

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



BTTH TUẦN 4: KHAI THÁC DỮ LIỆU

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

Ngày 28 tháng 4 năm 2024

1. Load thư viện và data

```
1 import numpy as np
2 import pandas as pd
3 from mlxtend.frequent_patterns import apriori, association_rules
4 from mlxtend.preprocessing import TransactionEncoder
5
6 #Load data
7 df = pd.read_csv('data.csv', header = None)
8 df
```

	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 df về dạng list

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

Trong list có các giá trị 'nan' đây không phải là hạng mục của giao dịch mà là giá trị 'NaN' khi chuyển sang dạng chuỗi

```
1 # Ta loại bỏ nó
2 records = [[item for item in i if item != 'nan'] for i in records]
3 records
```

Kết quả:

```
[['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']]
```

Từ bộ data trên ta tiến hành áp dụng thuật toán apriori

2. Hàm Apriori từ thư viện

a) *Apriori*

Để áp dụng được thư viện ta cần mã hoá nó thành dạng ‘transaction’

```
1 te = TransactionEncoder()
2 te_ary = te.fit(records).transform(records)
3 df1 = pd.DataFrame(te_ary, columns = te.columns_)
4 df1
```

Kết quả:

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

Kết quả thuật toán Apriori với `min_support = 0.6`

```
1 frequent_itemsets = apriori(df1, min_support = 0.6, use_colnames = True)
2 frequent_itemsets
```

	support	itemsets
0	0.681818	(Apple)
1	0.727273	(Bread)
2	0.681818	(Butter)
3	0.636364	(Chips)
4	0.772727	(Milk)
5	0.727273	(Wine)
6	0.636364	(Milk, Wine)

b) Association Rules

Luật kết hợp với `metric = "support"`

```
1 rules = association_rules(frequent_itemsets, metric = "support", support_only = True,
2                           min_threshold=0.1)
3 rules = rules[['antecedents', 'consequents', 'support', 'confidence', 'lift']]
4 print(rules)
```

	antecedents	consequents	support	confidence	lift
0	(Milk)	(Wine)	0.636364	NaN	NaN
1	(Wine)	(Milk)	0.636364	NaN	NaN

Với `metric = "confidence"`

```
1 rules = association_rules(frequent_itemsets, metric = "confidence",
2                           min_threshold=0.1)
3 rules = rules[['antecedents', 'consequents', 'support', 'confidence', 'lift']]
4 print(rules)
```

	antecedents	consequents	support	confidence	lift
0	(Milk)	(Wine)	0.636364	0.823529	1.132353
1	(Wine)	(Milk)	0.636364	0.875000	1.132353

3. Hàm Apriori từ đầu (không dùng thư viện)

Ta có cơ sở lý thuyết như sau:

- C_k : k-item candidate itemset (tập ứng viên có k-hạng mục)

- F_k : k-item frequent itemset (tập hạng mục thường xuyên được lọc ra bởi C_k với 'support \geq minimum_support')

4. Apriori

Vì tập C_1 (tập ứng viên một hạng mục) là các item duy nhất của itemsets nên ta tạo nó trước

```
1 def create_candidate_1_itemsets(X):
2     """
3     Duyệt lần lượt các giao dịch, lấy ra các hạng mục duy nhất thêm vào list c1
4     """
5
6     c1 = []
7     for transaction in X:
8         for t in transaction:
9             # frozenset là tập không thay đổi (immutable set) tức nó không cho phép chỉnh sửa, thay d
10            ổi hay giá trị lặp
11            t = frozenset([t])
12            if t not in c1:
13                c1.append(t)
14
15     return c1
```

Hàm tính F_k

```
1 def create_frequent_itemsets(X, ck, min_support):
2     """
3     Lọc từ candidate itemset (ck) ra frequent itemset (fk)
4     với min_support
5     """
6
7     # Lặp từng giao dịch và đếm số lần xuất hiện của mỗi hạng mục
8     item_count = {}
9     for transaction in X:
10        for item in ck:
11            if item.issubset(transaction):
12                if item not in item_count:
13                    item_count[item] = 1
14            else:
15                item_count[item] += 1
16
17     n_row = len(X)
18     freq_item = []
19     item_support = {}
20
21     # Nếu support của mỗi hạng mục lớn hơn hoặc bằng
22     # min_support thì thêm nó vào freq_item (fk)
23     # Lưu lại support của các hạng mục (để so sánh)
24
25     for item in item_count:
26         support = item_count[item] / n_row
27         if support >= min_support:
28             freq_item.append(item)
29
30     item_support[item] = support
31
32     return freq_item, item_support
```

Hàm tính C_k (C_1 tính riêng)

```
1 # Hàm 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
14    # Với k != 1
15    else:
16        for f1, f2 in combinations(freq_item, 2):
17
18            # Nếu k-itemset có k - 1 phần tử chung
19            # sẽ lấy phần hợp của chúng vào k+1-itemset
20            intersection = f1 & f2
21            if len(intersection) == k - 1:
22                item = f1 | f2
```

```

22         if item not in ck:
23             ck.append(item)
24     return ck

```

Từ các hàm con tạo hàm ‘apriori_from_scratch’ tổng hợp với hai tham số

- ‘X’ là tập hạng mục (transactions)
- ‘min_support’ là support tối thiểu

```

1 def apriori_from_scratch(X, min_support):
2
3     # tạo trước 1-item candidate itemsets (c1)
4     c1 = create_candidate_1_itemsets(X)
5     # tạo f1 (frequent itemsets) từ c1
6     freq_item, item_support_dict = create_frequent_itemsets(X, c1, min_support)
7     # thêm f1 vào freq_items (tập hạng mục thường xuyên)
8     freq_items = [freq_item]
9
10    # lặp từ k = 1
11    k = 1
12    while True:
13        # lấy ra fk từ tập freq_items (ở đây là list nên idx bắt đầu từ 0)
14        freq_item = freq_items[k-1]
15        # tính ck+1 từ fk
16        ck = create_candidate_k_itemsets(freq_item, k)
17        # tính fk+1 từ ck+1 và support dict
18        freq_item, item_support = create_frequent_itemsets(X, ck, min_support)
19        # nếu fk+1 trống => dừng
20        # tạo dk dừng ở đây để tránh thêm list rỗng vào freq_items
21        if len(freq_item) == 0:
22            break
23        # thêm fk+1 vào freq_items
24        freq_items.append(freq_item)
25        # cập nhật support dict tổng
26        item_support_dict.update(item_support)
27        # tăng k lên một tiếp tục lặp
28        k += 1
29
30    return freq_items, item_support_dict

```

So sánh kết quả Chạy ‘apriori_from_scratch’ với ‘min_support = 0.6’ (giống với hàm thư viện)

```

1 freq_items, item_support_dict = apriori_from_scratch(records, 0.6)
2 freq_items

```

```

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

```

Ta có thể thấy kết quả trùng với thư viện

Nhằm kiểm chứng tính đúng đắn, ta in ra các giá trị support sắp xếp tăng dần

```

1 dict(sorted(item_support_dict.items(), key=lambda item: item[1], reverse=True))

```

```
{frozenset({'Milk'}): 0.7727272727272727,
 frozenset({'Wine'}): 0.7272727272727273,
 frozenset({'Bread'}): 0.7272727272727273,
 frozenset({'Butter'}): 0.6818181818181818,
 frozenset({'Apple'}): 0.6818181818181818,
 frozenset({'Chips'}): 0.6363636363636364,
 frozenset({'Milk', 'Wine'}): 0.6363636363636364,
 frozenset({'Bread', 'Wine'}): 0.5909090909090909,
 frozenset({'Bread', 'Butter'}): 0.5909090909090909,
 frozenset({'Bread', 'Milk'}): 0.5909090909090909,
 frozenset({'Butter', 'Milk'}): 0.5909090909090909,
 frozenset({'Apple', 'Bread'}): 0.5454545454545454,
 frozenset({'Butter', 'Wine'}): 0.5,
 frozenset({'Apple', 'Wine'}): 0.5,
 frozenset({'Apple', 'Butter'}): 0.5,
 frozenset({'Apple', 'Milk'}): 0.5,
 frozenset({'Chips', 'Milk'}): 0.4545454545454545,
 frozenset({'Apple', 'Chips'}): 0.4545454545454545,
 frozenset({'Chips', 'Wine'}): 0.4090909090909091,
 frozenset({'Bread', 'Chips'}): 0.4090909090909091,
 frozenset({'Butter', 'Chips'}): 0.4090909090909091}
```

Chương trình em viết chỉ lấy những set có support ≥ 0.6
Để dễ nhìn hơn em có viết hàm chuyển nó sang df

```
1 def tidy(freq_items, item_support_dict):
2     supports = []
3     itemsets = []
4     for itemset in freq_items:
5         for set in itemset:
6             supports.append(item_support_dict[set])
7             itemsets.append(set)
8     return pd.DataFrame({'Itemset': itemsets, 'Support': supports})
9
10 tidy(freq_items, item_support_dict)
```

	Itemset	Support
0	(Wine)	0.727273
1	(Chips)	0.636364
2	(Bread)	0.727273
3	(Butter)	0.681818
4	(Milk)	0.772727
5	(Apple)	0.681818
6	(Wine, Milk)	0.636364

a) Association Rules

Từ 'freq_items' và 'item_dict_support' trả về từ hàm 'apriori' trên ta tạo hàm tính metric chuẩn bị cho hàm tạo luật kết hợp

Hàm `compute_metric` này nhận vào các tham số như sau:

‘`freq_items`’: Các tập hạng mục thường xuyên.

‘`item_support_dict`’: Từ điển lưu trữ hỗ trợ của từng tập hạng mục.

‘`freq_set`’: Tập hạng mục hiện tại cần xem xét để tạo ra các luật.

‘`subsets`’: Các phần tử con của `freq_set`.

‘`min_threshold`’: Ngưỡng tối thiểu cho metric được chỉ định.

‘`metric`’: Loại metric để đánh giá luật kết hợp (Confidence, Support, Lift).

```
1 def compute_metric(freq_items, item_support_dict, freq_set, subsets, min_threshold, metric):
2
3     # Trả lỗi nếu metric không hợp lệ
4     if metric not in ("Confidence", "Support", 'Lift'):
5         raise ValueError("Invalid metric value")
6
7     rules = []
8     right_hand_side = []
9
10    for rhs in subsets:
11        # Tạo phần trái của quy tắc
12        lhs = freq_set - rhs
13
14        # Kiểm tra xem phần tử hợp nhất của quy tắc có đủ phổ biến không
15        if lhs in item_support_dict and rhs in item_support_dict:
16
17            # Tính các metric
18            supp = item_support_dict[freq_set]
19            conf = item_support_dict[freq_set] / item_support_dict[lhs]
20            lift = conf / item_support_dict[rhs]
21
22            # Chọn metric để so sánh với min_threshold
23            if metric == "confidence":
24                metric_value = conf
25            elif metric == "lift":
26                metric_value = lift
27            else:
28                metric_value = supp
29
30            # Thêm quy tắc vào danh sách nếu metric vượt qua ngưỡng
31            if metric_value >= min_threshold:
32                rules_info = lhs, rhs, supp, conf, lift
33                rules.append(rules_info)
34                right_hand_side.append(rhs)
35
36    return rules, right_hand_side
```

Hàm ‘`association_rules_from_scratch`’ tạo ra luật kết hợp từ đầu (không dùng thư viện)

Các tham số lần lượt là ‘`(freq_items, item_support_dict)`’ từ hàm ‘`apriori`’, ‘`min_threshold`’ là ngưỡng giữ lại luật kết hợp tùy thuộc vào ‘`metric`’ cụ thể

```
1 def association_rules_from_scratch(freq_items, item_support_dict, min_threshold, metric= "confidence"
2 ):
3
4     association_rules = []
5
6     # Lập trong tập hạng mục thường xuyên (freq_items)
7     # Lập từ phần tử thứ 2, vì phần tử đầu là c1 chứa các
8     # hạng mục đơn không thể tạo luật kết hợp
9
10    for idx, freq_item in enumerate(freq_items[1:(len(freq_items))]):
11        for freq_set in freq_item:
12
13            # Tạo luật cho các hạng mục đơn trước
14            # Tạo tập con không lặp
15            subsets = [frozenset([item]) for item in freq_set]
16
17            # Tính rules (đã qua sàng lọc theo ngưỡng và metric) và tập hạng mục phía phải (
18            # consequents)
19            rules, right_hand_side = compute_metric(freq_items, item_support_dict,
20                                                    freq_set, subsets, min_threshold, metric)
21            association_rules.extend(rules)
22
23            # tạo luật từ 3-itemset trở đi
24            if idx != 0:
25                k = 0
26                while len(right_hand_side[0]) < len(freq_set) - 1:
27                    ck = create_candidate_k_itemsets(right_hand_side, k = k)
```



```

25         rules, right_hand_side = compute_metric(freq_items, item_support_dict,
26                                                  freq_set, ck, min_threshold, metric)
27         association_rules.extend(rules)
28         k += 1
29
30     return association_rules

```

Kiểm tra kết quả

```

1 association_rules = association_rules_from_scratch(freq_items, item_support_dict, min_threshold =
    0.1, metric="confidence")

1 def tidy(list):
2     return pd.DataFrame(list, columns=['Antecedents ', 'Consequents', 'Support', 'Confidence', 'Lift'
    ])
3
4 tidy(association_rules)

```

	Antecedents	Consequents	Support	Confidence	Lift
0	(Milk)	(Wine)	0.636364	0.823529	1.132353
1	(Wine)	(Milk)	0.636364	0.875000	1.132353

Giống với hàm thư viện