**Data structures:**

Each transaction is stored as an ArrayList of Integers. The database of transactions is therefore an ArrayList of ArrayLists of Integers. All projected databases are also stored in this manner.

Frequent itemsets are represented as space-delimited strings; those strings are hashed to their corresponding frequencies via a HashMap. These strings provide the path of the recursive mining algorithm.

**Projection Algorithm Implementation:**

The algorithm recursively calls generateFrequentItemsets() in a depth-first manner. At each step, the two arguments passed to this function are the projected database, and the current itemset (represented as a String so that it can be hashed). Initially it is passed the entire database and an empty string.

The function first scans the projected database for any items that are frequent; the frequency of items is tracked with a HashMap. Then it constructs a new database which is a copy of the original, but only keeps items which have a frequency greater or equal to minsup. This process filters out any itemsets that are infrequent. At the same time, any item that has a frequency greater than minsup is added to an ArrayList of candidates. The candidates are sorted in lexicographic order so that they will be expanded in order.

The function proceeds to iterate over and project the database onto all of the candidates (this is done in a depth-first manner to avoid large space complexity). The projection is done by iterating over the original database, and simultaneously constructing a new one. Any row that contains the item being projected onto is sliced at that item's index. This is possible to do because all transactions are sorted in lexicographic order. Any row not containing the item is ignored.

Since the algorithm can only see the item it's currently projecting on, it needs to keep track of the path it took to get to this candidate. Therefore, the candidate is added to the current itemset (it is appended to the string in a space-delimited manner) to construct a new itemset. Also, this candidate was frequent, which means the newly-constructed itemset is frequent. The itemset string is hashed to the frequency of the candidate. With this updated itemset, and the projected database, the function is called recursively. The itemset grows with the depth of the algorithm, and provides the algorithm its current path.
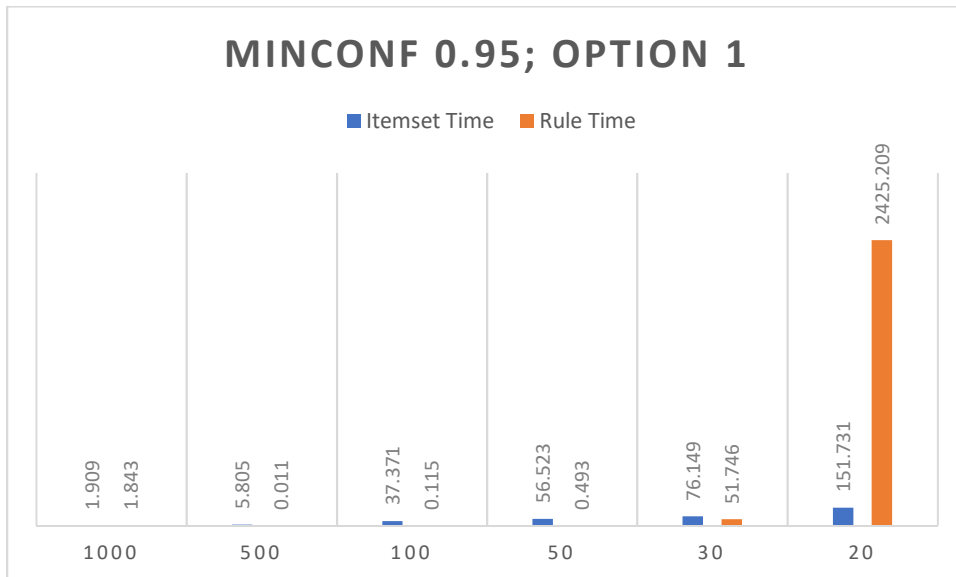
**Algorithm Run Times:**

A few notes about my data:

> *Though I've run the program for virtually all Minsup/Options/Minconf combinations at one point or another for testing purposes, unfortunately I underestimated how much time it would take to run and record all combinations. As a result, some values are not filled in: in particular, I was unable to generate association rules for some combinations of Minsup 15 and 20. This is because I understood from the instructions that I wouldn't have to generate these, but after reading the discussion forum, I found that I was mistaken (and left with not enough time to generate them).*

*Also, it should be noted that some of these times ended up quite inconsistent. Though all tests were run on CSE Lab machines, I had to switch machines a few times. I found that some machines performed better than others. On the run I eventually recorded, generating frequent itemsets for minsup 15 took a very long time. However, I've run this many times before and gotten times less than 20 minutes. As a result, I'm not sure how reliable much of this data is anyways.*
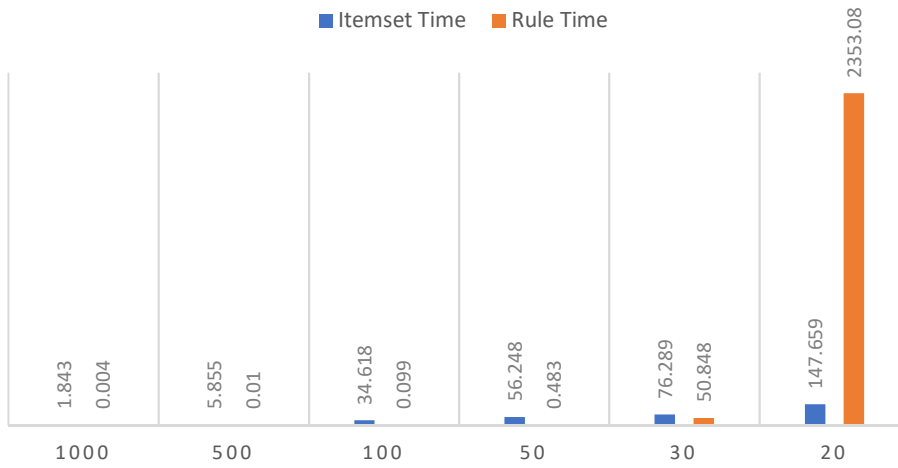
*Finally, I wasn't exactly clear on how many bar charts to make. Though the instructions say to make two bar charts for each minconf and options combination (so 3 \* 3 \* 2 = 18), this doesn't make much sense to me. The second bar chart is supposed to be the number of rules and itemsets generated for different minsups. However, it doesn't matter what option is chosen; the number of rules and itemsets will stay the same. As a result, I generated 12 charts total. 9 for the first bar chart, and then 3 charts for the different minconf values.*

*In the bar charts, the x axis represents minsup.*



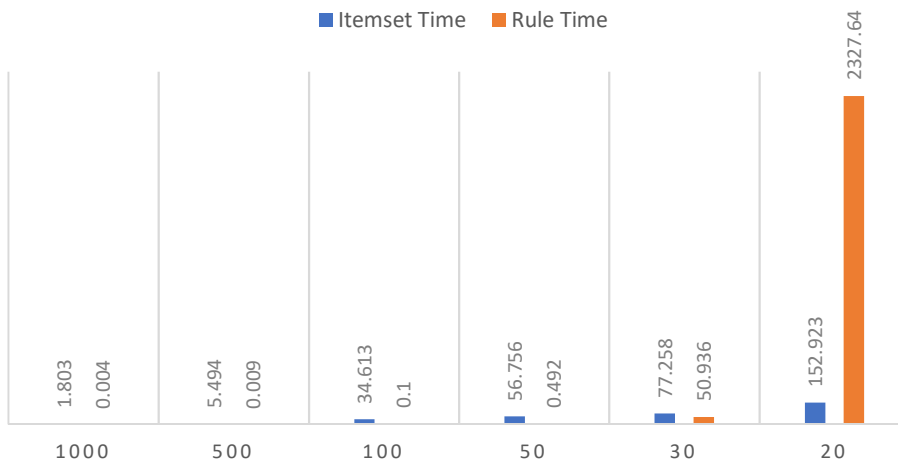**MINCONF 0.95; OPTION 1**

■ Itemset Time   ■ Rule Time

| minsup | Itemset Time | Rule Time |
|---|---|---|
| 1000 | 1.909 | 1.843 |
| 500 | 5.805 | 0.011 |
| 100 | 37.371 | 0.115 |
| 50 | 56.523 | 0.493 |
| 30 | 76.149 | 51.746 |
| 20 | 151.731 | 2425.209 |

# MINCONF 0.95; OPTION 2

■ Itemset Time ■ Rule Time

| | 1000 | 500 | 100 | 50 | 30 | 20 |
|---|---|---|---|---|---|---|
| Itemset Time | 1.843 | 5.855 | 34.618 | 56.248 | 76.289 | 147.659 |
| Rule Time | 0.004 | 0.01 | 0.099 | 0.483 | 50.848 | 2353.08 |

# MINCONF 0.95; OPTION 3

■ Itemset Time ■ Rule Time

| | 1000 | 500 | 100 | 50 | 30 | 20 |
|---|---|---|---|---|---|---|
| Itemset Time | 1.803 | 5.494 | 34.613 | 56.756 | 77.258 | 152.923 |
| Rule Time | 0.004 | 0.009 | 0.1 | 0.492 | 50.936 | 2327.64 |

**MINCONF 0.9; OPTION 1**

Itemset Time ■ Rule Time ■

| | 1000 | 500 | 100 | 50 | 30 |
|---|---|---|---|---|---|
| Itemset Time | 1.839 | 5.576 | 35.696 | 56.523 | 88.04 |
| Rule Time | 0.005 | 0.009 | 0.095 | 0.493 | 162.592 |



**MINCONF 0.9; OPTION 2**

Itemset Time ■ Rule Time ■

| | 1000 | 500 | 100 | 50 | 30 |
|---|---|---|---|---|---|
| Itemset Time | 1.88 | 5.496 | 36.282 | 56.076 | 82.428 |
| Rule Time | 0.004 | 0.009 | 0.11 | 0.626 | 163.878 |

**MINCONF 0.9; OPTION 3**

■ Itemset Time  ■ Rule Time

| | 1000 | 500 | 100 | 50 | 30 |
|---|---|---|---|---|---|
| Itemset Time | 1.952 | 5.531 | 36.511 | 56.789 | 85.264 |
| Rule Time | 0.004 | 0.009 | 0.102 | 0.631 | 159.538 |



**MINCONF 0.8; OPTION 1**

■ Itemset Time  ■ Rule Time

| | 1000 | 500 | 100 | 50 | 30 |
|---|---|---|---|---|---|
| Itemset Time | 1.928 | 5.834 | 34.985 | 56.573 | 70.854 |
| Rule Time | 0.005 | 0.01 | 0.118 | 0.902 | 145.21 |

## MINCONF 0.8; OPTION 2

■ Itemset Time  ■ Rule Time

| Value | Itemset Time | Rule Time |
|-------|--------------|-----------|
| 1000  | 1.787        | 0.005     |
| 500   | 5.695        | 0.009     |
| 100   | 35.12        | 0.129     |
| 50    | 53.739       | 0.868     |
| 30    | 66.479       | 143.725   |

## MINCONF 0.8; OPTION 3

■ Itemset Time  ■ Rule Time

| Value | Itemset Time | Rule Time |
|-------|--------------|-----------|
| 1000  | 1.885        | 0.005     |
| 500   | 5.458        | 0.018     |
| 100   | 35.059       | 0.103     |
| 50    | 56.901       | 0.898     |
| 30    | 72.117       | 143.676   |

# RULES AND PATTERNS GENERATED FOR MINCONF 0.95

■ Itemsets Mined   ■ Rules Generated

| | 1000 | 500 | 100 | 50 | 30 | 20 |
|---|---|---|---|---|---|---|
| Itemsets Mined | 637 | 2311 | 39979 | 166565 | 1433744 | 17565168 |
| Rules Generated | 0 | 1 | 204 | 83918 | 16035375 | 1075765687 |

# RULES AND PATTERNS GENERATED FOR MINCONF 0.9

■ Itemsets Mined   ■ Rules Generated

| | 1000 | 500 | 100 | 50 | 30 |
|---|---|---|---|---|---|
| Itemsets Mined | 637 | 2311 | 39979 | 166565 | 1433744 |
| Rules Generated | 3 | 15 | 809 | 204503 | 35587019 |

## RULES AND PATTERNS GENERATED AT MINCONF 0.8

■ Itemsets Mined   ■ Rules Generated

| | 1000 | 500 | 100 | 50 | 30 |
|---|---|---|---|---|---|
| Itemsets Mined | 637 | 2311 | 39979 | 166565 | 1433744 |
| Rules Generated | 37 | 170 | 5590 | 425785 | 79951807 |

**Data Analysis:**

Though the data are not entirely complete, conclusions can still be drawn about the data. For high minsup values, there doesn't appear to me much correlation between the option chosen and the time taken. However, a pattern starts to emerge when minsup gets smaller (though it's not incredibly pronounced). In particular, option 2 seems to perform the best when minsup gets small. While we were told to look for this, it honestly wasn't obvious to me why this is happening. If I had implemented the algorithm via an FP tree, then this would make a lot of sense; the tree would be a lot more compact with option 2. However, in my algorithm there is no tree, any speed increase can only occur from the change in the way my algorithm expands the data.

My hypothesis is that if an item is common, it's likely to have a lot of supersets. When the lexicographic ordering is done by item frequency, then all of its potential supersets are expanded from one node in the enumeration tree. That means we can quickly determine all its candidates, and thus all of its immediate supersets. Alternatively, if the lexicographic ordering is not done based on item frequency, then a frequent item may reappear at many places in the enumeration tree. Nodes that would otherwise have no children, will end up having a frequent itemset as a child. Furthermore, to find the supersets of a frequent item, we'll have to find all instances of the itemset in the enumeration tree and count its frequency. In essence, **we're expanding this same item multiple times instead of once, leading to a less efficient algorithm.**