

# Øving 2 Datavarehus

Alexander Fredheim

February 2020

# 1 Apriori Algorithm

## 1.a

Here we are going to use the apriori algorithm to generate frequent itemsets, which we are later going to use to find association rules.

We have minsupport = 0.5 which for our table means that every itemset must have a minimum support of 5 in our initial dataset.

| C1 | Itemset | Supportcount |
|----|---------|--------------|
|    | {A}     | 6            |
|    | {B}     | 8            |
|    | {C}     | 10           |
|    | {D}     | 4            |
|    | {E}     | 5            |
|    | {F}     | 3            |
|    | {G}     | 8            |
|    | {H}     | 7            |

Table 1: our initial 1-element set

| L1 | Itemset | Supportcount |
|----|---------|--------------|
|    | {A}     | 6            |
|    | {B}     | 8            |
|    | {C}     | 10           |
|    | {E}     | 5            |
|    | {G}     | 8            |
|    | {H}     | 7            |

Table 2: frequent 1-element sets with supportcount  $\geq$  minsupport

| C2 | Itemset | Supportcount |
|----|---------|--------------|
|    | {A, B}  | 5            |
|    | {A, C}  | 6            |
|    | {A, E}  | 1            |
|    | {A, G}  | 6            |
|    | {A, H}  | 5            |
|    | {B, C}  | 8            |
|    | {B, E}  | 4            |
|    | {B, G}  | 7            |
|    | {B, H}  | 5            |
|    | {C, E}  | 5            |
|    | {C, G}  | 8            |
|    | {C, H}  | 7            |
|    | {E, G}  | 3            |
|    | {E, H}  | 3            |
|    | {G, H}  | 5            |

Table 3: 2-element sets generated from frequent 1-element sets

| L2 | Itemset | Supportcount |
|----|---------|--------------|
|    | {A, B}  | 5            |
|    | {A, C}  | 6            |
|    | {A, G}  | 6            |
|    | {A, H}  | 5            |
|    | {B, C}  | 8            |
|    | {B, G}  | 7            |
|    | {B, H}  | 5            |
|    | {C, E}  | 5            |
|    | {C, G}  | 8            |
|    | {C, H}  | 7            |
|    | {G, H}  | 5            |

Table 4: frequent 2-element sets with supportcount  $\geq$  minsupport

| C3 | Itemset   | Supportcount |
|----|-----------|--------------|
|    | {A, B, C} | 5            |
|    | {A, B, G} | 5            |
|    | {A, B, H} | 4            |
|    | {A, C, G} | 6            |
|    | {A, C, H} | 5            |
|    | {A, G, H} | 5            |
|    | {B, C, G} | 7            |
|    | {B, C, H} | 5            |
|    | {B, G, H} | 4            |
|    | {C, G, H} | 5            |

Table 5: 3-element sets generated from frequent 2-element sets

Note that in table 5 we don't include {C, E, G} and {C, E, H} as they have subsets {E, G} and {E, H} which doesn't meet the minimum support requirement

| L3 | Itemset   | Supportcount |
|----|-----------|--------------|
|    | {A, B, C} | 5            |
|    | {A, B, G} | 5            |
|    | {A, C, G} | 6            |
|    | {A, C, H} | 5            |
|    | {A, G, H} | 5            |
|    | {B, C, G} | 7            |
|    | {B, C, H} | 5            |
|    | {C, G, H} | 5            |

Table 6: frequent 3-element sets with supportcount  $\geq$  minsupport

| C4 | Itemset      | Supportcount |
|----|--------------|--------------|
|    | {A, B, C, G} | 5            |
|    | {A, C, G, H} | 5            |

Table 7: 4-element sets generated from frequent 3-element sets

Note that in table 7 we don't include {B, C, G, H} as it has subset {B, G, H} which doesn't meet the minimum support requirement

| L4 | Itemset      | Supportcount |
|----|--------------|--------------|
|    | {A, B, C, G} | 5            |
|    | {A, C, G, H} | 5            |

Table 8: frequent 4-element sets with supportcount  $\geq$  minsupport

## 1.b

For this task we are going to generate some rules for our frequent 4-element sets  $\{A, B, C, G\}$  and  $\{A, C, G, H\}$ . We use  $\text{minconf} = 0.8$ .

Rule generation for element set  $\{A, B, C, G\}$ :

| Rules                           | Confidence |
|---------------------------------|------------|
| $\{B, C, G\} \Rightarrow \{A\}$ | 0.71       |
| $\{A, C, G\} \Rightarrow \{B\}$ | 0.83       |
| $\{A, B, G\} \Rightarrow \{C\}$ | 1.00       |
| $\{A, B, C\} \Rightarrow \{G\}$ | 1.00       |

Table 9: possible rules generated from  $\{A, B, C, G\}$  and their corresponding confidence

We see in table 9 that  $\{B, C, G\} \Rightarrow \{A\}$  does not meet our confidence threshold. We can therefore prune all possible rules containing A on the right hand side

| Rules                           | Confidence |
|---------------------------------|------------|
| $\{A, G\} \Rightarrow \{B, C\}$ | 0.83       |
| $\{A, C\} \Rightarrow \{B, G\}$ | 0.83       |
| $\{A, B\} \Rightarrow \{C, G\}$ | 1.00       |

Table 10: possible rules further generated from the rules of  $\{A, C, G\}$ ,  $\{A, B, G\}$  and  $\{A, B, C\}$  and their corresponding confidence. After pruning

| Rules                           | Confidence |
|---------------------------------|------------|
| $\{A\} \Rightarrow \{B, C, G\}$ | 0.83       |

Table 11: possible rules further generated from the rules of  $\{A, G\}$ ,  $\{A, C\}$  and  $\{A, B\}$  and its corresponding confidence. After pruning

Rule generation for element set  $\{A, C, G, H\}$ :

| Rules                           | Confidence |
|---------------------------------|------------|
| $\{C, G, H\} \Rightarrow \{A\}$ | 1.00       |
| $\{A, G, H\} \Rightarrow \{C\}$ | 1.00       |
| $\{A, C, H\} \Rightarrow \{G\}$ | 1.00       |
| $\{A, C, G\} \Rightarrow \{H\}$ | 0.83       |

Table 12: possible rules generated from  $\{A, C, G, H\}$  and their corresponding confidence

| Rules                           | Confidence |
|---------------------------------|------------|
| $\{G, H\} \Rightarrow \{A, C\}$ | 1.00       |
| $\{C, H\} \Rightarrow \{A, G\}$ | 0.71       |
| $\{C, G\} \Rightarrow \{A, H\}$ | 0.63       |
| $\{A, H\} \Rightarrow \{C, G\}$ | 1.00       |
| $\{A, G\} \Rightarrow \{C, H\}$ | 0.83       |
| $\{A, C\} \Rightarrow \{G, H\}$ | 0.83       |

Table 13: possible rules further generated from the rules of  $\{G, H\}$ ,  $\{C, H\}$ ,  $\{C, G\}$ ,  $\{A, H\}$ ,  $\{A, G\}$  and  $\{A, C\}$  and its corresponding confidence.

We see in table 13 that  $\{C, H\} \Rightarrow \{A, G\}$  and  $\{C, G\} \Rightarrow \{A, H\}$  does not meet our confidence threshold. We can therefore prune all possible rules containing subsets of either  $\{A, G\}$  or  $\{A, H\}$  on the right hand side of the rule.

| Rules                           | Confidence |
|---------------------------------|------------|
| $\{A\} \Rightarrow \{C, G, H\}$ | 0.83       |

Table 14: the only rule derived from the rules in table 13, not containing supersets of  $\{A, G\}$  or  $\{A, H\}$

we end up with the following rules for  $\{A, B, C, G\}$  and  $\{A, C, G, H\}$ :

| Rules                           | Confidence |
|---------------------------------|------------|
| $\{A, C, G\} \Rightarrow \{B\}$ | 0.83       |
| $\{A, B, G\} \Rightarrow \{C\}$ | 1.00       |
| $\{A, B, C\} \Rightarrow \{G\}$ | 1.00       |
| $\{A, G\} \Rightarrow \{B, C\}$ | 0.83       |
| $\{A, C\} \Rightarrow \{B, G\}$ | 0.83       |
| $\{A, B\} \Rightarrow \{C, G\}$ | 1.00       |
| $\{A\} \Rightarrow \{B, C, G\}$ | 0.83       |
| $\{C, G, H\} \Rightarrow \{A\}$ | 1.00       |
| $\{A, G, H\} \Rightarrow \{C\}$ | 1.00       |
| $\{A, C, H\} \Rightarrow \{G\}$ | 1.00       |
| $\{A, C, G\} \Rightarrow \{H\}$ | 0.83       |
| $\{G, H\} \Rightarrow \{A, C\}$ | 1.00       |
| $\{A, H\} \Rightarrow \{C, G\}$ | 1.00       |
| $\{A, G\} \Rightarrow \{C, H\}$ | 0.83       |
| $\{A, C\} \Rightarrow \{G, H\}$ | 0.83       |
| $\{A\} \Rightarrow \{C, G, H\}$ | 0.83       |

Table 15: all rules generated from  $\{A, B, C, G\}$  and  $\{A, C, G, H\}$

## 2 FP-Growth

We start by scanning the items and sorting it by support (decreasing). Our minsupport is still 0.5 which means minimum support is still 5

| Item | Supportcount |
|------|--------------|
| C    | 10           |
| B    | 8            |
| G    | 8            |
| H    | 7            |
| A    | 6            |
| E    | 5            |

Table 16: our elements sorted by frequency. Elements not meeting support requirements removed

| TID | Initial Items | Ordered Items |
|-----|---------------|---------------|
| 110 | A,C,F,G,H     | C,G,H,A       |
| 111 | B,C,E,D,G     | C,B,G,E       |
| 112 | B,C,E,F,G     | C,B,G,E       |
| 113 | A,B,C,G       | C,B,G,A       |
| 114 | C,D,E,H       | C,H,E         |
| 115 | A,B,C,G,H     | C,B,G,H,A     |
| 116 | A,B,C,D,G,H   | C,B,G,H,A     |
| 117 | B,C,E,G       | C,B,G,E       |
| 118 | A,B,C,F,G,H   | C,B,G,H,A     |
| 119 | A,B,C,D,E,G,H | C,B,G,H,A,E   |

Table 17: Transactions with the prioritized ordering



we now create the FP-tree, transaction by transaction, following the order in "ordered items from table 17"

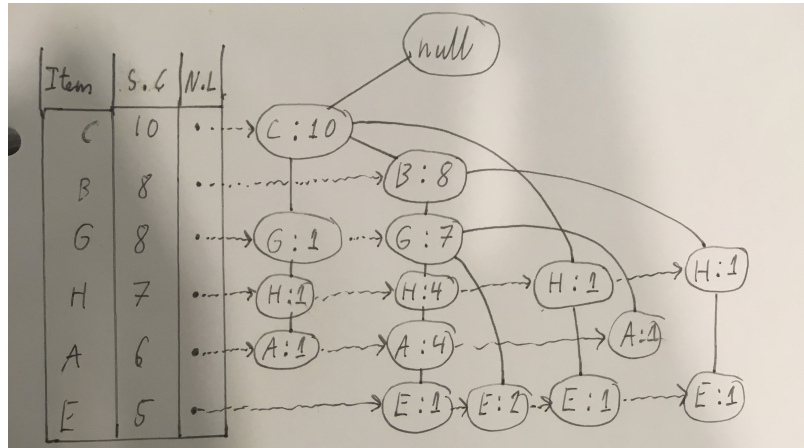


Figure 1: FP-tree

we now need to find the frequent itemsets by constructing a conditional FP-tree. To do this we start bottom up for the items in table 16 and look at all trees ending with the given item. We then use the divide and conquer approach where necessary to generate the needed trees. we get the following table

| Item | Conditional Pattern Base                                   | Cond-itional FP-tree   | Frequent Patterns Generated   |
|------|--|--|---|
| E    | $\{\{C,B,H:1\}, \{C,B,G:2\}, \{C,H:1\}, \{C,B,G,H,A:1\}\}$ | $\langle C:5 \rangle$  | $\{C,E:5\}$   |
| A    | $\{\{C,G,H:1\}, \{C,B,G:1\}, \{C,B,G,H:4\}\}$              | $\langle C:6,G:6,H:1 \rangle, \langle C:6,G:6,B:5,H:4 \rangle$ | $\{C,A:6\}, \{G,A:6\}, \{H,A:5\}, \{B,A:5\}, \{C,G,A:6\}, \{C,B,A:5\}, \{G,B,A:5\}, \{C,H,A:5\}, \{G,H,A:5\}, \{C,G,H,A:5\}, \{C,G,B,A:5\}$ |
| H    | $\{\{C,G:1\}, \{C,B:1\}, \{C:1\}, \{C,B,G:4\}\}$           | $\langle C:7,G:1 \rangle, \langle C:7,B:5,G:4 \rangle$         | $\{C,H:7\}, \{B,H:5\}, \{G,H:5\}, \{C,B,H:5\}, \{C,G,H:5\}$   |
| G    | $\{\{C:1\}, \{C,B:7\}\}$                                   | $\langle C:8,B:7 \rangle$                                      | $\{C,G:8\}, \{B,G:7\}, \{C,B,G:7\}$   |
| B    | $\{\{C:8\}\}$  | $\langle C:8 \rangle$  | $\{C,B:8\}$   |

Table 18: Table showing the generated frequent patterns

### 3 KNIME

The node setup from KNIME:

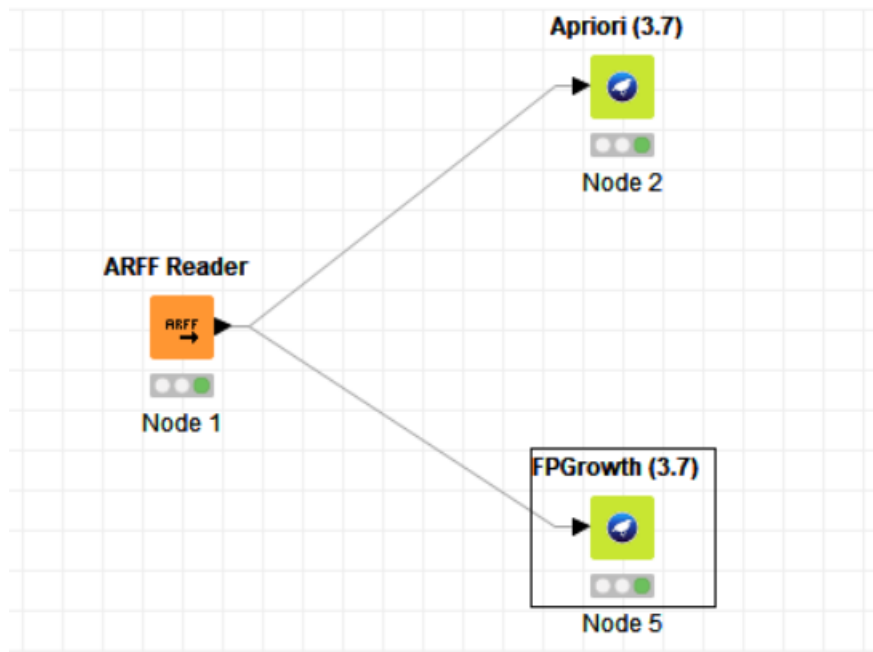


Figure 2: Nodes in KNIME

Apriori rules:

```
File

Apriori
=====

Minimum support: 0.5 (5 instances)
Minimum metric <confidence>: 0.8
Number of cycles performed: 10

Generated sets of large itemsets:

Size of set of large itemsets L(1): 6
Size of set of large itemsets L(2): 11
Size of set of large itemsets L(3): 8
Size of set of large itemsets L(4): 2

Best rules found:

1. B=t 8 ==> C=t 8 <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
2. G=t 8 ==> C=t 8 <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
3. H=t 7 ==> C=t 7 <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
4. B=t G=t 7 ==> C=t 7 <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
5. A=t 6 ==> C=t 6 <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
6. A=t 6 ==> G=t 6 <conf:(1)> lift:(1.25) lev:(0.12) [1] conv:(1.2)
7. A=t G=t 6 ==> C=t 6 <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
8. A=t C=t 6 ==> G=t 6 <conf:(1)> lift:(1.25) lev:(0.12) [1] conv:(1.2)
9. A=t 6 ==> C=t G=t 6 <conf:(1)> lift:(1.25) lev:(0.12) [1] conv:(1.2)
10. E=t 5 ==> C=t 5 <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
11. A=t B=t 5 ==> C=t 5 <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
12. A=t B=t 5 ==> G=t 5 <conf:(1)> lift:(1.25) lev:(0.1) [0] conv:(1)
13. A=t H=t 5 ==> C=t 5 <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
14. G=t H=t 5 ==> A=t 5 <conf:(1)> lift:(1.67) lev:(0.2) [2] conv:(2)
15. A=t H=t 5 ==> G=t 5 <conf:(1)> lift:(1.25) lev:(0.1) [0] conv:(1)
16. B=t H=t 5 ==> C=t 5 <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
17. G=t H=t 5 ==> C=t 5 <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
18. A=t B=t G=t 5 ==> C=t 5 <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
19. A=t B=t C=t 5 ==> G=t 5 <conf:(1)> lift:(1.25) lev:(0.1) [0] conv:(1)
20. A=t B=t 5 ==> C=t G=t 5 <conf:(1)> lift:(1.25) lev:(0.1) [0] conv:(1)
21. C=t G=t H=t 5 ==> A=t 5 <conf:(1)> lift:(1.67) lev:(0.2) [2] conv:(2)
22. A=t G=t H=t 5 ==> C=t 5 <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
23. A=t C=t H=t 5 ==> G=t 5 <conf:(1)> lift:(1.25) lev:(0.1) [0] conv:(1)
24. G=t H=t 5 ==> A=t C=t 5 <conf:(1)> lift:(1.67) lev:(0.2) [2] conv:(2)
25. A=t H=t 5 ==> C=t G=t 5 <conf:(1)> lift:(1.25) lev:(0.1) [0] conv:(1)
26. G=t 8 ==> B=t 7 <conf:(0.88)> lift:(1.09) lev:(0.06) [0] conv:(0.8)
27. B=t 8 ==> G=t 7 <conf:(0.88)> lift:(1.09) lev:(0.06) [0] conv:(0.8)
28. C=t G=t 8 ==> B=t 7 <conf:(0.88)> lift:(1.09) lev:(0.06) [0] conv:(0.8)
29. B=t C=t 8 ==> G=t 7 <conf:(0.88)> lift:(1.09) lev:(0.06) [0] conv:(0.8)
30. G=t 8 ==> B=t C=t 7 <conf:(0.88)> lift:(1.09) lev:(0.06) [0] conv:(0.8)
31. B=t 8 ==> C=t G=t 7 <conf:(0.88)> lift:(1.09) lev:(0.06) [0] conv:(0.8)
32. A=t 6 ==> B=t 5 <conf:(0.83)> lift:(1.04) lev:(0.02) [0] conv:(0.6)
33. A=t 6 ==> H=t 5 <conf:(0.83)> lift:(1.19) lev:(0.08) [0] conv:(0.9)
34. A=t C=t 6 ==> B=t 5 <conf:(0.83)> lift:(1.04) lev:(0.02) [0] conv:(0.6)
35. A=t 6 ==> B=t C=t 5 <conf:(0.83)> lift:(1.04) lev:(0.02) [0] conv:(0.6)
```

Figure 3: the 35 "best" rules generated by the apriori node in KNIME

FP rules:

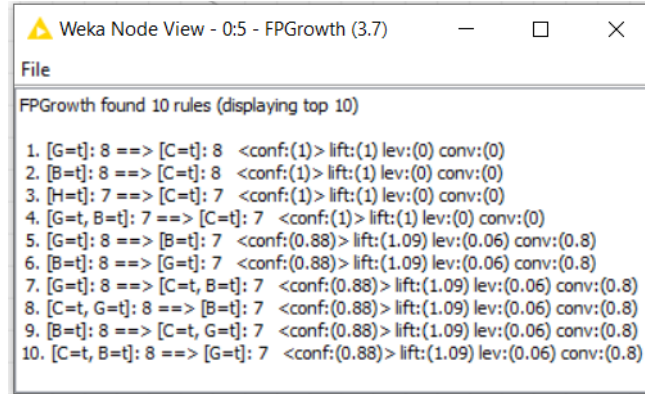


Figure 4: 10 rules generated by the FP-Growth node in KNIME

## 4 Compact Representation of Frequent Itemsets

The algorithm can be described by to steps:

1. generate frequent element sets by taking the union of all subsets of the closed frequent element sets
2. evaluate the support of any non-closed frequent k-element set by using the principle of k+1 frequent supersets. The support of any non frequent k-elementset is the maximum support of the frequent supersets.

This gives us the following table:

| Frequent Element set | Support |
|----------------------|---------|
| {a}                  | 11      |
| {b}                  | 10      |
| {c}                  | 6       |
| {d}                  | 13      |
| {e}                  | 8       |
| {a,b}                | 7       |
| {a,c}                | 6       |
| {a,d}                | 11      |
| {a,e}                | 7       |
| {b,d}                | 7       |
| {b,e}                | 8       |
| {c,d}                | 6       |
| {c,e}                | 5       |
| {d,e}                | 6       |
| {a,b,e}              | 7       |
| {a,c,d}              | 6       |
| {a,c,e}              | 5       |
| {a,d,e}              | 5       |
| {b,d,e}              | 4       |
| {c,d,e}              | 5       |
| {a,c,d,e}            | 5       |

Table 19: frequent item sets and their support