

**TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN, ĐHQG - HCM**  
**KHOA CÔNG NGHỆ THÔNG TIN**



## **LAB 02**

### **CÀI ĐẶT VÀ ỨNG DỤNG 2 THUẬT TOÁN: APRIORI VÀ TREE PROJECTION CHO VIỆC KHAI THÁC TẬP PHỔ BIẾN TRÊN CÁC TẬP DỮ LIỆU**

**Học phần:** Khai thác dữ liệu và ứng dụng

**Lớp:** 19\_21

**Giáo viên hướng dẫn:** Nguyễn Ngọc Đức

**Sinh viên:** Lê Kiệt

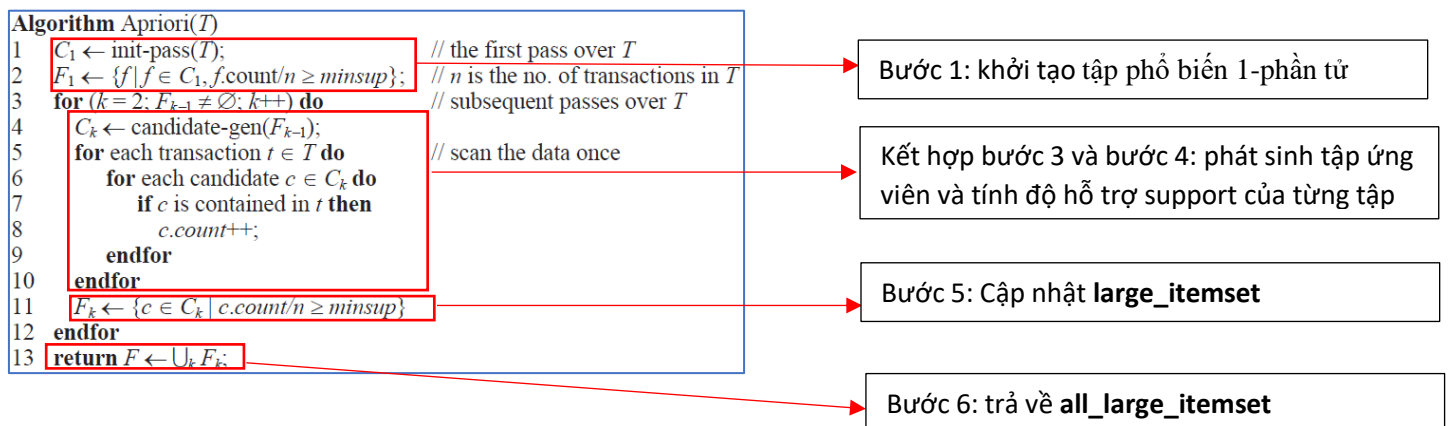
**MSSV:** 19120554

# MỤC LỤC

I.	THUẬT TOÁN APRIORI.....	1
II.	THUẬT TOÁN TREE PROJECTION .....	2
1.	Ý tưởng.....	2
2.	Định nghĩa cấu trúc dữ liệu phù hợp.....	3
3.	Các bước thực hiện .....	3
III.	ỨNG DỤNG 2 THUẬT TOÁN TRÊN CÁC TẬP DATASET.....	4
1.	Dataset foodmart.txt .....	4
2.	Dataset chess.txt .....	4
3.	Dataset mushrooms.txt .....	4
4.	Dataset retail.txt.....	4
IV.	SO SÁNH 2 THUẬT TOÁN .....	5
V.	ỨNG DỤNG TRONG KHAI THÁC DỮ LIỆU THỰC TẾ.....	5
VI.	THAM KHẢO.....	5

## I. THUẬT TOÁN APRIORI

- Sử dụng hàm **apriori(database, min\_sup)** để thực hiện thuật toán. Hàm nhận 2 tham số là **database** và **min\_sup** lần lượt là cơ sở dữ liệu giao tác và độ hỗ trợ tối thiểu; trả về **all\_large\_itemset** là 1 danh sách (Array) tất cả các tập phổ biến thỏa **min\_sup**. Sau đây là các bước mà hàm này thực hiện
  - Bước 1:** khởi tạo tập **large\_itemset** là danh sách các tập phổ biến 1-phần tử **one\_itemset** đã được sắp xếp theo thứ tự từ điển thỏa **min\_sup**,
  - Bước 2:** lưu **large\_itemset** này vào biến kết quả **all\_large\_itemset**
  - Bước 3:** phát sinh **candidate\_list** là danh sách các tập ứng viên sinh ra từ tập phổ biến **large\_itemset** nhờ hàm **candidate\_gen(large\_itemset)**
  - Bước 4:** cập nhật support của các tập ứng viên trong với mỗi dòng/giao tác trong database, kiểm tra:
    - với mỗi ứng viên trong tập ứng viên **candidate\_list**, kiểm tra:
      - if (ứng viên  $\subseteq$  giao tác):
      - support(ứng viên) += 1
  - Bước 5:** kết thúc bước 4, lọc lại các tập ứng viên thỏa **min\_sup**. Cập nhật **large\_itemset** là các tập ứng viên mới lọc này, nếu **large\_itemset** là tập rỗng (tức không còn tập ứng viên nào thỏa **min\_sup**) thì dừng thuật toán; ngược lại, quay lại bước 2
  - Bước 6:** trả về **all\_large\_itemset**
- Trong hàm trên có sử dụng thuật toán phát sinh các tập ứng viên là hàm **candidate\_gen(large\_itemset)**. Hàm nhận 1 **large\_itemset** là danh sách các tập phổ biến k-phần tử có thứ tự từ điển và trả về **candidate\_list** là danh sách các tập ứng viên k+1 phần tử. Sau đây là các bước mà hàm này thực hiện
  - Bước 1:** với mỗi cặp tập phổ biến k-phần tử (**large\_itemset[i]**, **large\_itemset[j]**),  $i \neq j$  sao cho k-1 phần tử đầu của từng tập phổ biến trong cặp phải giống nhau  $\rightarrow$  tiến hành bước 2
  - Bước 2:** dùng toán tử hợp (union) để hợp cặp này lại, kết quả sau khi hợp lưu vào **u**.
  - Bước 3:** tia nhánh: xét tất cả tập con (k-1)-phần tử của **u**, có tập con nào là tập không phổ biến, tức không nằm trong **large\_itemset** hay không. Nếu tồn tại 1 tập con như vậy thì không chấp nhận tập ứng viên **u**; ngược lại, nếu không tồn tại bất kỳ tập con nào không phổ biến thì thêm **u** vào **candidate\_list** (Dựa vào quy tắc: tập con không phổ biến  $\rightarrow$  tập cha chứa nó không phổ biến)
    - Lưu ý ở bước 3:** khi xét tập con có thuộc vào tập **large\_itemset** hay không. Về mặt lập trình trong Julia, cần sắp xếp lại tập con theo thứ tự bảng chữ cái (nếu là chữ) hoặc từ bé tới lớn (nếu là số). Ví dụ,  $['C', 'T'] \in [ ['A', 'B'], ['C', 'T'] ]$  trong khi  $['T', 'C'] \notin [ ['A', 'B'], ['C', 'T'] ]$
  - Bước 4:** Sau khi xét xong tất cả các cặp **large\_itemset[i]** và **large\_itemset[j]** thỏa bước 1 thì **candidate\_list** cũng đã chứa các tập ứng viên (k+1)-phần tử hợp lệ  $\rightarrow$  trả về **candidate\_list**
- 2 hàm trên được cài đặt dựa vào mã giả sau (trích từ slide bài giảng trên lớp)



```

Function candidate-gen( $F_{k-1}$ )
1   $C_k \leftarrow \emptyset$ ; // initialize the set of candidates
2  forall  $f_1, f_2 \in F_{k-1}$  // find all pairs of frequent itemsets
3     with  $f_1 = \{i_1, \dots, i_{k-2}, i_{k-1}\}$  // that differ only in the last item
4     and  $f_2 = \{i_1, \dots, i_{k-2}, i'_{k-1}\}$ 
5     and  $i_{k-1} < i'_{k-1}$  do // according to the lexicographic order
6          $c \leftarrow \{i_1, \dots, i_{k-2}, i'_{k-1}\}$ ; // join the two itemsets  $f_1$  and  $f_2$ 
7          $C_k \leftarrow C_k \cup \{c\}$ ; // add the new itemset  $c$  to the candidates
8  for each  $(k-1)$ -subset  $s$  of  $c$  do
9     if ( $s \notin F_{k-1}$ ) then
10        delete  $c$  from  $C_k$ ; // delete  $c$  from the candidates
11    endfor
12 endfor
13 return  $C_k$ ; // return the generated candidates

```

Kết hợp bước 1 và bước 2

Bước 3: tỉa nhánh

## II. THUẬT TOÁN TREE PROJECTION

### 1. Ý tưởng

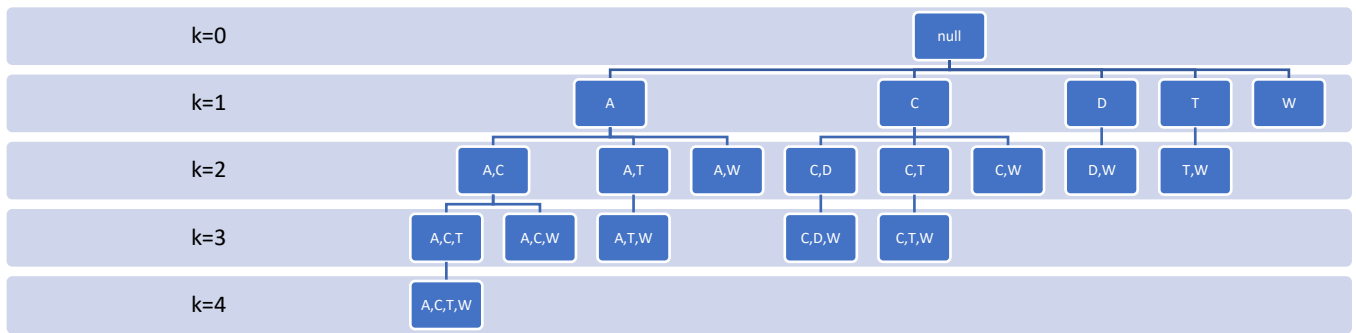
- Thuật toán Tree Projection được dựa theo thuật toán Apriori để tạo ra 1 cây được cấu trúc theo thứ tự từ điển (lexicographic tree) với mỗi node của cây là 1 tập phổ biến. Ta định nghĩa 1 cấp  $k$  trong cây ( $k$  nguyên dương,  $k \geq 0$ ) là cấp mà ở đó các tập phổ biến (cũng chính là các node) có  $k$  phần tử. Như vậy node gốc (root node) ở cấp 0 sẽ có 0 phần tử hay chính là node rỗng; các node ở cấp 1 sẽ có 1 phần tử; các node ở cấp 2 sẽ có 1 phần tử; v.v.
  - o Minh họa: 19 tập phổ biến trích từ slide “**FrequentItemsetMining.pdf**” trang 22 được cung cấp

```

1  [nothing]
2  └─ ["A"]
3     └─ ["A", "C"]
4        └─ ["A", "C", "T"]
5           └─ ["A", "C", "T", "W"]
6              └─ ["A", "C", "W"]
7                 └─ ["A", "T"]
8                    └─ ["A", "T", "W"]
9                       └─ ["A", "W"]
10 └─ ["C"]
11    └─ ["C", "D"]
12       └─ ["C", "D", "W"]
13          └─ ["C", "T"]
14             └─ ["C", "T", "W"]
15                └─ ["C", "W"]
16 └─ ["D"]
17    └─ ["D", "W"]
18 └─ ["T"]
19    └─ ["T", "W"]
20 └─ ["W"]

```

Trích từ file *sample\_tree.txt*



Cây được vẽ lại từ hình trên

- Các node là con của 1 node khi giữa node cha và mỗi node con có chung các phần tử của node cha cũng như node con phải nhiều hơn node cha 1 phần tử. Ví dụ như hình minh họa trên, node con của A là [A,W], [A,T], và [A,C] vì có chung 1 phần tử A và 3 node con này nhiều hơn 1 phần tử so với node cha A; node con của [A,W] là [A,T,W] và [A,C,W] vì có chung [A,W] và 2 node con nhiều hơn 1 phần tử so với node cha
- Tất cả các node đều có thể có nhiều con nhưng chỉ có 1 node cha, node nào đã có cha rồi thì không thêm cha khác vào nữa. Ví dụ node [A,W] đã có cha là A thì không cần thêm cha là W nữa

## 2. Định nghĩa cấu trúc dữ liệu phù hợp

- Để thực hiện vẽ cây với mỗi node là 1 tập phổ biến, ta cần định nghĩa 1 kiểu cấu trúc mới mang tên **Node**. 1 thể hiện Node (gọi là 1 node) sẽ có 4 thuộc tính để lưu trữ thông tin:
  - o **itemset**: lưu trữ tập phổ biến
  - o **children**: lưu trữ các tập phổ biến con của 1 node
  - o **sup**: độ hỗ trợ (support) của tập phổ biến **itemset**
  - o **parent**: lưu trữ đúng 1 node cha của tập phổ biến hiện tại
- Ngoài ra cấu trúc Node có constructor mặc định được định nghĩa như hình dưới

```

1 mutable struct Node
2   itemset # ex: ['a', 'b']
3   children::Vector{Any}
4   sup
5   parent
6
7   # constructor
8   Node(its, chi = [], su = 0, par = nothing) = new(its, chi, su, par)
9 end
  
```

## 3. Các bước thực hiện

- Thuật toán Tree Projection có các bước về bản chất giống với thuật toán Apriori đã nêu và điều chỉnh để thích hợp với kiểu dữ liệu struct Node; điểm khác biệt giữa cài đặt 2 thuật toán được mô tả bằng chữ **xanh**. Ta có hàm **apriori\_tree(database, min\_sup)** nhận 2 tham số là **database** và **min\_sup** lần lượt là cơ sở dữ liệu giao tác và độ hỗ trợ tối thiểu; trả về **tree** là cây với các Node lồng nhau theo quan hệ cha-con như ở mục 1 - Ý tưởng. Các bước thực hiện
  - o Bước 1:
    - Khởi tạo tập phổ biến 1-phần tử **one\_itemset** thỏa min\_sup. Sort **one\_itemset** theo thứ tự từ điển
    - Khởi tạo biến kết quả **tree** = Node(itemset = [nothing], children = one\_itemset). Ban đầu **tree** được khởi tạo chỉ có root node = null và con của root node là tập one\_itemset thỏa min\_sup
    - Khởi tạo biến con trở chạy dùng để duyệt qua các cấp trong cây **itr** = tree.children
  - o Bước 2: phát sinh danh sách các tập ứng viên **candidate\_nodes** nhờ hàm **candidate\_nodes\_gen(itr)**
  - o Bước 3: cập nhật support của các node ứng viên trong **candidate\_nodes** với mỗi **transaction** trong database, kiểm tra:
    - với mỗi **candidate** trong tập **candidate\_nodes**, kiểm tra:

if (candidate.itemset  $\subseteq$  transaction):

candidate.sup += 1

- Bước 4: kết thúc bước 3, lọc lại các node ứng viên thỏa **min\_sup** và lưu vào biến tạm là **large\_itemset**. Ngược lại, đối với các node ứng viên không thỏa **min\_sup** thì “xóa” các node này khỏi cây. Việc “xóa” node thực chất là việc cắt đứt quan hệ **con-cha** của node đó với node cha (node.parent = []) và cắt đứt quan hệ **cha-con** của node cha và node đó (node\_cha.children = node\_cha.children \ {node\_con})
  - Bước 5: **cập nhật con trở itr = large\_itemset**. Nếu itr rỗng thì thuật toán kết thúc; ngược lại quay lại bước 2
- Trong hàm trên có sử dụng thuật toán phát sinh các tập ứng viên là hàm **candidate\_nodes\_gen (large\_itemset\_nodes)**. Hàm nhận 1 **large\_itemset\_nodes** là danh sách các node phổ biến k-phần tử và trả về **candidate\_nodes** là danh sách các node ứng viên (k+1)-phần tử. Hàm này cũng tương tự hàm **candidate\_gen** của thuật toán Apriori; điểm khác biệt chính nằm ở bước 3 khi thêm 1 node ứng viên hợp lệ sẽ có 2 thao tác chính liên quan tới cây. Sau đây là các bước mà hàm này thực hiện
- **Bước 1:** với mỗi cặp tập phổ biến k-phần tử (large\_itemset[i], large\_itemset[j]),  $i \neq j$  sao cho k-1 phần tử đầu của từng tập phổ biến trong cặp phải giống nhau  $\rightarrow$  tiến hành bước 2
  - **Bước 2:** dùng toán tử hợp (union) để hợp cặp này lại, kết quả sau khi hợp lưu vào **u**
  - **Bước 3:** tia nhánh: xét tất cả tập con (k-1)-phần tử của **u**, có tập con nào là tập không phổ biến, tức không nằm trong **large\_itemset\_nodes** hay không. Nếu tồn tại 1 tập con như vậy thì không chấp nhận node ứng viên **u**. Ngược lại, chấp nhận **u** và thực hiện 2 thao tác liên quan tới cây:
    - Tạo quan hệ **cha-con**: Bổ sung node **u** vào tập node children của node **large\_itemset\_nodes[i]** (ưu tiên lưu vào **large\_itemset\_nodes[i]** hơn vào **large\_itemset\_nodes[j]** để thỏa mãn thứ tự từ điển)
    - Tạo quan hệ **con-cha**, tức **node\_u.parent = large\_itemset\_nodes[i]**. Vì 1 node chỉ được phép có 1 cha nên không cập nhật **node\_u.parent = large\_itemset\_nodes[j]**
    - Bổ sung node **u** vào **candidate\_nodes**
  - **Bước 4:** Sau khi xét xong tất cả các cặp large\_itemset[i] và large\_itemset[j] thỏa bước 1 thì **candidate\_nodes** cũng đã chứa các node ứng viên (k+1)-phần tử hợp lệ  $\rightarrow$  trả về **candidate\_nodes**

### III. ỨNG DỤNG 2 THUẬT TOÁN TRÊN CÁC TẬP DATASET

- Mỗi dataset có 2 file output: **[tên\_dataset].tree.txt** (output dùng thuật toán tree projection) và **[tên\_dataset].apriori.txt** (output dùng thuật toán apriori). Ví dụ: dataset foodmart.txt có 2 file output là: **foodmart\_tree.txt** và **foodmart\_apriori.txt**

#### 1. Dataset foodmart.txt

- Với min\_sup = 0.002421 = 0.2421%  $\Leftrightarrow$  sup.count  $\geq 10 \rightarrow$  cho ra 980 tập phổ biến, nhưng toàn bộ 980 tập này là tập one\_itemset (check lại)
- Với min\_sup = 0.000965 = 0.0965%  $\Leftrightarrow$  sup.count  $\geq 4 \rightarrow$  cho ra 1557 tập phổ biến, trong đó có 2 tập 2-itemset, còn lại là 1-itemset. min\_sup này cũng cho thấy sự khác biệt về thời gian chạy của 2 thuật toán: thuật toán Apriori tốn ~21 phút để chạy trong khi thuật toán Tree Projection tốn ~17 phút để chạy

#### 2. Dataset chess.txt

- Với min\_sup = 90%  $\rightarrow$  thuật toán tree projection tạo được sâu nhất tới cấp 7 (level=7); cả 2 thuật toán đều cho ra 622 tập phổ biến

#### 3. Dataset mushrooms.txt

- Với min\_sup = 40%  $\rightarrow$  thuật toán tree projection tạo được sâu nhất tới cấp 6 (level=6); cả 2 thuật toán đều cho ra 505 tập phổ biến

#### 4. Dataset retail.txt

- Với min\_sup = 9% = 0.09  $\rightarrow$  thuật toán tree projection tạo được sâu nhất tới cấp 2 (level=2); cả 2 thuật toán đều cho ra 12 tập phổ biến
- Với min\_sup = 5%  $\rightarrow$  thuật toán tree projection tạo được sâu nhất tới cấp 3 (level=3); cả 2 thuật toán đều cho ra 16 tập phổ biến

#### IV. SO SÁNH 2 THUẬT TOÁN

	Apriori	Tree Projection
Điểm mạnh	<ul style="list-style-type: none"><li>- Dễ hiểu và dễ cài đặt</li><li>- Dữ liệu đầu vào không cần gán nhãn</li><li>- Tìm toàn bộ các luật từ dataset ban đầu</li><li>- Dễ biến thể theo từng trường hợp cụ thể</li></ul>	<ul style="list-style-type: none"><li>- Đảm bảo không tạo node trùng trong quá trình phát sinh ứng viên → Tiết kiệm 1 vài bước tính toán so với Apriori nên nhanh hơn Apriori</li><li>- Thừa kế các điểm mạnh của Apriori</li></ul>
Điểm yếu	<ul style="list-style-type: none"><li>- Tốn nhiều thời gian vì 1 trong các lý do sau:<ul style="list-style-type: none"><li>o Ngưỡng min_sup nhỏ → có thể rất nhiều tập phổ biến phát sinh thỏa min_sup</li><li>o Số lượng tập itemset nạp vào là lớn</li><li>o Duyệt lặp đi lặp lại toàn bộ cơ sở dữ liệu nhiều lần và kiểm tra một tập lớn các ứng viên bằng phương pháp so khớp mẫu</li><li>o Số lượng tập hạng mục và luật được phát sinh là khổng lồ → khó khăn cho việc phân tích thông tin</li></ul></li></ul>	<ul style="list-style-type: none"><li>- Cần định nghĩa 1 cấu trúc dữ liệu mới để lưu trữ các node</li><li>- Phức tạp hơn Apriori do cần quan tâm tới quan hệ cha-con giữa các node trong cây với mỗi thao tác thêm/xóa node</li></ul>

#### V. ỨNG DỤNG TRONG KHAI THÁC DỮ LIỆU THỰC TẾ

- Nguồn dữ liệu: <https://www.kaggle.com/irfanasrullah/groceries/version/2>
- File groceries.csv là file miêu tả từng giỏ hàng của người mua tại 1 cửa hàng. Từng giỏ hàng tương ứng với từng dòng trong file. Ta sẽ tiến hành rút các tập phổ biến với min\_sup = 1%
- Output: được thể hiện trong 2 file groceries\_apriori.txt và groceries\_tree.txt

#### VI. THAM KHẢO

- [1] [https://www.academia.edu/12469140/Tree\\_Projection\\_Based\\_Frequent\\_Itemset\\_Mining\\_on\\_Multicore\\_CPUs\\_and\\_GPUs](https://www.academia.edu/12469140/Tree_Projection_Based_Frequent_Itemset_Mining_on_Multicore_CPUs_and_GPUs)
- [2] [https://www.engati.com/glossary/apriori-algorithm?utm\\_content=apriori-algorithm](https://www.engati.com/glossary/apriori-algorithm?utm_content=apriori-algorithm)
- [3] [https://www-users.cse.umn.edu/~kumar001/dmbook/ch5\\_association\\_analysis.pdf](https://www-users.cse.umn.edu/~kumar001/dmbook/ch5_association_analysis.pdf)