

TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN
THÀNH PHỐ HỒ CHÍ MINH



BÀI TẬP 2: KHAI THÁC LUẬT KẾT HỢP

I. THÔNG TIN SINH VIÊNHọ và tên: **TRẦN NHẬT HUY**Mssv: **1612272**Email: nhathuy13598@gmail.comSđt: **0354 878 677****II. BẢNG BÁO CÁO CÔNG VIỆC**

STT	CÁC CÂU HỎI	MỨC ĐỘ HOÀN THÀNH	GHI CHÚ
A	1. Tìm hiểu về phương pháp cải tiến quá trình tìm luật kết hợp từ tập phổ biến	100%	
	2.a. Sử dụng thuật toán Apriori và FP-Growth. Liệt kê tập phổ biến tối đại và tập phổ biến đóng	100%	
	2.b. Tìm tất cả các luật kết hợp	100%	
	2.c. Ứng dụng cải tiến ở câu a. vào câu b.	100%	
B	1. Chuyển dữ liệu trong plants.data sang dạng nhị phân	100%	
	2. Trả lời các câu hỏi	100%	
	3. Áp dụng thuật toán Apriori	100%	
	4. Khai thác tập phổ biến	100%	
	5. Khai thác luật kết hợp	100%	

III. CHI TIẾT BÀI LÀM

A. LÝ THUYẾT

1. Tìm hiểu về một phương pháp cải tiến khai thác luật kết hợp từ tập phổ biến

Nhược điểm của phương pháp khai thác luật kết hợp dựa trên tập phổ biến là khi số lượng tập phổ biến lớn thì luật kết hợp sinh ra sẽ rất lớn. Trong đó, một luật kết hợp này có thể là con của luật kết hợp khác do đó chúng ta sẽ phải tiến hành loại bỏ các luật là con của luật khác

Các phương pháp được đề nghị sử dụng thay thế cho tập phổ biến là khai thác luật kết hợp trên tập phổ biến đóng và tập phổ biến tối đại

Tập phổ biến đóng là tập phổ biến mà không có tập nào bao nó có cùng độ phổ biến

Tập phổ biến tối đại là tập phổ biến mà không có tập nào bao nó cũng là tập phổ biến

Từ định nghĩa ta có : tập phổ biến tối đại \subseteq tập phổ biến đóng \subseteq tập phổ biến tối đại

Tập phổ biến đóng thể hiện đầy đủ thông tin của tập phổ biến cùng với độ hỗ trợ chính xác của nó. Luật kết hợp được lấy ra từ tập phổ biến đóng sẽ nhỏ gọn hơn, dễ quản lý và phân tích. Giả sử, từ tập phổ biến ta tạo ra được tập phổ biến đóng, tuy nhiên tập phổ biến đóng này vẫn còn quá lớn thì khi đó ta sẽ tìm tập phổ biến tối đại. Khai thác tập phổ biến tối đại thích hợp với cơ sở dữ liệu dày đặc

Để khai thác tập phổ biến đóng, ta sử dụng thuật toán **CHARM**. Để khai thác tập phổ biến tối đại thì ta sử dụng thuật toán **GenMax**. Cả 2 thuật toán trên đều sử dụng cây được tạo ra từ thuật toán **Eclat**

2. Cài đặt các thuật toán

a. *Sử dụng thuật toán **Apriori** và **FP-Growth** để tìm luật kết hợp. So sánh. Tìm tập phổ biến tối đại và tập phổ biến đóng*

Thuật toán Apriori và FP-Growth được cài đặt trong các file lần lượt là **Apriori.py** và **FP_Growth.py**

Tập phổ biến khi chạy **Apriori.py** và **FP_Growth.py**

```

def Apriori_Run(T,L,minsup):...
L = Apriori_Run(T,L,minsup)
print("Tập phổ biến là: \n")
for i in L.values():
    print(i)

```

Run: "C:\Users\nhat huy\AppData\Local\Programs\Python\Python37\python.exe" "C:\Users\nhat huy\Desktop\1612272_02\Source\LY THUYET\Apriori.py"

Nhap minsup: 2
 Nhap minconf (0<= minconf <=1): 1
 Tap pho bien la:

```

{'A': 3, 'B': 3, 'C': 2, 'D': 3, 'E': 2, 'F': 2, 'G': 2, 'H': 2, 'J': 2}
{'AB': 2, 'AD': 2, 'AG': 2, 'BD': 2, 'BE': 2, 'BH': 2, 'BJ': 2, 'CD': 2, 'CF': 2, 'DE': 2, 'DH': 2, 'DJ': 2, 'HJ': 2}
{'BDH': 2, 'BDJ': 2, 'BHJ': 2, 'CDF': 2, 'DHJ': 2}
{'BDHJ': 2}

```

Process finished with exit code 0

Figure 1: Tập phổ biến khi sử dụng Apriori.py

```

Tap pho bien la: (['J'], 2)
(['J', 'B'], 2)
(['J', 'D'], 2)
(['J', 'D', 'B'], 2)
(['J', 'H'], 2)
(['J', 'H', 'B'], 2)
(['J', 'H', 'D'], 2)
(['J', 'H', 'D', 'B'], 2)
(['H'], 2)
(['H', 'B'], 2)
(['H', 'D'], 2)
(['H', 'D', 'B'], 2)
(['G'], 2)
(['G', 'A'], 2)
(['F'], 2)
(['F', 'D'], 2)
(['F', 'C'], 2)
(['F', 'C', 'D'], 2)
(['E'], 2)
(['E', 'B'], 2)
(['C'], 2)
(['C', 'D'], 2)
(['D'], 3)
(['D', 'B'], 2)
(['D', 'A'], 2)
(['B'], 3)
(['B', 'A'], 2)
(['A'], 3)

```

Figure 2: Tập phổ biến khi chạy FP_Growth.py

Kết quả của 2 phương pháp này là giống nhau

K-ITEM SET	ITEM
1-item set	A:3, B:3, C:2, D:3, E:2, F:2, G:2, H:2, J:2
2-item set	AB:2, AD:2, AG:2, BD:2, BE:2, BH:2, BJ:2, CD:2, CF:2, DF:2, DH:2, DJ:2, HJ:2
3-item set	BDH:2, BDJ:2, BHJ:2, CDF:2, DHJ:2
4-item set	BDHJ:2

Ta chạy file **Apriori.py** để tìm tập phổ biến tối đại và tập phổ biến đóng

```

"C:\Users\nhat huy\AppData\Local\Programs\Python\Python37\python.exe" "C:/Users/nhat huy/Desktop/1612272_02/Source/LY THUYET/Apriori.py"
Nhap minsup: 2
Nhap minconf (0<= minconf <=1): 1
Tap pho bien la:
{'A': 3, 'B': 3, 'C': 2, 'D': 3, 'E': 2, 'F': 2, 'G': 2, 'H': 2, 'J': 2}
{'AB': 2, 'AD': 2, 'AG': 2, 'BD': 2, 'BE': 2, 'BH': 2, 'BJ': 2, 'CD': 2, 'CF': 2, 'DF': 2, 'DH': 2, 'DJ': 2, 'HJ': 2}
{'BDH': 2, 'BDJ': 2, 'BHJ': 2, 'CDF': 2, 'DHJ': 2}
{'BDHJ': 2}
Luat khai thac bang tap pho bien la:
Tap pho bien dong la:
{'A': 3, 'B': 3, 'C': 2, 'D': 3, 'E': 2, 'F': 2, 'G': 2, 'H': 2, 'J': 2}
{'AB': 2, 'AD': 2, 'AG': 2, 'BE': 2}
{'CDF': 2}
{'BDHJ': 2}
Tap pho bien toi dai la:
{'A': 3, 'B': 3, 'C': 2, 'D': 3, 'E': 2, 'F': 2, 'G': 2, 'H': 2, 'J': 2}
{'AB': 2, 'AD': 2, 'AG': 2, 'BE': 2}
{'CDF': 2}
{'BDHJ': 2}
Process finished with exit code 0
  
```

Figure 3: Tập phổ biến tối đại và tập phổ biến đóng

Ta giữ lại tập **1-item set** để dễ tính **minconf**

Tập phổ biến đóng

K-ITEM SET	ITEM
1-item set	A:3, B:3, C:2, D:3, E:2, F:2, G:2, H:2, J:2
2-item set	AB:2, AD:2, AG:2, BD:2, BE:2
3-item set	CDF:2
4-item set	BDHJ:2

Tập phổ biến tối đại

K-ITEM SET

ITEM

1-item set

A:3, B:3, C:2, D:3, E:2, F:2, G:2, H:2, J:2

2-item set

AB:2, AD:2, AG:2, BD:2, BE:2

3-item set

CDF:2

4-item set

BDHJ:2

b. Tìm tất cả luật kết hợp thỏa minsup và minconf Chạy file **Apriori.py** để tìm các luật kết hợp thỏa minconf

```

1612272_02 [C:\Users\hat huy\Desktop\1612272_02] - _\Source\LY THUYET\Apriori.py [1612272_02] - PyCharm
File Edit View Navigate Code Refactor Run Tools VCS Window Help
1612272_02 LY THUYET Apriori.py
Run: Apriori
Luật khai thác bằng tập phổ biến là:
G => A with conf 100.0 %
E => B with conf 100.0 %
H => B with conf 100.0 %
J => B with conf 100.0 %
C => D with conf 100.0 %
F => C with conf 100.0 %
C => F with conf 100.0 %
F => D with conf 100.0 %
H => D with conf 100.0 %
J => D with conf 100.0 %
J => H with conf 100.0 %
H => J with conf 100.0 %
DH => B with conf 100.0 %
BH => D with conf 100.0 %
BD => H with conf 100.0 %
H => BD with conf 100.0 %
DJ => B with conf 100.0 %
BJ => D with conf 100.0 %
BD => J with conf 100.0 %
J => BD with conf 100.0 %
HJ => B with conf 100.0 %
BJ => H with conf 100.0 %
BH => J with conf 100.0 %
J => BH with conf 100.0 %
H => BJ with conf 100.0 %
DF => C with conf 100.0 %
CF => D with conf 100.0 %

```

Figure 4: Luật khai thác bằng tập phổ biến

```

DF => C with conf 100.0 %
CF => D with conf 100.0 %
CD => F with conf 100.0 %
F => CD with conf 100.0 %
C => DF with conf 100.0 %
HJ => D with conf 100.0 %
DJ => H with conf 100.0 %
DH => J with conf 100.0 %
J => DH with conf 100.0 %
H => DJ with conf 100.0 %
Time for frequent items: 0.0
Tap pho bien dong la:

{'A': 3, 'B': 3, 'C': 2, 'D': 3, 'E': 2, 'F': 2, 'G': 2, 'H': 2, 'I': 2}
{'AB': 2, 'AD': 2, 'AG': 2, 'BE': 2}
{'CDF': 2}
{'BDHJ': 2}

Luật khai thác bang tap pho bien dong:
G => A with conf 100.0 %
E => B with conf 100.0 %
DF => C with conf 100.0 %
CF => D with conf 100.0 %
CD => F with conf 100.0 %
F => CD with conf 100.0 %
C => DF with conf 100.0 %
Time for closed frequent items: 0.0
Tap pho bien toi dai la:

```

Figure 5: Luật khai thác bằng tập phổ biến (tt)

```

{'BDHJ': 2}

Luật khai thác bang tap pho bien dong:
G => A with conf 100.0 %
E => B with conf 100.0 %
DF => C with conf 100.0 %
CF => D with conf 100.0 %
CD => F with conf 100.0 %
F => CD with conf 100.0 %
C => DF with conf 100.0 %
Time for closed frequent items: 0.0
Tap pho bien toi dai la:

{'A': 3, 'B': 3, 'C': 2, 'D': 3, 'E': 2, 'F': 2, 'G': 2, 'H': 2, 'I': 2}
{'AB': 2, 'AD': 2, 'AG': 2, 'BE': 2}
{'CDF': 2}
{'BDHJ': 2}

Luật khai thác bang tap pho bien toi dai:
G => A with conf 100.0 %
E => B with conf 100.0 %
DF => C with conf 100.0 %
CF => D with conf 100.0 %
CD => F with conf 100.0 %
F => CD with conf 100.0 %
C => DF with conf 100.0 %
Time for maxima frequent items: 0.0

Process finished with exit code 0

```

Figure 6: Luật khai thác từ tập phổ biến đóng và tối đại

Ta thấy tập luật khai thác từ tập phổ biến nhiều nhưng có những luật là con của các luật khác.

Các luật thỏa yêu cầu bài toán là:

- DF => C with conf 100.0 %
- CF => D with conf 100.0 %
- CD => F with conf 100.0 %

c. Ứng dụng ở câu a. vào khai thác luật kết hợp. So sánh về hiệu quả

Việc thực hiện ứng dụng khai thác luật kết hợp với tập phổ biến đóng và tập phổ biến tối đại đã thực hiện ở câu trên. Vì tập phổ biến đóng và tập phổ biến tối đại là như nhau nên luật sinh ra là như nhau. Các luật sinh ra ở tập phổ biến đóng sẽ gọn hơn tập luật sinh ra ở tập phổ biến. Về thời gian thực hiện thì ta chú ý dòng “Time for ...” là như nhau đối với dataset này vì tập phổ biến ít nên chúng ta không thấy sự khác biệt. Tuy nhiên, đối với dataset lớn thì thời gian giữa các phương pháp sẽ có cách biệt đáng kể

B. THỰC HÀNH

1. CÂU 1

Ta tải các file trong [link](#) được cung cấp gồm có: **plants.data**, **plants.names**, **stateabbr.txt**.

Để chuyển đổi dữ liệu ta cần xem nội dung các file

```

1 State Abbreviations
2
3 U.S. States:
4 ab Alabama
5 ak Alaska
6 ar Arkansas
7 az Arizona
8 ca California
9 co Colorado
10 ct Connecticut
11 de Delaware
12 dc District of Columbia
13 fl Florida
14 ga Georgia
15 hi Hawaii
16 id Idaho
17 il Illinois
18 in Indiana
19 ia Iowa
20 ks Kansas
21 ky Kentucky
22 la Louisiana
23 me Maine
24 md Maryland
25 ma Massachusetts
26 mi Michigan
27 mn Minnesota
28 ms Mississippi
29 mo Missouri
30 mt Montana
31 ne Nebraska
32 nv Nevada
33 nh New Hampshire
34 nj New Jersey
35 nm New Mexico
36 ny New York
37 nc North Carolina
38 nd North Dakota

```

Figure 7: Nội dung file stateabbr.txt

Ta tiến hành loại bỏ các dòng không cần thiết để phục vụ cho các xử lý tiếp theo. File sau khi xử lý sẽ được lưu thành **stateabbr_process.txt**. Trong file **stateabbr.txt** có **state Prince Edward Island** không có ký hiệu viết tắt tuy nhiên trong file **plants.data** lại sử dụng **pe** thay cho **Prince Edward Island**

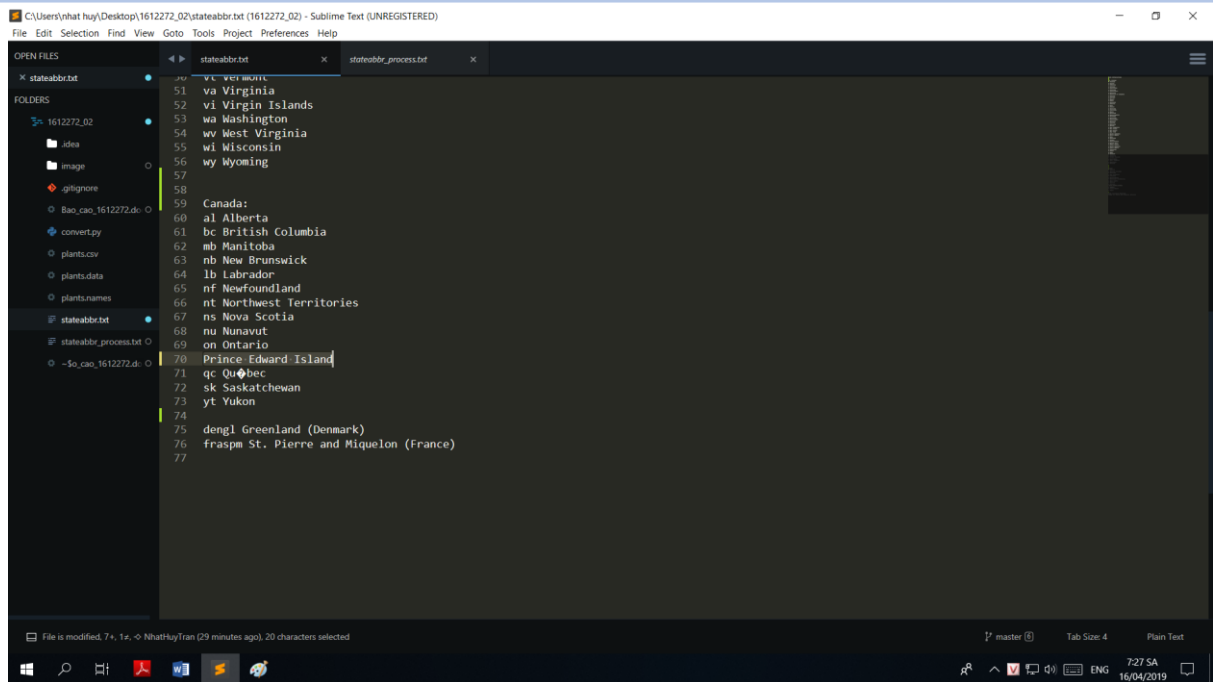


Figure 8: Thành phố Prince Edward Island trong stateabbr.txt

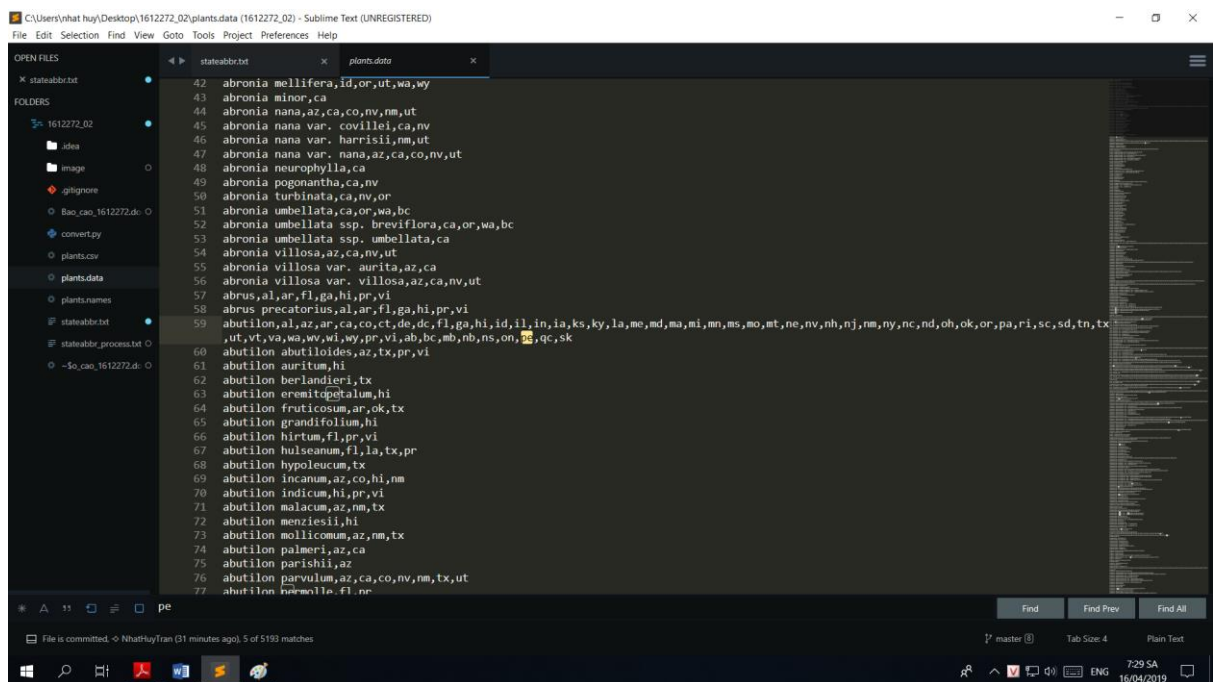


Figure 9: Ký hiệu pe được sử dụng trong plants.data

Ta sẽ thêm ký hiệu **pe** cho **Prince Edward Island** và xóa hết các dòng không cần thiết. Lưu lại thành file **stateabbr_process.txt**

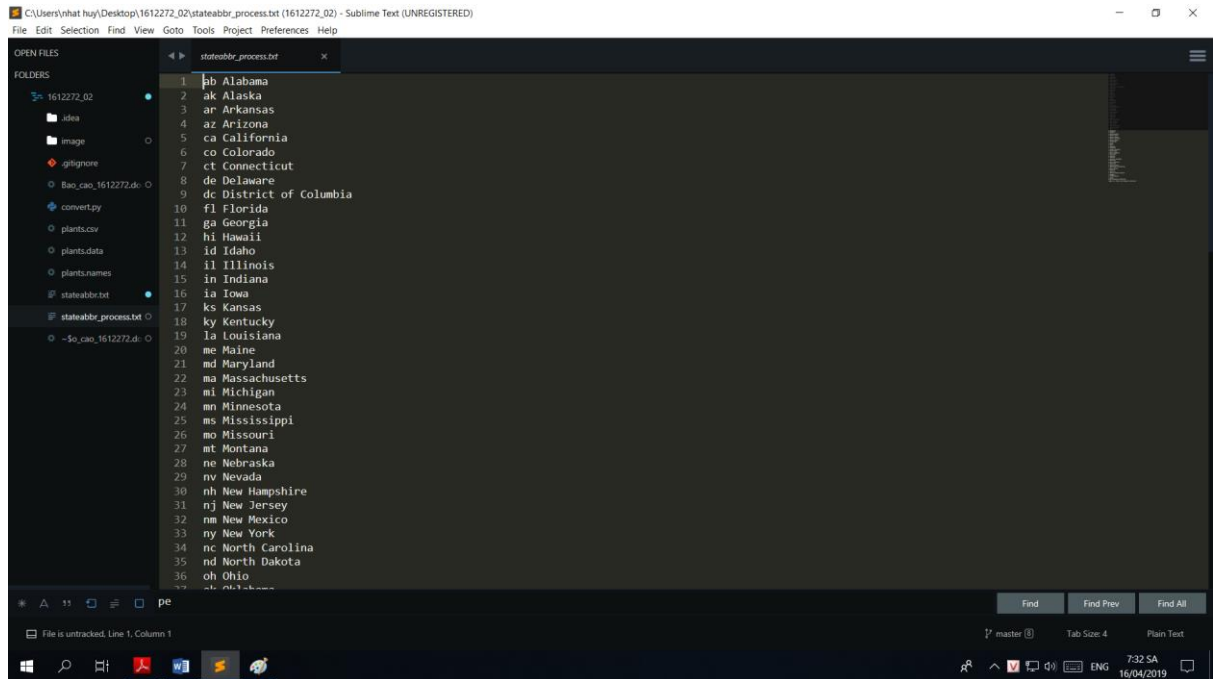


Figure 10: File stateabbr_process.txt

Để chuyển file plants.data thành file plant.csv ta chạy code file convert.py

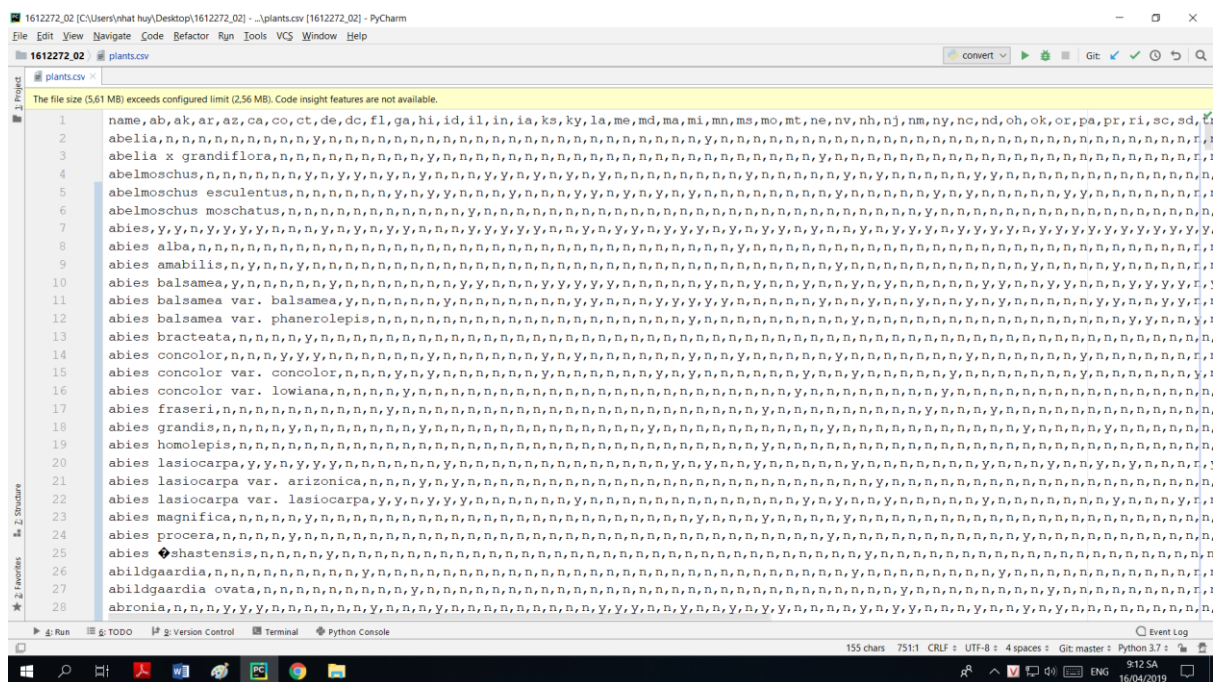


Figure 11: File plants.data sau khi chuyển

2. CÂU 2

Ta tiến hành mở file plants.csv bằng Weka

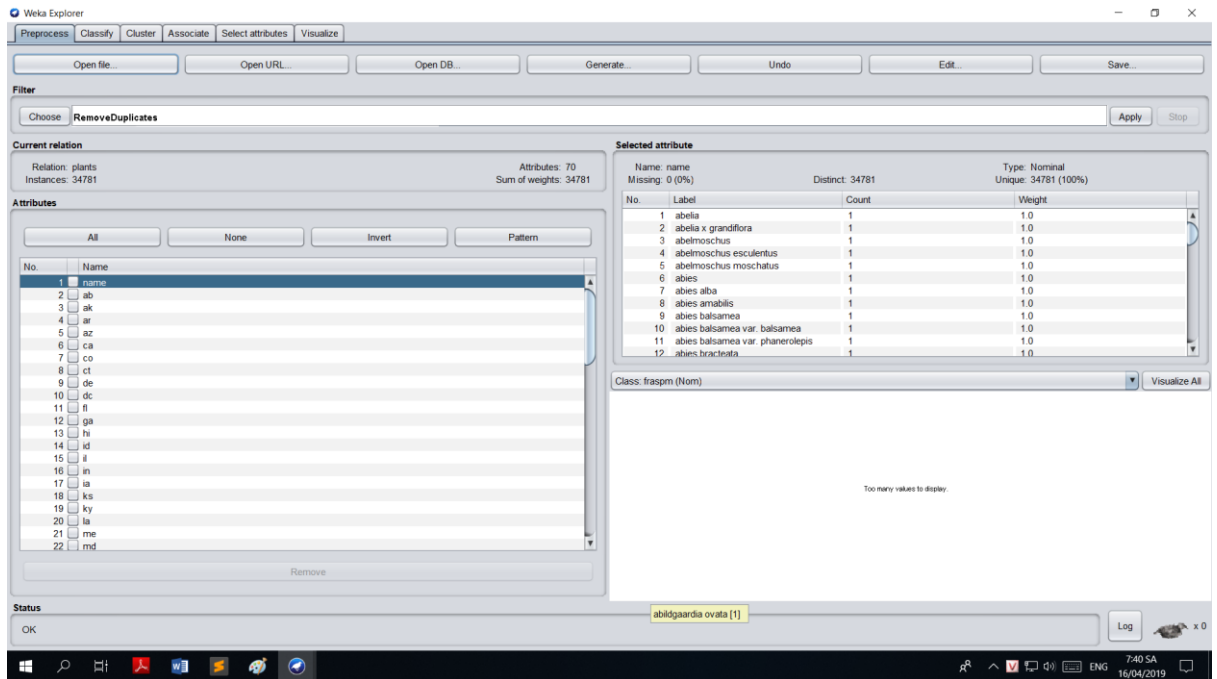


Figure 12: File plants.csv mở bằng Weka

Ta có tổng cộng 34781 cây khác nhau

Dựa vào hình ta có tổng cộng 70 thuộc tính trong đó có 1 thuộc tính name và 69 thuộc tính vùng phân bố. Vậy chúng ta có tổng cộng 69 vùng phân bố

Để xác định mỗi vùng có bao nhiêu loại cây, ta ấn vào nút **Visualize all** để xem



Figure 13: 31 vùng đầu tiên



Figure 14: 32 vùng tiếp theo

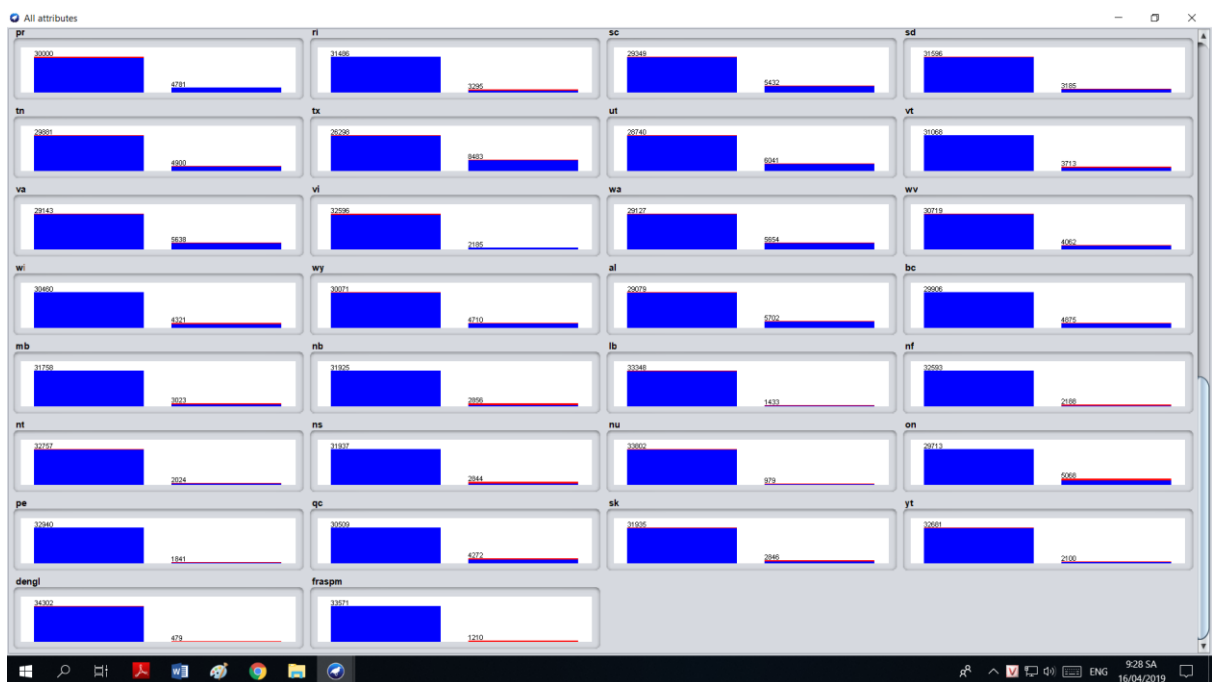


Figure 15: 6 vùng còn lại

Để kiểm tra xem vùng phân bố có ít/nhiều loài cây nhất, số lượng, tỉ lệ, trung bình một vùng phân bố bao nhiêu cây thì ta chạy file **min_max.py**

The screenshot displays the PyCharm IDE interface. The top toolbar shows the 'Run' button (a green play icon). Below the toolbar, the 'Run' tab is active, showing the execution output of the script `min_max.py`. The output text is as follows:

```
Dem xong!  
{ 'ab': 3408, 'ak': 2969, 'ar': 4610, 'az': 6778, 'ca': 11676, 'co': 5465, 'de': 3630, 'dc': 3080, 'fl': 6621, 'ga': 5942, 'hi':  
min_state: dengl with 479 instance and proportion is 1.3771886949771428  
max_state: ca with 11676 instances and proportion is 33.570052614933445  
average: 4370.333333333333  
  
Process finished with exit code 0
```

The bottom status bar of the IDE indicates the current file encoding is '46:94 CRLF', the text is in 'UTF-8' with '4 spaces' of indentation, the repository is 'Git: master', and the Python version is 'Python 3.7.7'. The system clock in the bottom right corner shows '11:49 SA 16/04/2019'.

Hình 1: Kết quả khi chạy file min_max.py

Vùng ít cây nhất là **Greenland Denmark (dengl)** với 479 loại cây và tỉ lệ là 1.37%

Vùng có nhiều cây nhất là **California (ca)** với 11676 loại cây và tỉ lệ là 33.57%

Trung bình mỗi vùng có khoảng 4370 loại

3. CÂU 3

Ta thay thế các giá trị “n” thành “?” bằng cách chạy file **change_value.py**. File sau khi đã thay thế được đặt tên là **plants_changed.csv**

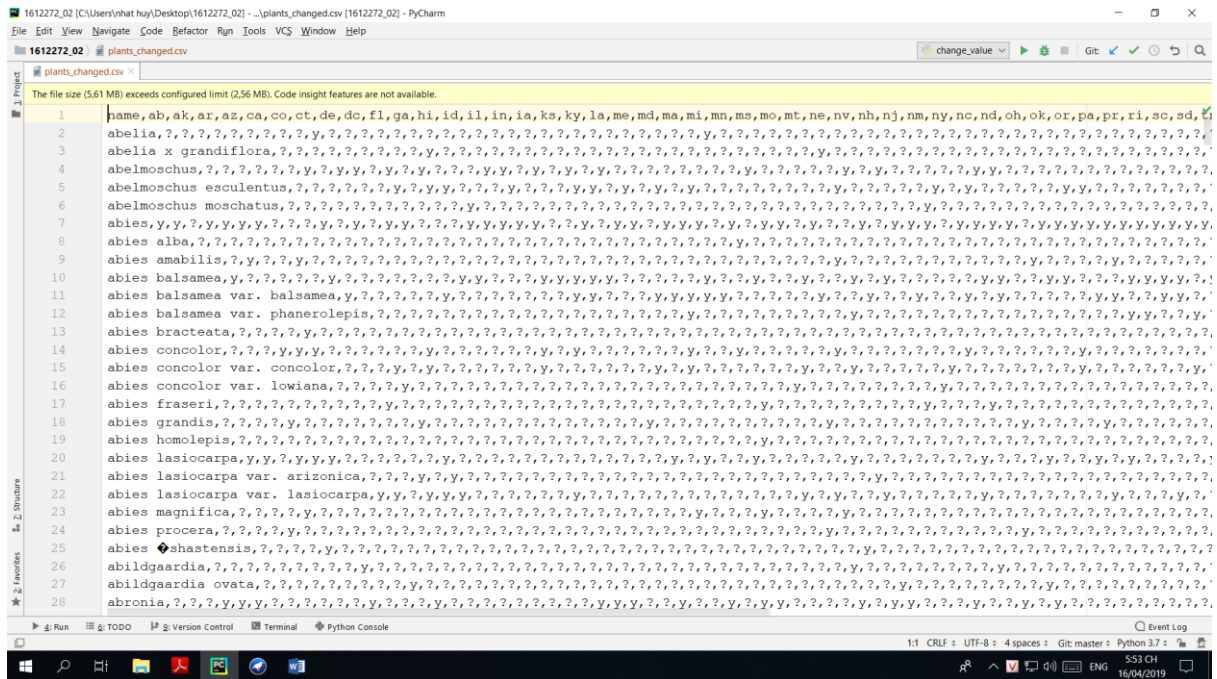


Figure 16: Nội dung file sau khi thay thế

Mở file plants_changed.csv bằng Weka và tiến hành các bước tiếp theo

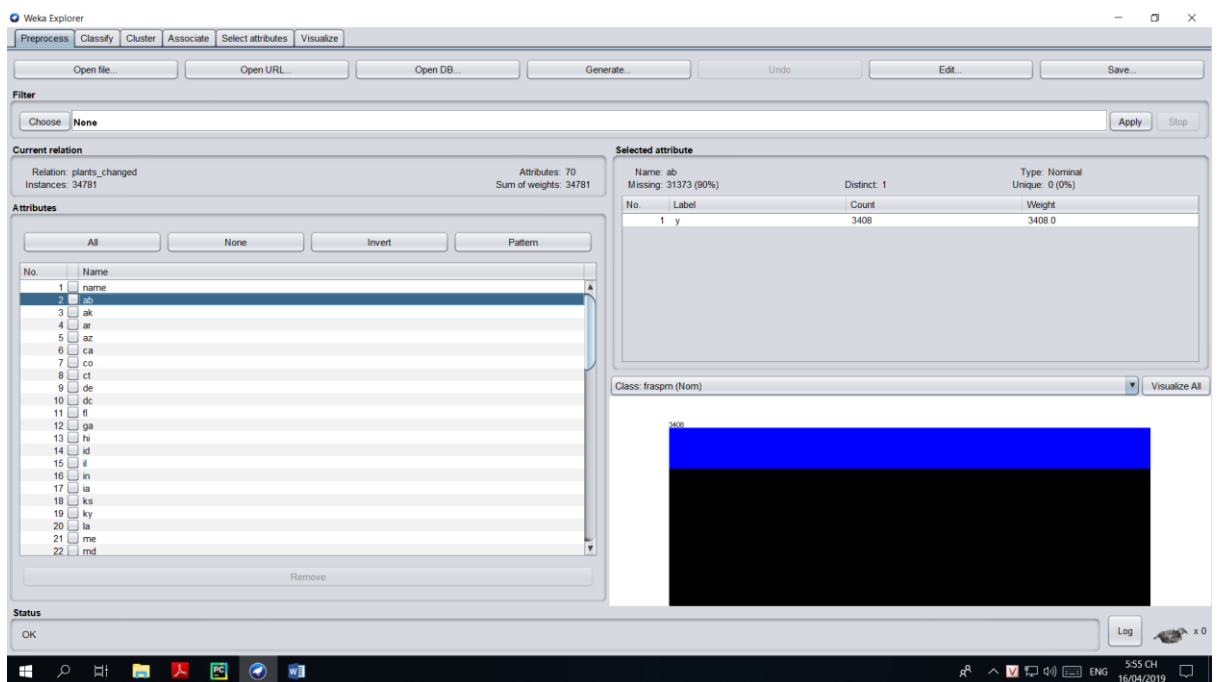


Figure 17: File plants_changed.csv được mở bằng Weka

Xóa thuộc tính name bằng cách click vào name và chọn Remove

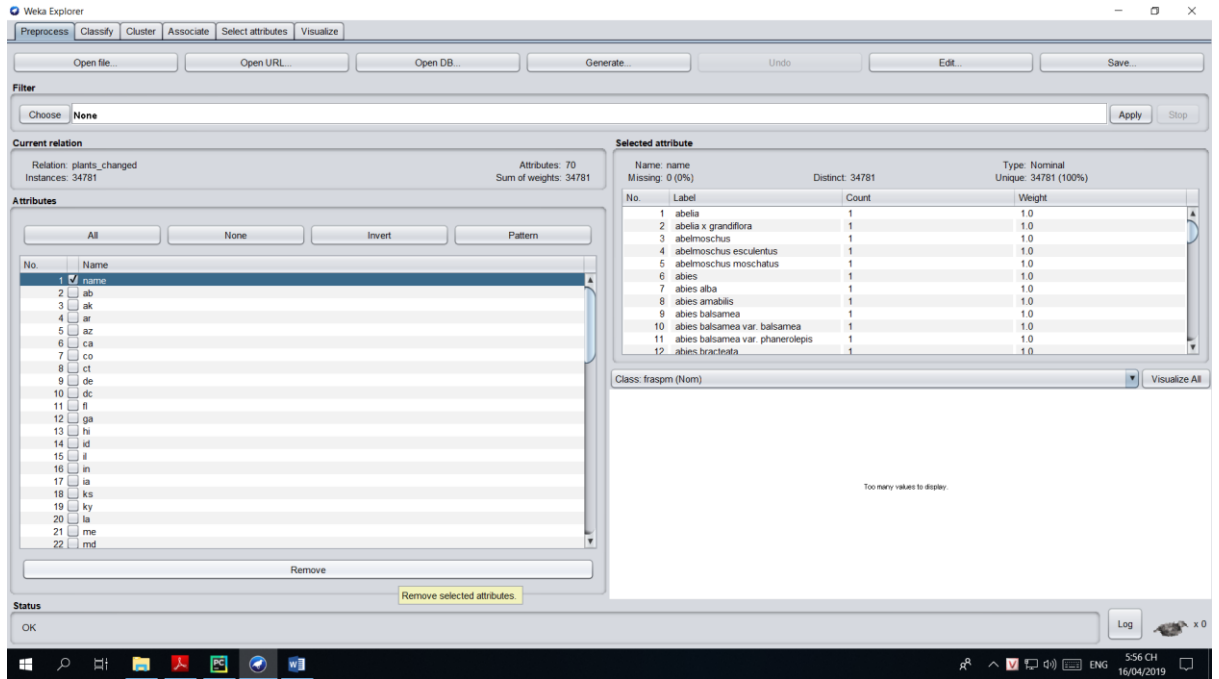


Figure 18: Chọn Remove để xóa thuộc tính name

Lưu file lại với tên là **plants.arff**

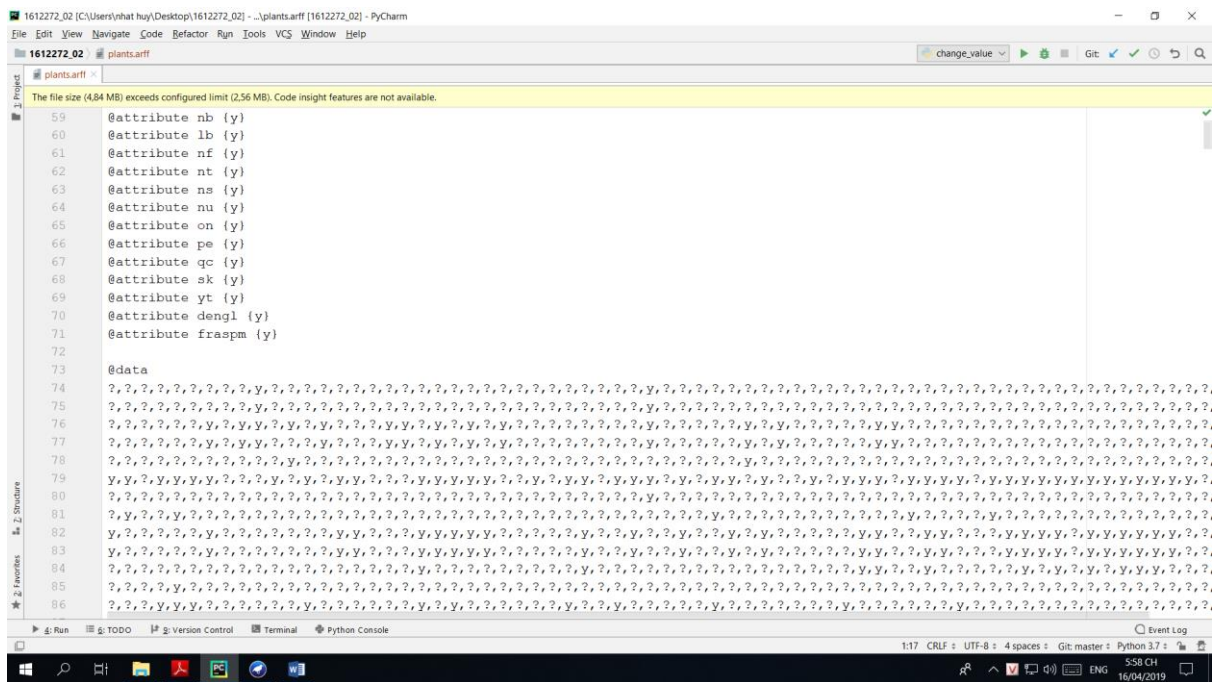


Figure 19: Nội dung file plants.arff

4. CÂU 4

Để khai thác tập phổ biến ta chuyển sang tab **Associate**

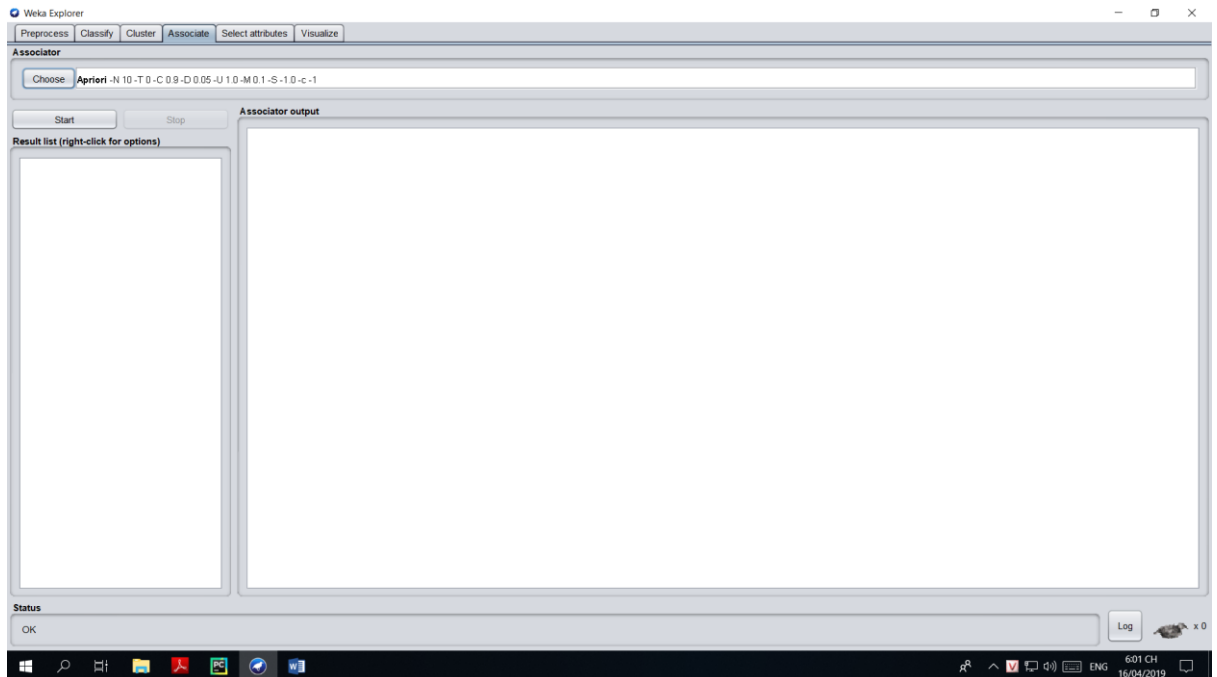


Figure 20: Tab Associate

Chọn thuật toán **Apriori** với các tham số sau

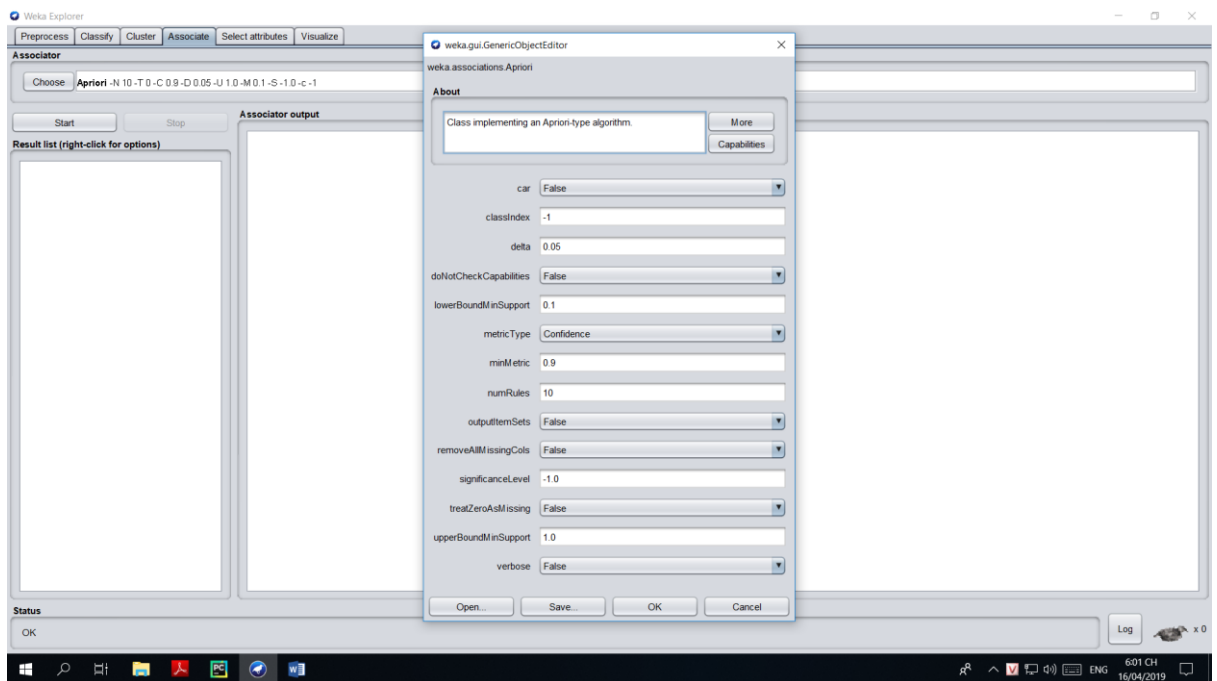


Figure 21: Thuật toán Apriori

Kết quả

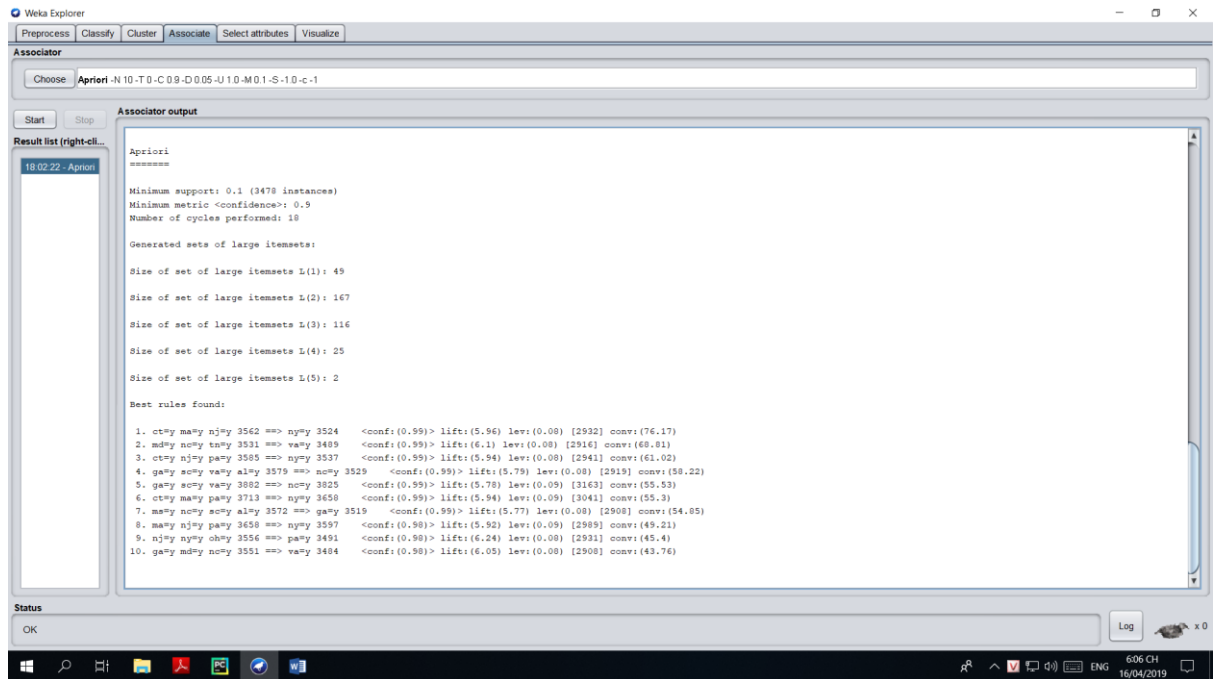


Figure 22: Kết quả chạy Apriori

KÍCH THƯỚC

SỐ LƯỢNG

1 hạng mục

49

2 hạng mục

167

3 hạng mục

116

4 hạng mục

25

5 hạng mục

2

5. CÂU 5

Chọn thuật toán **FP-Growth** với tham số **metricType** là **Confidence** và **minMetric** là **0.95**

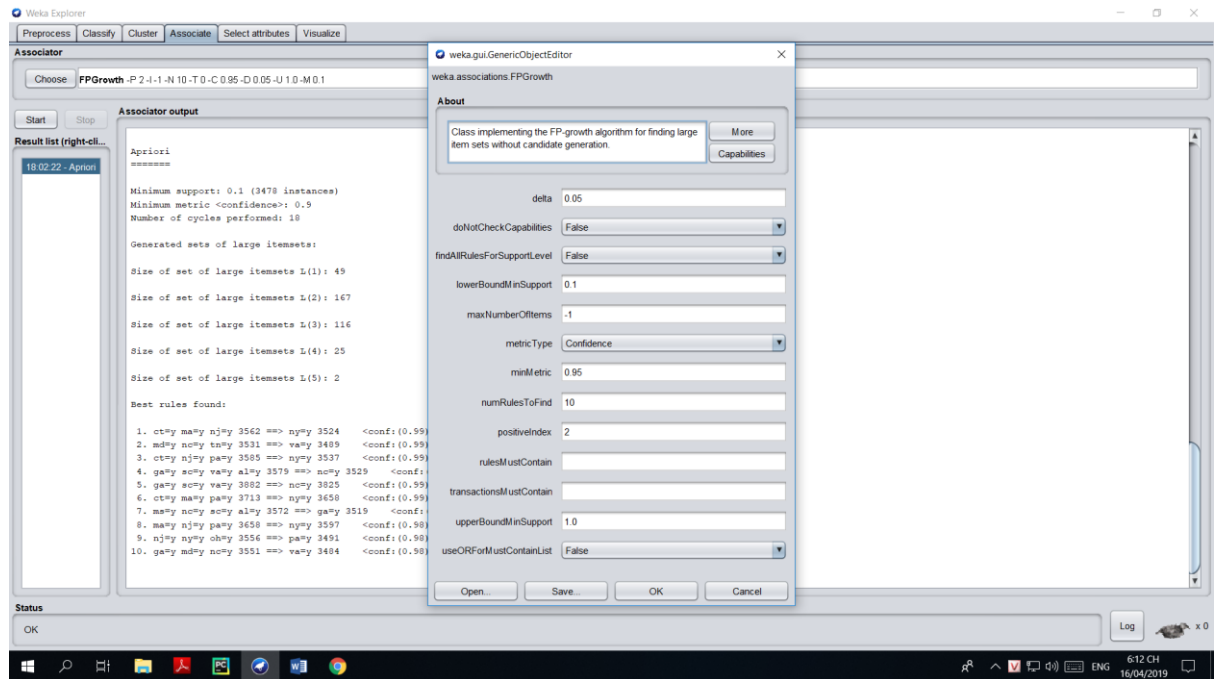


Figure 23: Tham số thuật toán FP-Growth

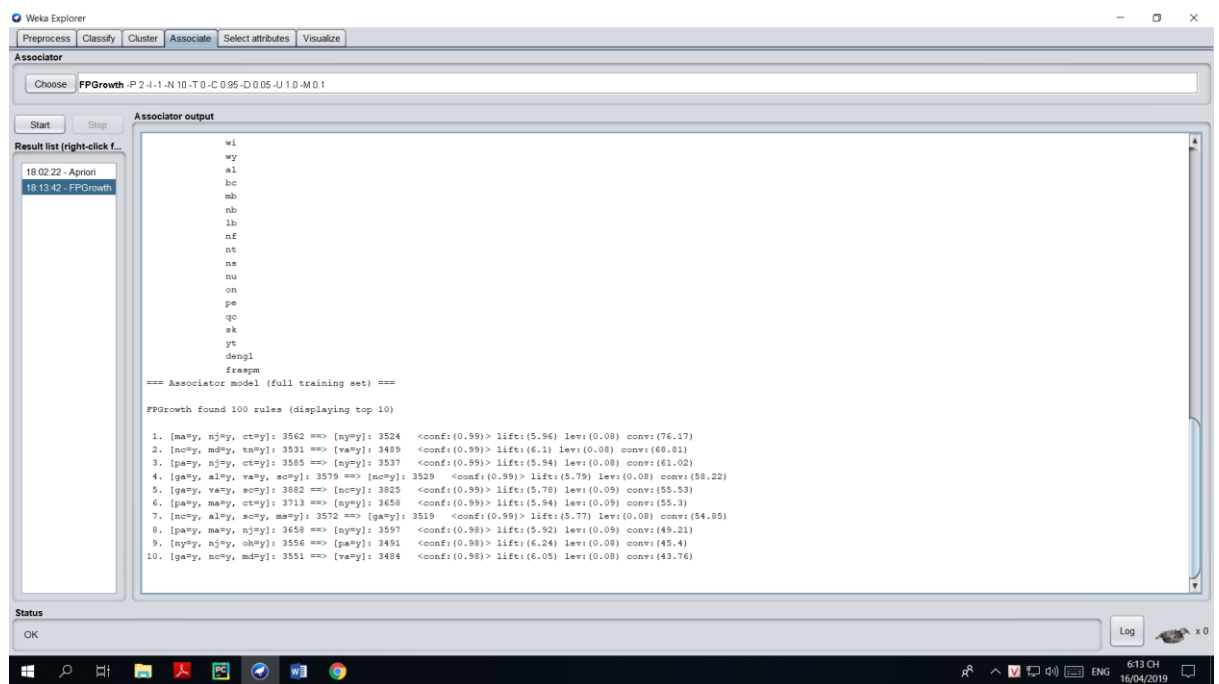


Figure 24: Kết quả chạy FP-Growth

TẬP HẠNG MỤC PHỔ BIẾN	SỐ LƯỢNG LUẬT
ma=y, nj=y, ct=y	3562
nc=y, md=y, tn=y	3531
pa=y, nj=y, ct=y	3585
ga=y, al=y, va=y, sc=y	3579
ga=y, va=y, sc=y	3882
pa=y, ma=y, ct=y	3713
nc=y, al=y, sc=y, ms=y	3572
pa=y, ma=y, nj=y	3658
ny=y, nj=y, oh=y	3556
ga=y, nc=y, md=y	3551