

# DATA SCIENCE MEI/1

**University of Beira Interior** Department of Informatics

Hugo Pedro Proença, <u>hugomcp@di.ubi.pt</u>, **2021/22** 

### Frequent Pattern Mining

#### **Frequent patterns** are

collections of items that appear in a dataset at a high frequency

- Such patterns reveal association rules and relations between variables.
- Frequent Pattern Mining is an important research area in data science applied to different domains:
  - □ Recommender systems;
  - Bioinformatics;

Decision making;



Source: http://www.dataception.in/blog/what-is-frequent-pattern-mining-association-and-how-does-it-support-business-analysis/

#### Frequent Pattern Mining

- □ Frequent patterns (a.k.a. Association Rule Mining) consists of first finding frequent item sets satisfying a minimum support threshold, or percentage from which strong association rules in the form A  $\rightarrow$  B of are generated.
- These rules also satisfy a minimum confidence threshold (a prespecified probability of satisfying B under the condition that A is satisfied).
- □Associations can be further analyzed to uncover correlation rules, which convey statistical correlations between item sets A and B.
- □Various algorithms have been developed for frequent itemset mining, from which association and correlation rules are further derived. These algorithms can be classified into three categories:
  - □Apriori-like algorithms;
  - □Frequent pattern growth –based algorithms such as FP-growth;
  - Algorithms that use the vertical data format.

□The **Apriori** algorithm is the most classical algorithm for mining frequent item sets for Boolean (A -> B) association rules.

□ It was proposed by Agrawal and Srikant in 1994.

- □ It explores the level-wise mining property that all nonempty subsets of a frequent itemset must also be frequent.
- □At the k<sup>th</sup> iteration (k >= 2), it forms frequent candidate sets of k elements based on the frequent (k 1)-item sets, and scans the database once to find the complete set of frequent k-item sets,  $L_k$ .
- Due to its popularity, several variations have been proposed, including hashing and transaction reduction, in order to make the procedure more efficient. Other variations include partitioning the data (mining on each partition and then combining the results) and sampling the data (mining on a data subset). These variations can reduce the number of data scans required to as little as two or even one.

Transaction	ltems
А	1, 3, 4
В	2, 3, 5
С	1, 2, 3, 5
D	2, 5

**Step 1:** Count frequency of 1-item sets

Set	Frequency
1	2
2	3
3	3
4	1
5	3

Step 2: Remove the least frequent set

Set	Frequency
1	3
2	3
3	3
4	1
5	4

Transaction	ltems
А	1, 3, 4
В	2, 3, 5
С	1, 2, 3, 5
D	2, 5

**Step 3:** Count frequency of 2-item sets (of remaining elements)

Set	Frequency
1, 2	1
1, 3	2
1, 5	1
2, 3	2
2, 5	3
3, 5	2

Set	Frequency
1, 2	1
1, 3	2
1, 5	1
2, 3	2
2, 5	3
3, 5	2

Step 4: Remove the least frequent set

A1, 3, 4Step 5: CountB2, 3, 5frequency of 3-itemC1, 2, 3, 5sets (of remaining elements)	Transaction	Items
B2, 3, 5frequency of 3-item sets (of remaining elements)D2, 5	А	1, 3, 4
C1, 2, 3, 5sets (of remaining elements)D2, 5	В	2, 3, 5
D 2, 5	С	1, 2, 3, 5
	D	2, 5

#### **Pseudocode:**

```
 \begin{array}{l} \mbox{Apriori}(T, \epsilon) \\ L_1 \leftarrow \{ large 1 - itemsets \} \\ k \leftarrow 2 \\ \mbox{while } L_{k-1} \mbox{ is not empty} \\ C_k \leftarrow \{ c = Union(a, \{b\}) : a \mbox{ in } L_{k-1} \mbox{ and } b \mbox{ not in } a, \{ s \subseteq c : |s| = k - 1 \} \subseteq L_{k-1} \} \\ \mbox{ for transactions t in } T \\ D_t \leftarrow \{ c \mbox{ in } C_k : c \subseteq t \} \\ \mbox{ for candidates } c \mbox{ in } D_t \\ \mbox{ count}[c] \leftarrow count[c] + 1 \\ \\ L_k \leftarrow \{ c \mbox{ in } C_k : count[c] \ge \epsilon \} \\ k \leftarrow k + 1 \\ \end{array} \right.
```

1, 2, 5

0

□ Frequent Pattern Growth is a method of mining frequent itemsets without candidate generation.

□This algorithm constructs a highly compact data structure (an FP-tree) to compress the original transaction database.

- An FP-tree is a tree-like structure made of the initial item-sets of the database. The purpose of the FP tree is to mine the most frequent pattern. Each node of the FP tree represents an item of the itemset.
- The root node represents *null* while the lower nodes represent the itemsets. The paths between the root and the leaves represent the item-sets

Instead employing the generate-and-test strategy of Apriori-like methods, it focuses on frequent pattern (fragment) growth, which avoids costly candidate generation, and enables greater efficiency.

Transaction	Items
А	3, 4
В	2, 3, 5
С	1, 2, 3, 5
D	2, 5
E	5,4
F	5, 3

**Step 1:** Count frequency of 1-item sets

Set	Frequency	
1	1	
2	3	
3	4	
4	2	
5	5	

Step 2: Sort 1-item sets in descending order

Set	Frequency
5	5
3	4
2	3
4	2
1	1

		ordered set $\{4, 3\} \rightarrow$
Transaction	Items	{3, 4}
А	4, 3	Freq(3)=4, Freq(4)=2
В	2, 3, 5	
С	1, 2, 3, 5	
D	2, 5	
E	5, 4	
F	5, 3	Step A. Undete
		Support

3: S1

4: S1

As there are no nodes in the tree, new elements should be created and inserted

NULL



Transaction	Items
А	4, 3
В	2, 3, 5
С	1, 2, 3, 5
D	2, 5
E	5,4
F	5, 3

Step 6: Update Support

Α

В

С

D

Ε

F











Items	Conditional Pattern Base
1	{5, 3, 2: S1}
2	{5: S5}, {5, 3: S2}
3	{5: S3}
4	{3: S1}, {5: S1}
5	

Step 11: Build the Conditional Frequent Pattern Tree, that identifies the common elements in each row of the Conditional Pattern Base and sums their supports

ltems	Conditional Pattern Base	Conditional Frequent Pattern Tree
1	{5, 3, 2: S1}	{5, 2, 3: S1}
2	{5: S1}, {5, 3: S2}	{5: S3}, {3: S2}
3	{5: S3}	{5: S3}
4	{3: S1}, {5: S1}	{3: S1}, {5: S1}
5		

Items	Conditional Pattern Base	Conditional Frequent Pattern Tree
1	{5, 3, 2: S1}	{5, 2, 3: S1}
2	{5: S1}, {5, 3: S2}	{5: S3}, {3: S2}
3	{5: S3}	{5: S3}
4	{3: S1}, {5: S1}	{3: S1}, {5: S1}
5		

Items	Frequent Patterns	
1	{5, 2, 3, <b>1</b> : S1}	
2	{5, <b>2</b> : S3}, {3, <b>2</b> : S2}	
3	{5, <b>3</b> : \$3}	
4	{3, <b>4</b> : S1}, {5, <b>4</b> : S1}	
5		

Step 12: From the Conditional Frequent Pattern Trees (CFPT), the Frequent Patterns are detected by pairing each CFPT element, to the item

#### Frequent Pattern Mining Exercise

- Consider the following set of transactions of a supermarket.
- Find, according to the APRIORI and the FP-GROWTH algorithms, the Frequent Patterns in such dataset, considering:

Support	= 20%
---------	-------

□Support = 50%

Transaction	ltems
А	1, 7, 4, 6,
В	2, 3, 5
С	1, 2, 3, 5
D	4, 6, 7
E	8, 7, 2, 3
F	2, 3, 5, 2, 1
G	4, 5, 6, 7, 9
Н	9, 1, 3, 2, 4
I	4, 6, 7, 5
J	7, 9, 5, 6, 1