv', header = None)
7 dataframe

```

	0	1	2	3	4	5
0	Wine	Chips	Bread	Butter	Milk	Apple
1	Wine	NaN	Bread	Butter	Milk	NaN
2	NaN	NaN	Bread	Butter	Milk	NaN
3	NaN	Chips	NaN	NaN	NaN	Apple
4	Wine	Chips	Bread	Butter	Milk	Apple
5	Wine	Chips	NaN	NaN	Milk	NaN
6	Wine	Chips	Bread	Butter	NaN	Apple
7	Wine	Chips	NaN	NaN	Milk	NaN
8	Wine	NaN	Bread	NaN	NaN	Apple
9	Wine	NaN	Bread	Butter	Milk	NaN
10	NaN	Chips	Bread	Butter	NaN	Apple
11	Wine	NaN	NaN	Butter	Milk	Apple
12	Wine	Chips	Bread	Butter	Milk	NaN
13	Wine	NaN	Bread	NaN	Milk	Apple
14	Wine	NaN	Bread	Butter	Milk	Apple
15	Wine	Chips	Bread	Butter	Milk	Apple
16	NaN	Chips	Bread	Butter	Milk	Apple
17	NaN	Chips	NaN	Butter	Milk	Apple
18	Wine	Chips	Bread	Butter	Milk	Apple
19	Wine	NaN	Bread	Butter	Milk	Apple
20	Wine	Chips	Bread	NaN	Milk	Apple
21	NaN	Chips	NaN	NaN	NaN	NaN

Chuyển dữ liệu thành lớp ECLAT và khởi tạo DataFrame nhị phân

```

1 eclat_instance = ECLAT(data = dataframe, verbose = True)
2 eclat_instance.df_bin

```

	Chips	Butter	Wine	Milk	Apple	Bread
0	1	1	1	1	1	1
1	0	1	1	1	0	1
2	0	1	0	1	0	1
3	1	0	0	0	1	0
4	1	1	1	1	1	1
5	1	0	1	1	0	0
6	1	1	1	0	1	1
7	1	0	1	1	0	0
8	0	0	1	0	1	1
9	0	1	1	1	0	1
10	1	1	0	0	1	1
11	0	1	1	1	1	0
12	1	1	1	1	0	1
13	0	0	1	1	1	1
14	0	1	1	1	1	1
15	1	1	1	1	1	1
16	1	1	0	1	1	1
17	1	1	0	1	1	0
18	1	1	1	1	1	1
19	0	1	1	1	1	1
20	1	0	1	1	1	1
21	1	0	0	0	0	0

Khởi tạo các luật kết hợp

```

1  # count items in each row
2  items_per_transaction = eclat_instance.df_bin.astype(int).sum(axis=1)
3  # the item should appear at least at 5% of transactions
4  min_support = 0.5
5  # start from transactions containing at least 2 items
6  min_combination = 2
7  # up to maximum items per transactions
8  max_combination = max(items_per_transaction)
9  rule_indices, rule_supports = eclat_instance.fit(min_support = min_support,
10                                                  min_combination = min_combination,
11                                                  max_combination = max_combination,
12                                                  separator = ' & ',
13                                                  verbose = True)
14  result = pd.DataFrame(rule_supports.items(), columns=['Item', 'Support'])
15  # result = result.sort_values(by=['Support', 'Item'], ascending = False)
16  result

```

	Item	Support
0	Butter & Wine	0.500000
1	Butter & Milk	0.590909
2	Butter & Apple	0.500000
3	Butter & Bread	0.590909
4	Wine & Milk	0.636364
5	Wine & Apple	0.500000
6	Wine & Bread	0.590909
7	Milk & Apple	0.500000
8	Milk & Bread	0.590909
9	Apple & Bread	0.545455
10	Butter & Milk & Bread	0.500000
11	Wine & Milk & Bread	0.500000

II. Cài đặt lại thuật toán Vertical Apriori

1. Xử lý dữ liệu

Chuyển DataFrame về dạng List

```

1 records = []
2 for i in range(0,dataframe.shape[0]):
3     records.append([str(dataframe.values[i,j]) for j in range (0, dataframe.shape[1])])

```

Sau đó loại bỏ chuỗi 'nan' (không phải là một item trong itemsets)

```

1 records = [[item for item in i if item != 'nan'] for i in records]
2 records

```

```

[['Wine', 'Chips', 'Bread', 'Butter', 'Milk', 'Apple'],
 ['Wine', 'Bread', 'Butter', 'Milk'],
 ['Bread', 'Butter', 'Milk'],
 ['Chips', 'Apple'],
 ['Wine', 'Chips', 'Bread', 'Butter', 'Milk', 'Apple'],
 ['Wine', 'Chips', 'Milk'],
 ['Wine', 'Chips', 'Bread', 'Butter', 'Apple'],
 ['Wine', 'Chips', 'Milk'],
 ['Wine', 'Bread', 'Apple'],
 ['Wine', 'Bread', 'Butter', 'Milk'],
 ['Chips', 'Bread', 'Butter', 'Apple'],
 ['Wine', 'Butter', 'Milk', 'Apple'],
 ['Wine', 'Chips', 'Bread', 'Butter', 'Milk'],
 ['Wine', 'Bread', 'Milk', 'Apple'],
 ['Wine', 'Bread', 'Butter', 'Milk', 'Apple'],
 ['Wine', 'Chips', 'Bread', 'Butter', 'Milk', 'Apple'],
 ['Chips', 'Bread', 'Butter', 'Milk', 'Apple'],
 ['Chips', 'Butter', 'Milk', 'Apple'],
 ['Wine', 'Chips', 'Bread', 'Butter', 'Milk', 'Apple'],
 ['Wine', 'Bread', 'Butter', 'Milk', 'Apple'],
 ['Wine', 'Chips', 'Bread', 'Milk', 'Apple'],
 ['Chips']]

```

Hàm **convert_vertical** chuyển dataset từ dạng ngang (horizontal) sang dạng thẳng đứng (vertical)

```

1 def convert_vertical(records):
2     transactions = {}
3     for transaction_index, transaction in enumerate(records):
4         for t in transaction:

```

```

5         # frozenset là tập không thay đổi (immutable set) tức nó không cho phép chỉnh sửa
        , thay đổi hay giá trị lặp
6         t = frozenset([t])
7         if t not in transactions:
8             transactions[t] = []
9             transactions[t].append(transaction_index)
10    return transactions

```

Kiểm tra kết quả

```

1    X = convert_vertical(records)
2    for i in X:
3        print(i, "\t:", X[i])

```

```

frozenset({'Wine'})      : [0, 1, 4, 5, 6, 7, 8, 9, 11, 12, 13, 14, 15, 18, 19, 20]
frozenset({'Chips'})    : [0, 3, 4, 5, 6, 7, 10, 12, 15, 16, 17, 18, 20, 21]
frozenset({'Bread'})    : [0, 1, 2, 4, 6, 8, 9, 10, 12, 13, 14, 15, 16, 18, 19, 20]
frozenset({'Butter'})   : [0, 1, 2, 4, 6, 9, 10, 11, 12, 14, 15, 16, 17, 18, 19]
frozenset({'Milk'})     : [0, 1, 2, 4, 5, 7, 9, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20]
frozenset({'Apple'})    : [0, 3, 4, 6, 8, 10, 11, 13, 14, 15, 16, 17, 18, 19, 20]

```

2. Cài đặt thuật toán

Tóm tắt lý thuyết

- C_k : k-item candidate itemset (tập ứng viên có k-hạng mục)
- F_k : k-item frequent itemset (được lọc ra từ C_k với các giá trị thỏa mãn ‘minimum_support’)

Để áp dụng thuật toán *Apriori* ta cần khởi tạo trước tập C_1 từ đó tạo tập F_1

Hàm `create_candidate_1_itemsets(X)` tạo tập C_1 từ data X

```

1    # Vì C_1 chỉ có một item duy nhất nên ta chỉ cần lấy các key của database X
2    # Vì X là database dạng vertical nên các key của nó chính là các itemsets (tập hạng mục)
3    def create_candidate_1_itemsets(X):
4        c1 = []
5        for transaction in X:
6            c1.append(transaction)
7        return c1

```

Hàm `find_num_transactions(X)` tính tập F_K

Khác với *Horizontal Apriori*, thuật toán *Vertical Apriori* tính support sử dụng chiều dài của danh sách tid của itemset (nếu itemset có nhiều hơn một hạng mục thì danh sách tid chính là tập giao của các item đơn lẻ trong itemset)

```

1    # Hàm trả về tổng số giao dịch (với dữ liệu ngang thì nó là số dòng)
2    # Vì đây là dữ liệu dọc nên nó sẽ là số giao dịch lớn nhất trong danh sách tid
3    def find_num_transactions(X):
4
5        num_transactions = 0
6        for item in X:
7            for tid in X[item]:
8                if tid > num_transactions:
9                    num_transactions = tid
10
11    # Trả về tổng số giao dịch + 1 (vì Python bắt đầu đếm từ 0)
12    return num_transactions + 1
13
14    def create_frequent_itemsets(X, ck, min_support):

```

```

15     # Danh sách tập hợp mục thường xuyên
16     freq_item = []
17     # Tạo từ điển để lưu trữ hỗ trợ của mỗi tập hợp mục
18     item_support = {}
19     # Tìm tổng số giao dịch trong tập dữ liệu
20     num_transactions = find_num_transactions(X)
21
22     # Duyệt qua từng tập hạng mục ứng viên (tập ck)
23     for items in ck:
24         # tid_list lưu trữ các tid chung của các tập hạng mục
25         tid_list = None
26         # Duyệt qua mỗi mục trong tập hạng mục
27         for item in items:
28             # Chuyển về dạng frozenset đảm bảo tính duy nhất và độc lập
29             item = frozenset([item])
30             # Nếu tid_list = None (tức là hạng mục đầu tiên)
31             if tid_list is None:
32                 # Khởi tạo set với các hạng mục đầu tiên
33                 tid_list = set(X[item])
34             else:
35                 # Thực hiện phép giao lấy tid chung của các hạng mục trong tập hạng mục
36                 tid_list = tid_list.intersection(X[item])
37             # Tính support = số tid chung / tổng giao dịch
38             item_support[items] = len(tid_list) / num_transactions
39
40     # Duyệt qua từng tập hạng mục và kiểm tra support
41     for item in item_support:
42         # Nếu hỗ trợ lớn hơn hoặc bằng min_support, thêm vào danh sách tập hạng mục thường xuyên
43         if item_support[item] >= min_support:
44             freq_item.append(item)
45
46     # Trả về danh sách tập hạng mục thường xuyên và từ điển hỗ trợ
47     return freq_item, item_support

```

Hàm `create_candidate_k_itemsets(freq_item, k` tính C_k từ F_{k-1}

- Nếu $k = 1$, hàm tính C_2
- Nếu $k \neq 1$, hàm tính C_i với $i \neq 1$ (C_1 đã được tính từ trước)

```

1  # Hàm hỗ trợ tính tổ hợp từ tập hợp các phần tử cho trước
2  from itertools import combinations
3
4  # Hàm tạo k-item candidate itemsets
5  def create_candidate_k_itemsets(freq_item, k):
6      ck = []
7
8      # Với k = 1, tạo 2-item candidate itemsets (c2)
9      if k == 1:
10         for f1, f2 in combinations(freq_item, 2):
11             item = f1 | f2 # union of two sets
12             ck.append(item)
13     # Với k != 1
14     else:
15         for f1, f2 in combinations(freq_item, 2):
16
17             # Nếu k-itemset có k - 1 phần tử chung
18             # sẽ lấy phần hợp của chúng vào k+1-itemset
19             intersection = f1 & f2
20             if len(intersection) == k - 1:

```

```

21         item = f1 | f2
22         if item not in ck:
23             ck.append(item)
24     return ck

```

Từ các hàm con trên, ta có hàm `apriori_from_scratch(X, min_support, min_combination=1, max_combination=None)`:

- **X** là tập dữ liệu chứa các giao dịch (vertical)
- **min_support** là ngưỡng hỗ trợ tối thiểu để một tập hạng mục được coi là thường xuyên.
- **max_combination** chỉ trả về các giao dịch có tối đa **max_combination** (mặc định bằng tất cả số hạng mục mà database có).

```

1  def apriori_from_scratch(X, min_support, min_combination=1, max_combination=None):
2
3      # Khởi tạo max_combination nếu không được cung cấp
4      if max_combination is None:
5          max_combination = len(X)
6
7      # Tạo tập hợp ứng viên 1-item (c1)
8      c1 = create_candidate_1_itemsets(X)
9      # Tạo tập hợp itemsets thường xuyên f1 từ c1
10     freq_item, item_support_dict = create_frequent_itemsets(X, c1, min_support)
11     # Thêm f1 vào freq_items (tập hợp các itemsets thường xuyên)
12     freq_items = [freq_item]
13
14     # Bắt đầu lặp từ k = 1
15     k = 1
16     while True:
17         # Lấy tập hợp itemsets thường xuyên fk từ freq_items (được lưu dưới dạng list, bắt
18         # đầu từ index 0)
19         freq_item = freq_items[k-1]
20         # Tính tập hợp ứng viên ck+1 từ fk
21         ck = create_candidate_k_itemsets(freq_item, k)
22         # Tính tập hợp itemsets thường xuyên fk+1 từ ck+1 và từ điển hỗ trợ
23         freq_item, item_support = create_frequent_itemsets(X, ck, min_support)
24         # Nếu fk+1 trống => dừng
25         # Kiểm tra điều kiện dừng để tránh thêm list rỗng vào freq_items
26         if len(freq_item) == 0:
27             break
28
29         # Thêm fk+1 vào freq_items
30         freq_items.append(freq_item)
31         # Cập nhật từ điển hỗ trợ tổng
32         item_support_dict.update(item_support)
33         # Tăng k để tiếp tục lặp
34         k += 1
35
36     # Chọn các tập hợp itemsets thường xuyên trong khoảng min_combination đến max_combination
37     freq_items = freq_items[min_combination-1:max_combination-1]
38
39     return freq_items, item_support_dict

```

3. So sánh kết quả

Khởi tạo một số tham số (giống với thông số sử dụng hàm thư viện)

```

1  X = convert_vertical(records)
2  min_support = 0.5

```



```

3 min_combination = 2
4 max_combination = len(X)

```

Tính kết quả

```

1 freq_items, item_support_dict = apriori_from_scratch(X, min_support, min_combination,
2 max_combination)
3 freq_items

```

```

[[frozenset({'Bread', 'Wine'}),
  frozenset({'Butter', 'Wine'}),
  frozenset({'Milk', 'Wine'}),
  frozenset({'Apple', 'Wine'}),
  frozenset({'Bread', 'Butter'}),
  frozenset({'Bread', 'Milk'}),
  frozenset({'Apple', 'Bread'}),
  frozenset({'Butter', 'Milk'}),
  frozenset({'Apple', 'Butter'}),
  frozenset({'Apple', 'Milk'})],
 [frozenset({'Bread', 'Milk', 'Wine'}),
  frozenset({'Bread', 'Butter', 'Milk'})]]

```

Để nhìn giống kết quả hàm thư viện ta format lại kết quả

```

1 def tidy(freq_items, item_support_dict):
2     supports = []
3     itemsets = []
4     for itemset in freq_items:
5         for item in itemset:
6             supports.append(item_support_dict[item])
7             itemsets.append(item)
8     df = pd.DataFrame({'Itemset': itemsets, 'Support': supports})
9     df_sorted = df.sort_values(by=['Support'], ascending=False)
10    return df_sorted
11
12 tidy(freq_items, item_support_dict)

```

	Itemset	Support
2	(Wine, Milk)	0.636364
0	(Bread, Wine)	0.590909
4	(Bread, Butter)	0.590909
5	(Bread, Milk)	0.590909
7	(Butter, Milk)	0.590909
6	(Bread, Apple)	0.545455
1	(Butter, Wine)	0.500000
3	(Apple, Wine)	0.500000
8	(Apple, Butter)	0.500000
9	(Apple, Milk)	0.500000
10	(Bread, Wine, Milk)	0.500000
11	(Bread, Butter, Milk)	0.500000

Tuy thứ tự có đôi chút khác nhưng nhìn chung kết quả trùng khớp với thư viện **pyECLAT**
Thuật toán tự cài đặt đã hoạt động đúng với mong đợi