

Final Report

1) Data Structures Used:

Reading: We used **List** to store the transactions read from the input file.

Preparing Dataset: We used **Data Frame** to store the transactions in the in the form of rows and items in the form of columns (representing in the form of asymmetric binary attributes) which we can later use to find the support in the following way. Sum of a single column will give the support of that item and sum of multiple columns will give the support of itemset.

Ex: If we want to find out the support of an itemset $\{ '1', '2', '3' \}$ then we can select the rows from $'1', '2', '3'$ columns such that all the columns have 1 value. And finally count the rows so we get the support of $\{ '1', '2', '3' \}$.

Generating Frequent 1-itemsets: We used a **set** to store each item and then used **data frame** to find its support by summing up the column. We have the item and its support eliminate the itemset not satisfying minimum support. And now we used a **dictionary** to store them as key and value pair.

Candidate Generation: In order to generate kth candidates, we used $F_{k-1} \times F_{k-1}$. In a loop We used **set** operation union and checked its length is equal to k if satisfied then we added it to the **List** and then performed pruning on it. But to store them in dictionary the items should be immutable, so we use **Frozenset()** function.

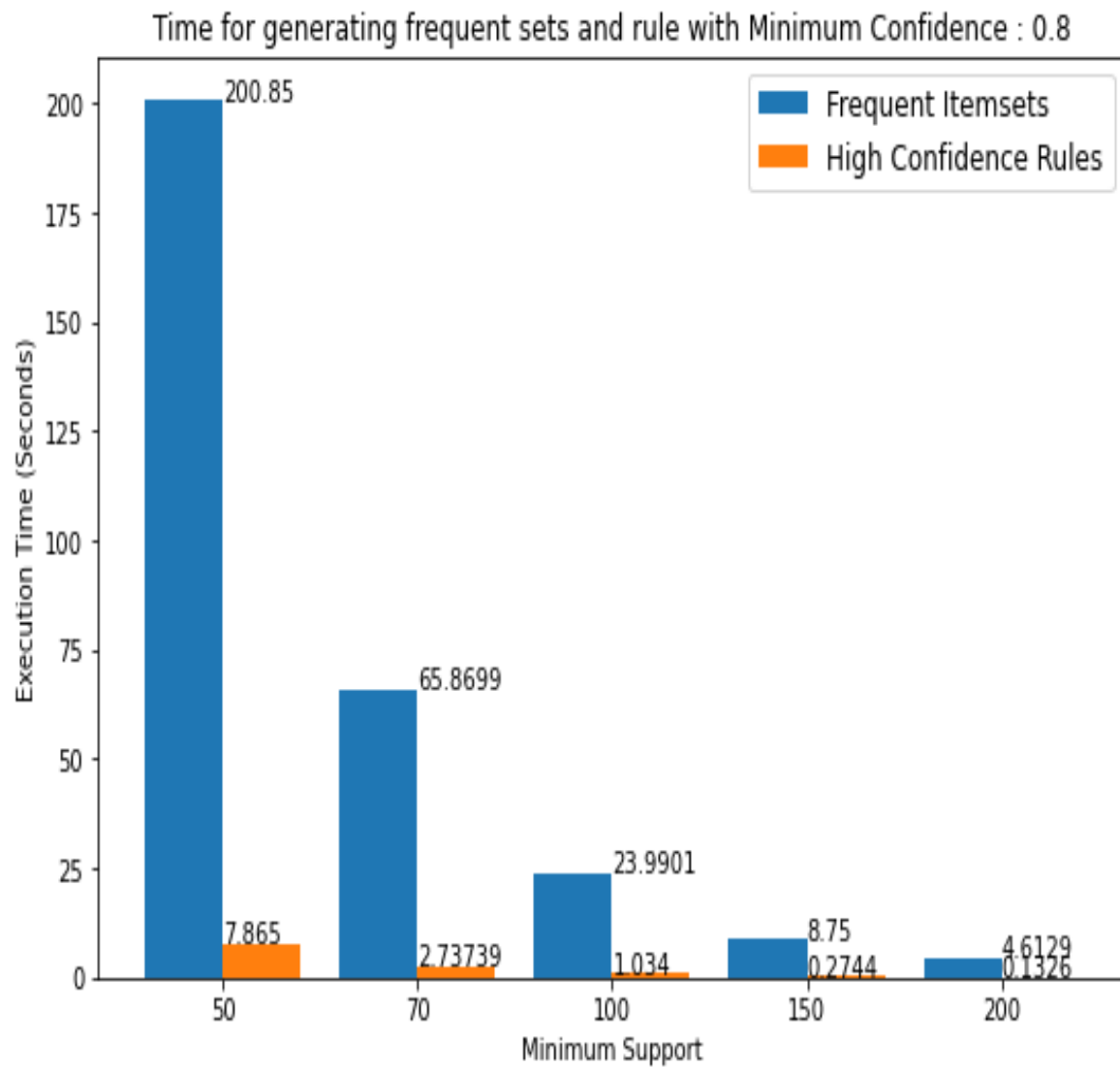
Pruning: After candidate generation we checked immediate subsets of each itemset in candidate generation from **List**. We generated k-1 all possible **subsets** of itemset using **itertools combination** and stored in a **List**. Now we checked every combination is frequent or not. If atleast one of them, is not frequent then we eliminate itemset.

Support Counting: After Pruning we get a **List** of **sets** so we now calculated the support of each itemset from the **Data Frame** and store them as key-value pair in a **dictionary**.

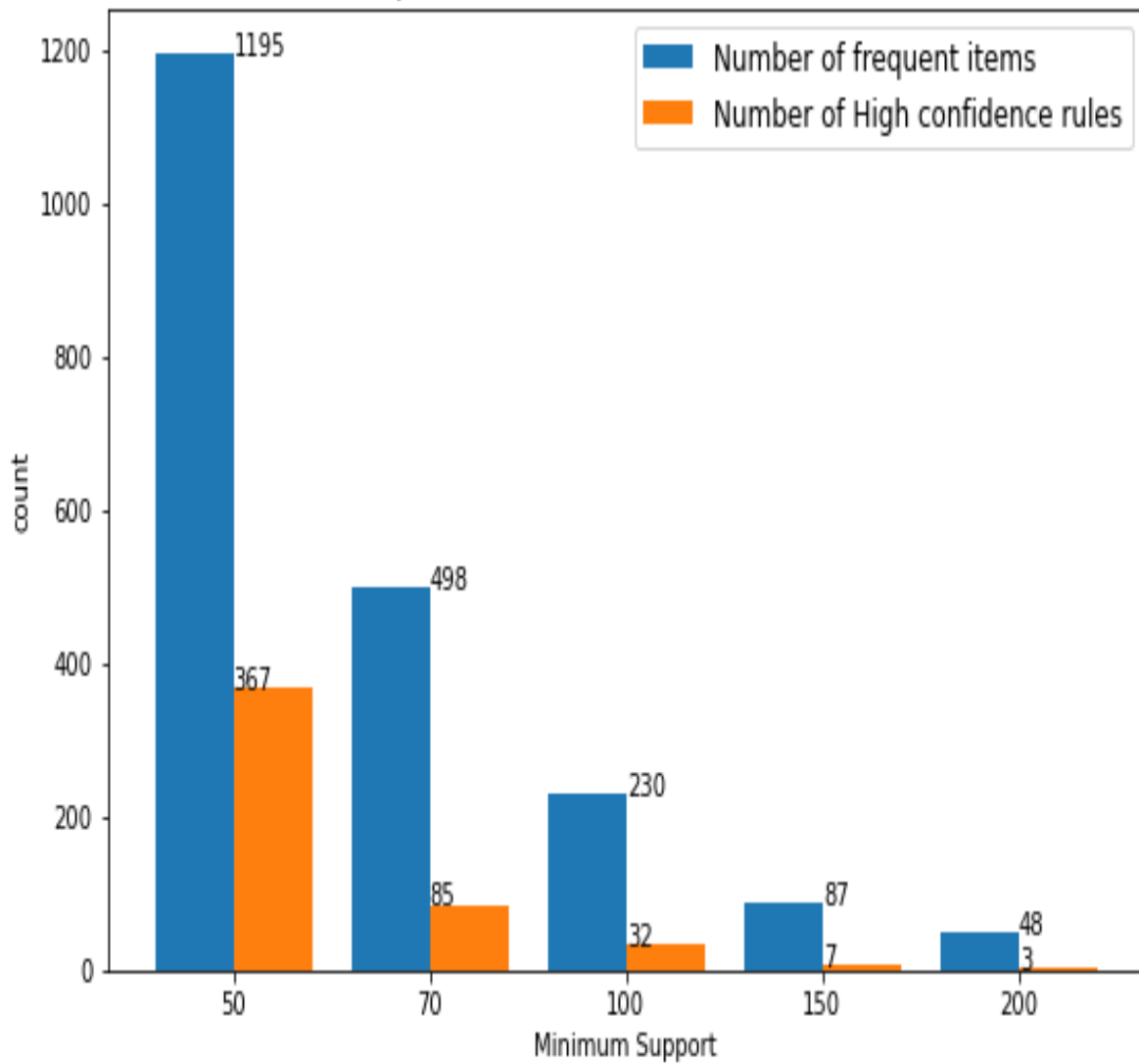
Candidate Elimination: Now we have keys(itemset) and values(support) in a **dictionary**. We check each key's support with minimum support if not satisfied then we eliminate them. Those which satisfy can be stored in a **List**. All such lists are stored in a **final List**.

Rule Generation: Now we have the **final List** consisting of all the frequent itemset. Now we calculate each itemset **combinations List** from 2 to n-1. Calculate each combinations confidence using support from **Data Frame**. And check the minimum confidence if it satisfies then add to output.

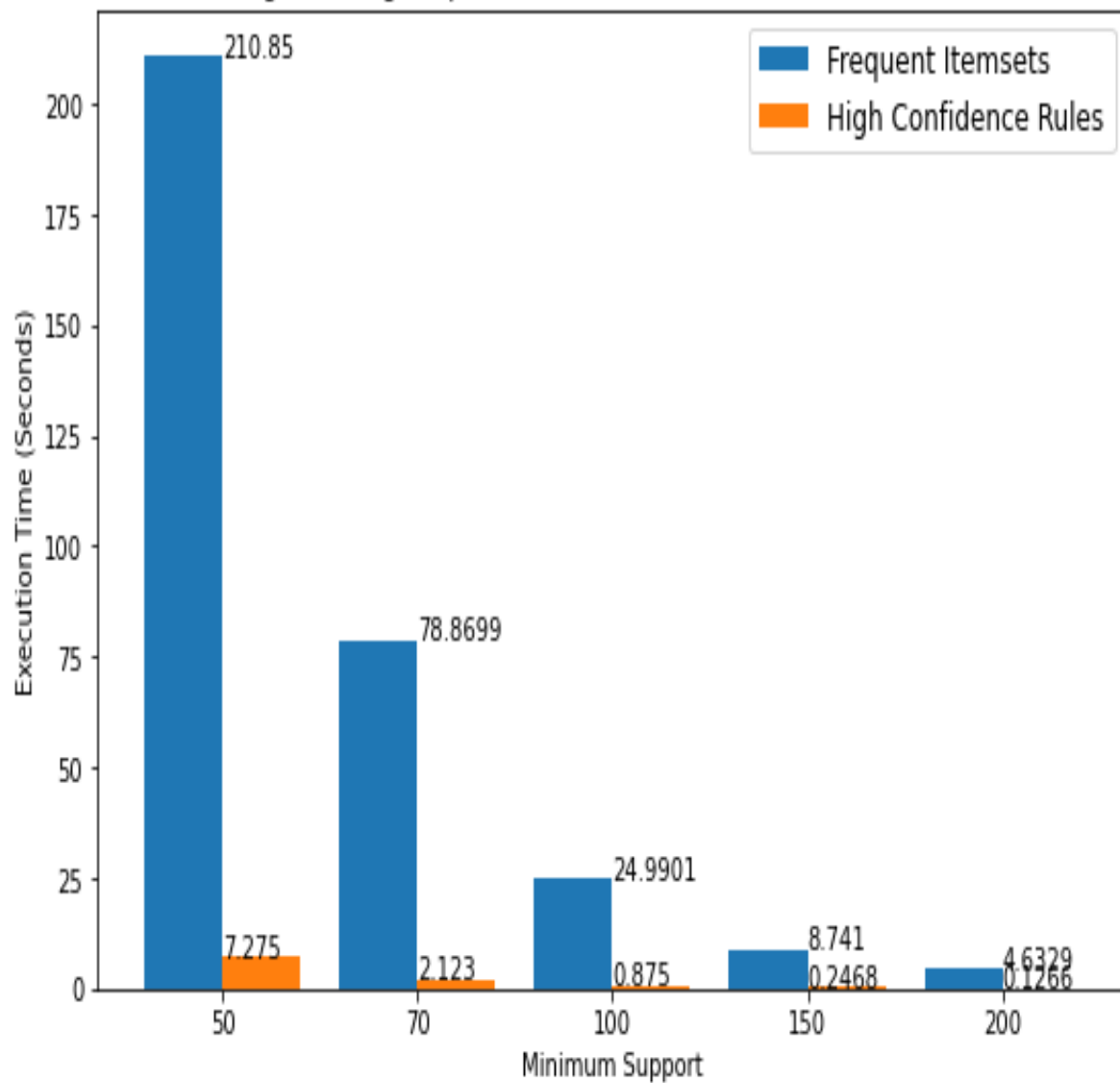
2)



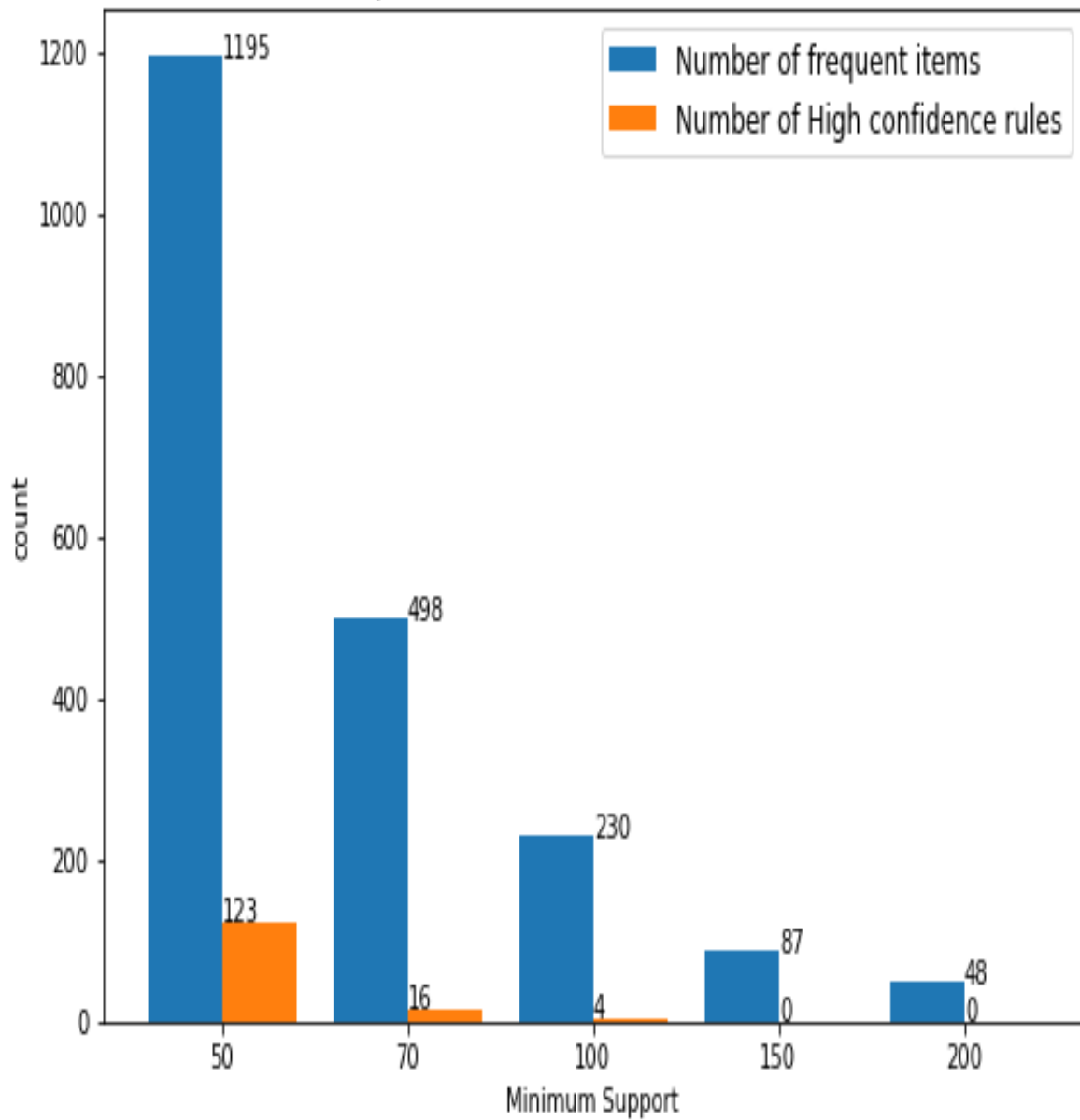
Number of frequent sets and rules with Minimum Confidence : 0.8



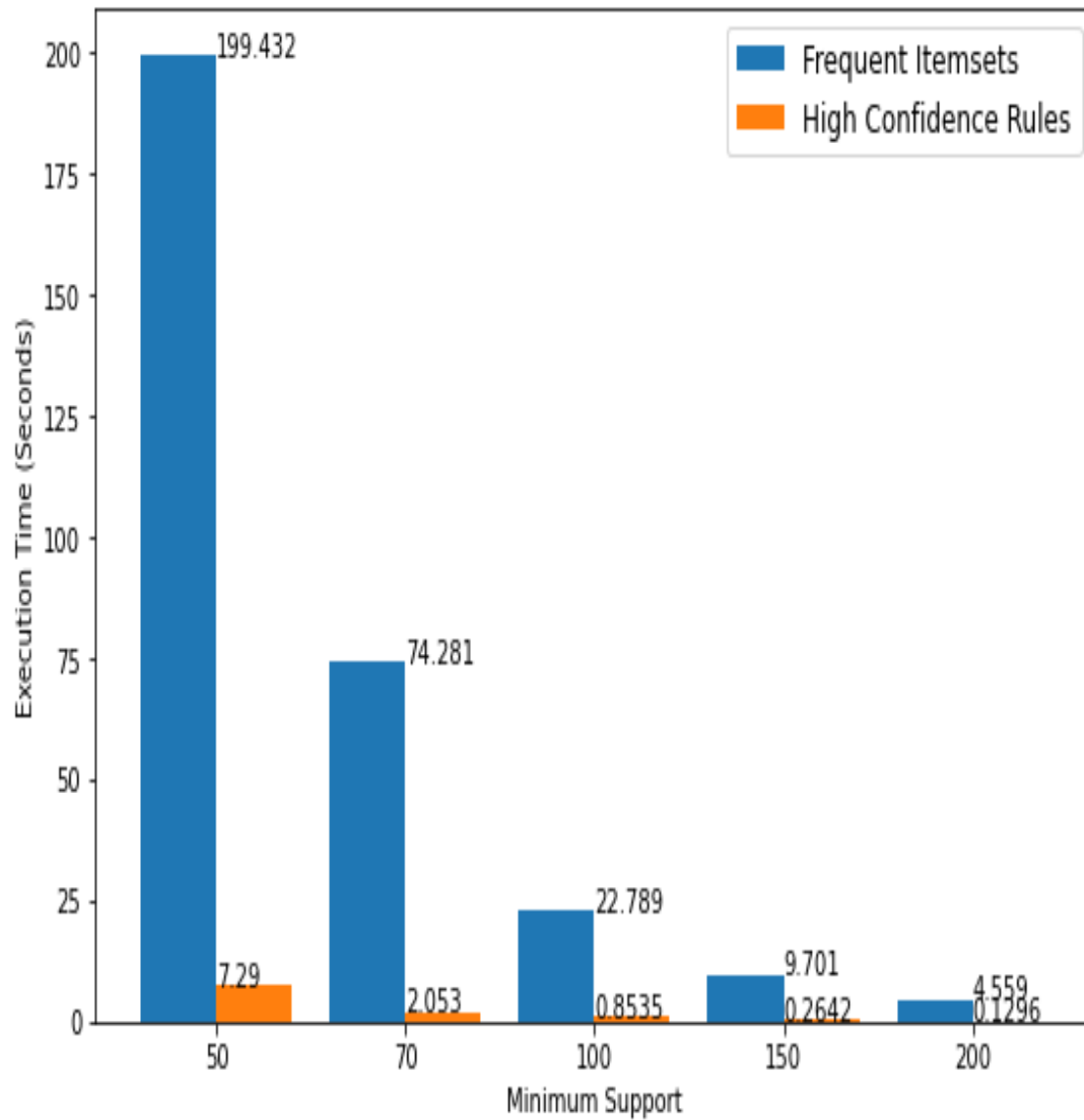
Time for generating frequent sets and rule with Minimum Confidence : 0.9



Number of frequent sets and rules with Minimum Confidence : 0.9



Time for generating frequent sets and rule with Minimum Confidence : 0.95



Number of frequent sets and rules with Minimum Confidence : 0.95

