

**Data Warehousing and Data Mining**  
**Assignment 2**  
**ASSOCIATION ANALYSIS**

**Alicja Jonczyk, Kacper Multan**

**1. Apriori Algorithm**

- a) Show thoroughly the steps on how the frequent itemsets are generated.

We started with counting the number of appearances for each of the 1-item itemsets: A, B, C, D, I, U. Then we continued with calculating the support using the formula:

$$s(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{N}$$

With red we have marked all those itemsets which support was below the minimum threshold. Then, using the itemsets that have passed (marked in green) we generated the 2-item itemsets. We repeated counting and calculating the support. Again, we pruned the not frequent itemsets and repeated the steps until we reached the state where all of the supersets were below the threshold.

Item	Count	Support	Item	Count	Support	Item	Count	Support
A	4	0,666667	AB	2	0,333333	ABD	0	0
B	2	0,333333	AC	1	0,166667	ABI	0	0
C	2	0,333333	AD	2	0,333333	ADI	2	0,333333
D	3	0,5	AI	2	0,333333			
I	4	0,666667	BC	1	0,166667			
U	1	0,166667	BD	0	0			
			BI	0	0			
			CD	0	0			
			CI	1	0,166667			
			DI	3	0,5			

- b) Find all association rules based on this set, given confidence threshold  $c = 60\%$ .

1. $\{A, D\} \Rightarrow \{I\}$	$c = \frac{2}{2} = 1 \quad +$
2. $\{A, I\} \Rightarrow \{D\}$	$c = \frac{2}{2} = 1 \quad +$
3. $\{D, I\} \Rightarrow \{A\}$	$c = \frac{2}{3} \quad +$
4. $\{A\} \Rightarrow \{D, I\}$	$c = \frac{2}{6} = \frac{1}{3} \quad -$ $1 - \frac{1}{3} = \frac{2}{3}$
5. $\{D\} \Rightarrow \{A, I\}$	$c = \frac{2}{3} \quad +$
6. $\{I\} \Rightarrow \{A, D\}$	$c = \frac{2}{4} = \frac{1}{2} \quad -$
↓	
Correct rules that satisfy confidence threshold:	
7. $\{A, D\} \Rightarrow \{I\}$	
8. $\{A, I\} \Rightarrow \{D\}$	
9. $\{D, I\} \Rightarrow \{A\}$	
10. $\{D\} \Rightarrow \{A, I\}$	

## 2. ECLAT Algorithm

2.1

Implement one of three variations of the frequent itemset mining algorithms: Apriori algorithm, FP-growth, and Eclat. Which variation will you implement and why?

Given the scale of data the most suitable algorithm to implement would be the FP-growth algorithm. The FP-growth algorithm is more efficient than ECLAT and the Apriori algorithm, especially when dealing with large datasets. It doesn't require multiple passes over the data to generate candidate itemsets, instead it builds a compact data structure - the FP-tree to represent the transactions and efficiently mines frequent itemsets from it. Additionally, the FP-growth algorithm tends to be faster than both Apriori and Eclat, especially as the dataset size increases.

2.2

Use ECLAT algorithm to compute the frequent itemsets of the Table2 with minimum support count of 3.

The frequent itemsets of Table 2 are: A, C, D, E, F, G, AE, AG, CD, CF, CG, DF, DG, FG.

### 3 FP-Growth Algorithm

#### 3.1 Assignment

Note: We have overseen the note about the support count ties and put the item G before item B, which has resulted in slightly different trees. But according to our analysis, it didn't change the exercise's logic, because the found frequent itemsets are the same ones we would get if we would change this order. (only difference happens in Rg and Rb trees where the mined frequent itemsets from Rb tree swapped places with the frequent itemsets from Rg tree).

**Ordering items and transactions by descending support count:**

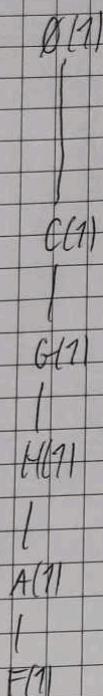
Item	Support
C	10
G	8
B	8
H	7
A	6
E	5
D	4
F	3

$A = 6$   
 $B = 8$   
 $C = 10$   
 $D = 4$   
 $E = 5$   
 $F = 3$   
 $G = 8$   
 $H = 7$

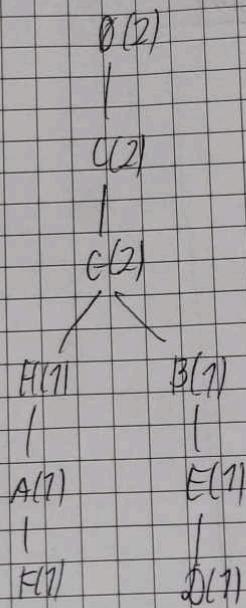
Transactions	Items
T10	CGHAF
T11	CGBED
T12	CBGEF
T13	CGBA
T14	CHED
T15	GBHA
T16	GBHAD
T17	GBE
T18	GBHAF
T19	GBHAED

Creating main FP-tree:

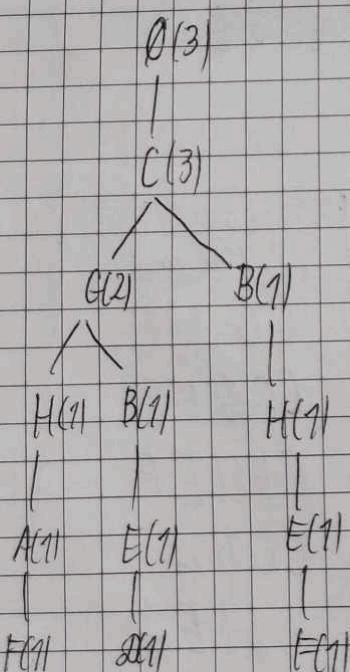
1. CGHAF



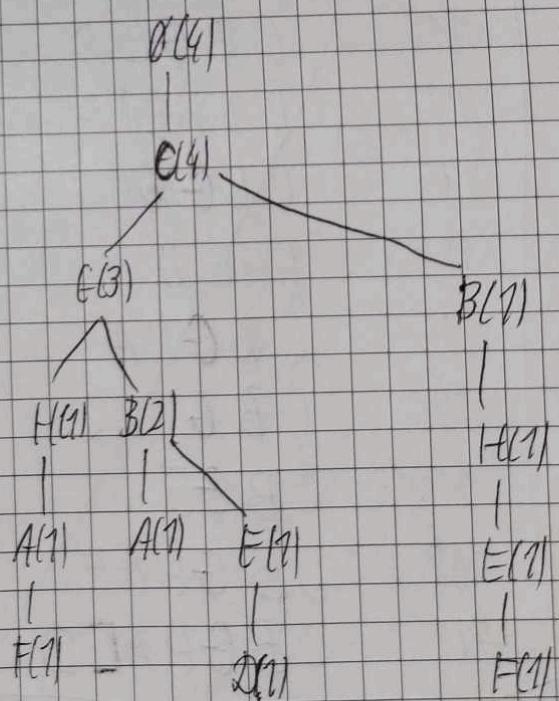
2. CGBED



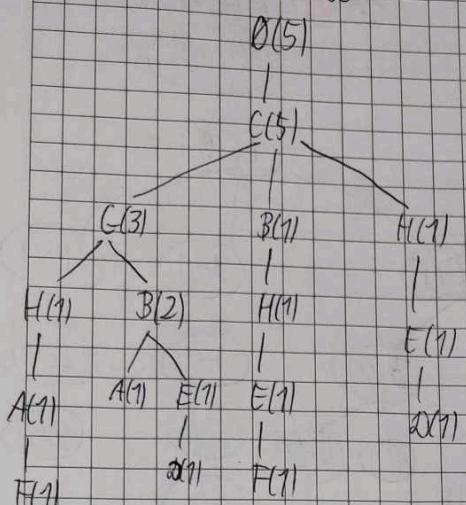
3. CBHEF



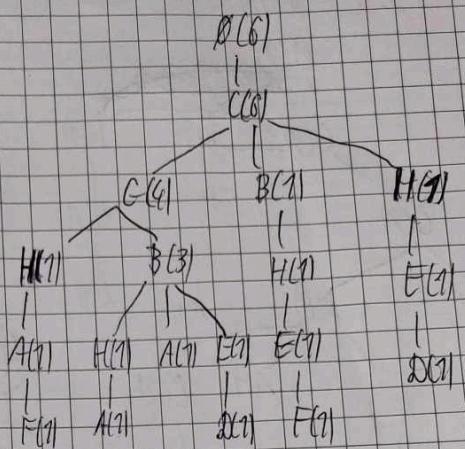
4. CGBA



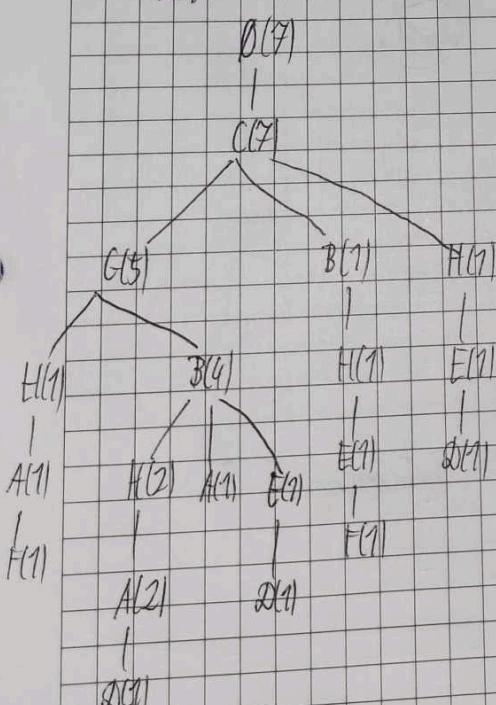
5. ~~CBHD~~ CHED



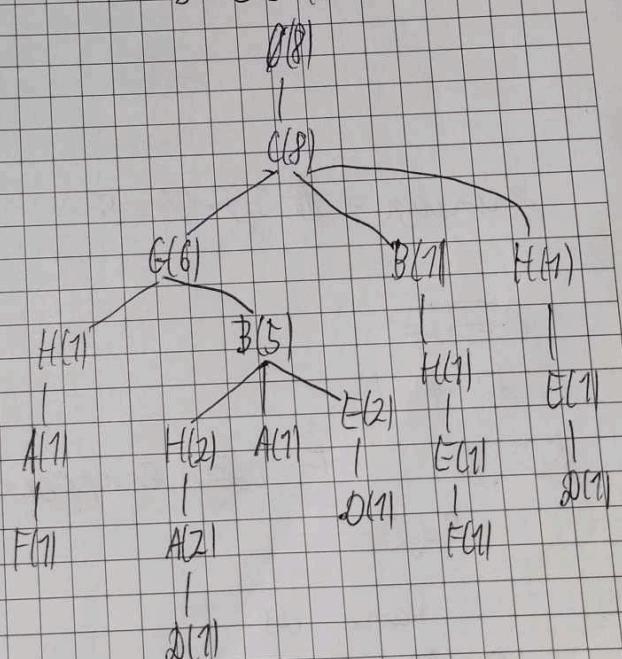
6. ~~CBHA~~ CGBAH

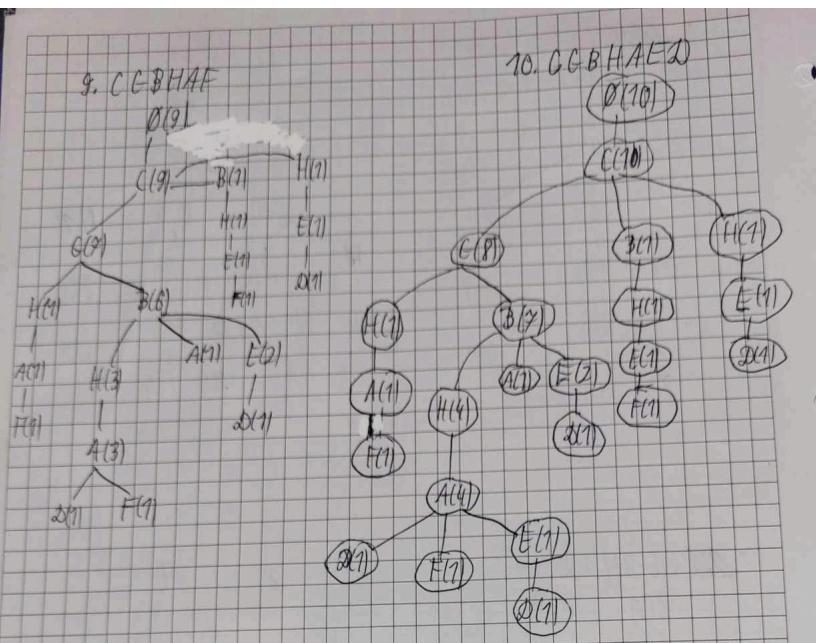


7. ~~CG-BHA~~ CG-BHA



8. CG-BE





$$\text{missup} = 0,5 \cdot 10 = 5$$

~~REH3~~

P=F:

• P F appears only 3 times < 5

P=D:

Same as F

P=E:

{EC(5)}

P=A:

{ACC(6), AG(6), AB(5), ABC(5), ACC(6), ABGC(5), AH(5)}

P=H:

{HC(5), HG(5), HB(5), HB, CC(5), HGC(5)}

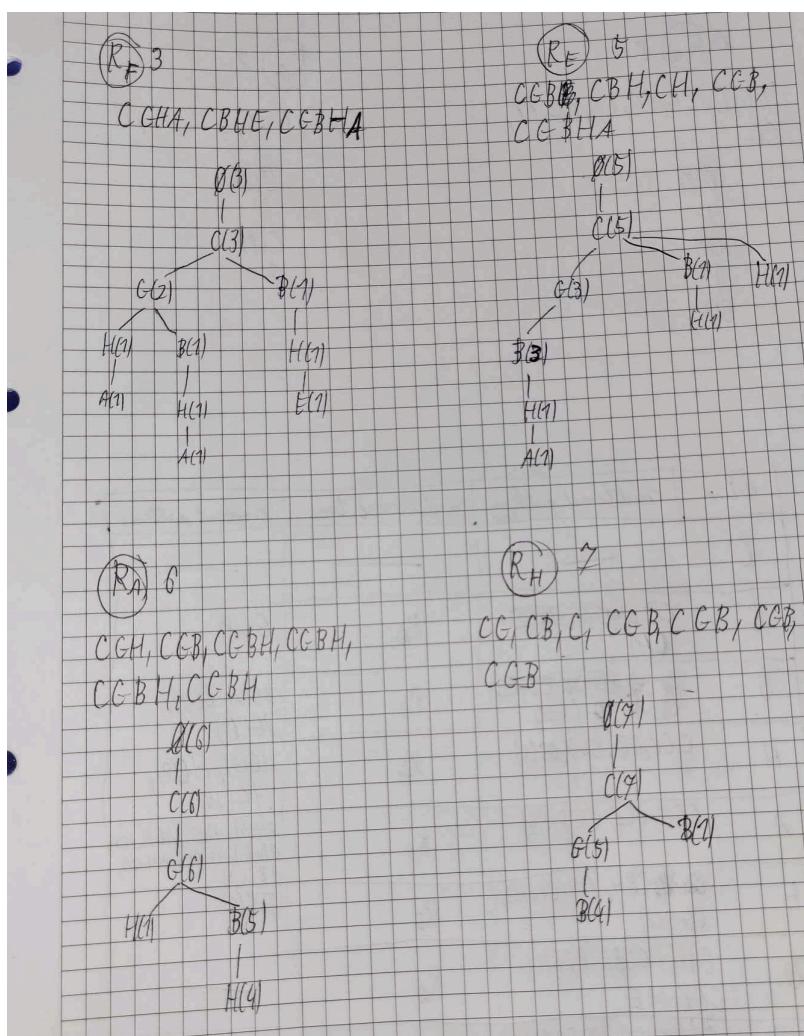
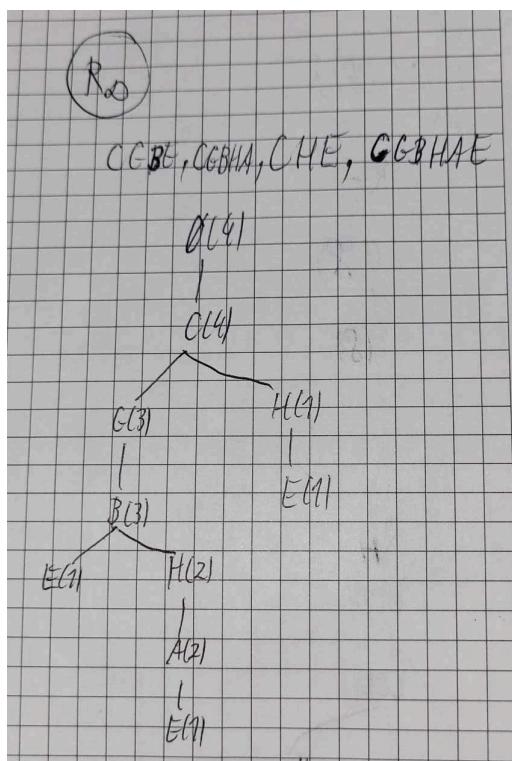
P=B:

{BC(8), BG(7), BGC(7)}

P=G:

{GC(8)}

## Conditional FP-trees:



$\emptyset(R_B)$  8

CG, C, CG, CG, CG, CG, CG,  
CG

$\emptyset(8)$

|  
C(8)

|  
G(8)

$R_G$  8

C, C, C, C, C, C, C, C

$\emptyset(8)$

|

C(8)

8

item	conditional pattern base	conditional tree	frequent patterns
C	—	—	—
G	<del>CG, C, CG, CG, CG, CG, CG, CG</del> CH(8)	$R_G$	GC(8)
B	<del>CG, C, CG, CG, CG, CG, CG, CG</del> C, CG(7)	$R_B$	BC(7), BC(8), BGC(7)
H	C, CG, CB, CGB <del>(6)</del> (6)	$R_H$	HG(5), HC(9), HGC(5)
A	CCH, CGB, CGBH(4)	$R_A$	AH(5), AB(5), AC(6), AC(6), ABG(5), AHG(5), AHG(5), ABC(5)
E	<del>CCB, CH, CCBHA, CGB(2)</del> CGB(2)	$R_E$	EC(5)
D	CHE, CGBE, CGBHA, CGBHAE	$R_D$	—
F	CGHA, CBHE, CGBHA	$R_F$	—

#### 4. Mixtend

FP-growth Result:		
	support	itemsets
0	1.0	(Kidney Beans)
1	0.8	(Eggs)
2	0.6	(Bread)
3	0.6	(Yogurt)
4	0.6	(Onion)
5	0.8	(Eggs, Kidney Beans)
6	0.6	(Bread, Kidney Beans)
7	0.6	(Yogurt, Kidney Beans)
8	0.6	(Eggs, Onion)
9	0.6	(Onion, Kidney Beans)
10	0.6	(Eggs, Onion, Kidney Beans)

Apriori Result:		
	support	itemsets
0	0.6	(Bread)
1	0.8	(Eggs)
2	1.0	(Kidney Beans)
3	0.6	(Onion)
4	0.6	(Yogurt)
5	0.6	(Bread, Kidney Beans)
6	0.8	(Eggs, Kidney Beans)
7	0.6	(Eggs, Onion)
8	0.6	(Onion, Kidney Beans)
9	0.6	(Yogurt, Kidney Beans)
10	0.6	(Eggs, Onion, Kidney Beans)

#### 5. Compact Representation of Frequent Itemsets

To solve this problem we decided to use this algorithm:

---

##### Algorithm 5.4 Support counting using closed frequent itemsets.

---

- 1: Let  $C$  denote the set of closed frequent itemsets and  $F$  denote the set of all frequent itemsets.
  - 2: Let  $k_{\max}$  denote the maximum size of closed frequent itemsets
  - 3:  $F_{k_{\max}} = \{f | f \in C, |f| = k_{\max}\}$  {Find all frequent itemsets of size  $k_{\max}$ .}
  - 4: **for**  $k = k_{\max} - 1$  down to 1 **do**
  - 5:    $F_k = \{f | f \in F, |f| = k\}$  {Find all frequent itemsets of size  $k$ .}
  - 6:   **for each**  $f \in F_k$  **do**
  - 7:     **if**  $f \notin C$  **then**
  - 8:        $f.support = \max\{f'.support | f' \in F_{k+1}, f \subset f'\}$
  - 9:     **end if**
  - 10:   **end for**
  - 11: **end for**
-

With yellow we have marked all the closed frequent itemset given in Table 4. Then, following the algorithm we denoted  $k_{max} = 4$ , the maximum size of the closed frequent itemsets. Then we presided to find all the itemsets of length  $k = k_{max} - 1$  where  $k$  equaled 3 in our case. For each of those itemsets, that were not part of the closed frequent itemsets, we performed the operation:  $f.\text{support} = \max\{f . \text{support}|f \in F_{k+1}, f \subset f\}$ . This meant, we assigned every itemset that was the subset of the closed frequent itemset from  $F_{k+1}$  the same support count value. We repeated the operation until  $k = 1$  and all the frequent itemsets were given their support count.

Items	Support Count						
a	11	ab	7	abc		abcd	
b	10	ac	6	abd		abce	
c	6	ad	11	abe	7	abde	
d	13	ae	7	acd	6	acde	5
e	8	bc		ace	5	bcde	
		bd	7	ade	5		
		be	8	bcd			
		cd	6	bce			
		ce	5	bde	4		
		de	6	cde	5		