

Exercise Sheet 3

Exercise 15 Eclat Algorithm

- a) Execute the Eclat algorithm (using the prefix tree scheme) for the transaction database shown on the right and a minimum support of $s_{\min} = 3!$
- | | |
|----------------|----------------|
| <i>a d f</i> | <i>a b d</i> |
| <i>a c d e</i> | <i>b d e</i> |
| <i>b d</i> | <i>b c e g</i> |
| <i>b c d</i> | <i>c d f</i> |
| <i>b c</i> | <i>a b d</i> |
- b) How is support counted in the Eclat algorithm?
How does this scheme differ from the scheme used in Apriori? What are the advantages/disadvantages of this counting scheme?
- c) Why does Eclat (usually) not exploit the Apriori property fully?
How could it be made to exploit the Apriori property fully?
- d) Why is it advantageous in the Eclat algorithm to process less frequent items first?
(Hint: Consider the average length of the transaction identifier lists.)

Exercise 16 Eclat Algorithm

- a) In the conditional transaction databases occurring in the search for frequent item sets, it is often possible to combine transactions that have become equal after some items have been eliminated. How can this be exploited in the Eclat algorithm?
- b) Why is a bit vector representation of the transaction identifier lists not necessarily better than an explicit list of transaction identifiers? (Hint: What happens to the transaction identifier lists in the recursive processing?)
- c) Show that the diffset for an item (w.r.t. some conditional database), that is, $D(a | I) = K(I) - K(I \cup \{a\})$, can be both larger and smaller than its cover (depending on the database).

Exercise 17 Item Order

- a) How can one exploit a local item order in the Apriori algorithm?
- b) How can one exploit a local item order in the Eclat algorithm?
Why is it particularly easy to use a local item order in it?
- c) How can one exploit the order in which a divide-and-conquer scheme visits the item sets to optimize the reporting of found frequent item sets?
Does this also work for a local item order? Why?

Exercise 18 LCM Algorithm

- a) How does the occurrence deliver scheme of LCM work?
How is it related to the Eclat algorithm? How is it different?

What additional data/information is needed for the LCM algorithm (compared to Eclat)?

b) Why does LCM allow mining with memory proportional to the database size and Eclat does not, even though their processing schemes are similar?

(Hint: Consider the order in which the subproblems are processed.)

c) Execute the LCM algorithm for the transaction database shown on the right and a minimum support of $s_{\min} = 3!$

<i>a d f</i>	<i>a b d</i>
<i>a c d e</i>	<i>b d e</i>
<i>b d</i>	<i>b c e g</i>
<i>b c d</i>	<i>c d f</i>
<i>b c</i>	<i>a b d</i>

Exercise 19 SaM and RElim Algorithms

a) How is support counted in the SaM algorithm?

What does SaM have to ensure for the transaction array for this to work?

b) Why can the merge step of the SaM algorithm be a problem?

c) How does the RElim algorithm improve over the SaM algorithm?

What sorting principle is exploited?

d) Construct the initial database representation for the RElim algorithm for the transaction database shown on the right!

<i>a d e</i>	<i>a c d</i>
<i>b c d</i>	<i>b c</i>
<i>a c e</i>	<i>a c d e</i>
<i>a c d e</i>	<i>c b e</i>
<i>a e</i>	<i>a d e</i>

e) How is a k -items machine related to the RElim algorithm?

What are the differences and the advantages/disadvantages?

Exercise 20 FP-growth Algorithm

a) What is an FP-tree? How is it constructed from a transaction database? Why does it combine a horizontal and a vertical representation?

b) How is support counted with an FP-tree?

c) How is an FP-tree projected? What does the projection yield?

d) How do we obtain the conditional transaction database for the second subproblem in the FP-growth algorithm?

e) Execute the FP-growth algorithm for the transaction database shown on the right and a minimum support $s_{\min} = 3!$

<i>a d e</i>	<i>a c d</i>
<i>b c d</i>	<i>b c</i>
<i>a c e</i>	<i>a c d e</i>
<i>a c d e</i>	<i>b c e</i>
<i>a e</i>	<i>a d e</i>

You may skip the steps of finding the item frequencies and resorting the items in the transactions w.r.t. the frequency order (that is: use the alphabetical order).