I use an enumeration tree structure to store the projected database, at each node, name P, I maintain candidate extension items for node P, and a transaction database that only contains transactions that contain the current itemset of P and in every transaction, only those who are lexicographically greater than the current item of P are preserved. The current itemset of every node is generated by add one element in the candidate extension of its parent node to the current itemset of its parent node. The count of every node is achieved by counting transactions of its parent node that contain the element extended.

In the code, I used a two-dimensional ArrayList, which is a space expensive and stupid way, to store every projected database.

I just realized that sparse matrix might be a better choice to store the database. Each line of the input file is stored as an element in a m*2 int array(m is the number of lines in the input file, for each line, index 0 stores the transaction number, and index 1 stores the item). For each node maintain an indicator array that contains pointers to the first element of each transaction that contains the current itemset of node P, and we start from the first element to search for the candidate extension now we are looking at. If we find the extension item, the count of the node P increments. This might be a quicker way to do the projection, but since I start the project assignment late, time is not enough for me to implement it.

In order to improve efficiency, I tried the buckets method mentioned in Charu's book in section 4.4.3.2. But it doesn't work better than the previous one. Maybe It's my bad

Because of the low efficiency of my code, I don't have the time to run all the combination of minconf and options value using large dataset before due time.

Therefore, the following charts are based on the **small dataset**. I guess it's enough to show the correctness of my code.

Option 2 results in faster execution time because ordering the items from least support to greatest support ensures that the computationally heavier branches of the enumeration tree have fewer relevant transactions. This helps maximizing the selectivity of projections and thus ensures better efficiency.

## conf:0.8;option:1

| | time(s) |
|---|---|
| 15 | 365 |
| 20 | 141.7 |
| 30 | 48 |
| 50 | 16.9 |
| 100 | 6.3 |
| 500 | 0.86 |
| 1000 | 0.48 |

Time

meta-chart.com

## conf:0.8;option:2

| | time(s) |
|---|---|
| 15 | 4.5 |
| 20 | 2.2 |
| 30 | 5.9 |
| 50 | 0.5 |
| 100 | 0.1 |
| 500 | 0.007 |
| 1000 | 0.004 |

Time

meta-chart.com

## conf:0.8;option:3

| | time(s) |
|---|---|
| 15 | 129.7 |
| 20 | 32.6 |
| 30 | 10.5 |
| 50 | 2 |
| 100 | 0.2 |
| 500 | 0.007 |
| 1000 | 0.004 |

time(s)

Time

## conf:0.9;option:1

| | time(s) |
|---|---|
| 15 | 365 |
| 20 | 141.7 |
| 30 | 45.9 |
| 50 | 18.6 |
| 100 | 6 |
| 500 | 1 |
| 1000 | 0.4 |

time(s)

Time

# conf:0.9;option:2

| | |
|---|---|
| 15 | 4.5 |
| 20 | 2.2 |
| 30 | 4.2 |
| 50 | 0.5 |
| 100 | 0.1 |
| 500 | 0.007 |
| 1000 | 0.003 |

time(s)

Time

meta-chart.com

# conf:0.9;option:3

| | |
|---|---|
| 15 | 129.7 |
| 20 | 32.6 |
| 30 | 9.5 |
| 50 | 1.8 |
| 100 | 0.2 |
| 500 | 0.01 |
| 1000 | 0.004 |

time(s)

Time

meta-chart.com

**conf:0.95;option:1**

| | time(s) |
|---|---|
| 15 | 365 |
| 20 | 141.7 |
| 30 | 47.2 |
| 50 | 17.5 |
| 100 | 6.2 |
| 500 | 0.9 |
| 1000 | 0.4 |

Time

meta-chart.com



**conf:0.95;option:2**

| | time(s) |
|---|---|
| 15 | 4.5 |
| 20 | 2.2 |
| 30 | 4.1 |
| 50 | 0.5 |
| 100 | 0.1 |
| 500 | 0.006 |
| 1000 | 0.003 |

Time

meta-chart.com

**conf:0.95;option:3**

| | time(s) |
|---|---|
| 15 | 129.7 |
| 20 | 32.6 |
| 30 | 8.5 |
| 50 | 1.8 |
| 100 | 0.3 |
| 500 | 0.008 |
| 1000 | 0.004 |

Time

meta-chart.com



**support count**

| | support/count |
|---|---|
| 15 | 214409 |
| 20 | 50762 |
| 30 | 7471 |
| 50 | 1195 |
| 100 | 230 |
| 500 | 7 |
| 1000 | 1 |

rule counts;conf:0.8

rule count

| 15 | 20 | 30 | 50 | 100 | 500 | 1000 |
|---|---|---|---|---|---|---|
| 0 | 0 | 16709 | 367 | 32 | 0 | 0 |



rule counts;conf:0.9

rule count

| 15 | 20 | 30 | 50 | 100 | 500 | 1000 |
|---|---|---|---|---|---|---|
| 0 | 0 | 6724 | 123 | 4 | 0 | 0 |

rule counts;conf:0.95