

Apriori implementation using Matlab

Luting Chen 50133507, Yuze Liu 50207903, Vicky Zheng 50037709

Introduction

The apriori algorithm is used to determine association rules. In this homework, we implemented the apriori algorithm by using Matlab. The dataset we used the algorithm on is gene data. The applications of this are straight forward: with association rules we can see how genes relate to each other and what the relation between genes can tell us about the data.

Data Preprocessing

The data provided is a .txt file that contains a 100 by 102 matrix. There are 100 samples each corresponding to the genes of a patient with a disease. The diseases are: ALL, AML, Breast Cancer and Colon Cancer. In order to optimize our program, we converted each of these values into integers: 1,2,3,4 respectively. The data set also gives us a sequence of 100 genes labeled as either up or down. Each gene at index i is treated as 1 if it is up and 0 if it is down. We did this to optimize our program by having fewer string comparisons since integer comparisons are much faster.

Implementation

For this homework, we choose to use Matlab because our input data is a matrix. Matlab has helpful functions for matrix manipulation like determining whether a set of numbers is in a row

of the matrix. It also has helpful functions such as the union of two sets which helped us do the pruning necessary for an efficient apriori algorithm. In our implementation, we enumerate through all possible tuples where a tuple can be comprised of genes or genes and diseases. We start at the shortest length of tuples, and work our way up. We based our implementation off the diagram provided in class which is shown in Fig. 1.

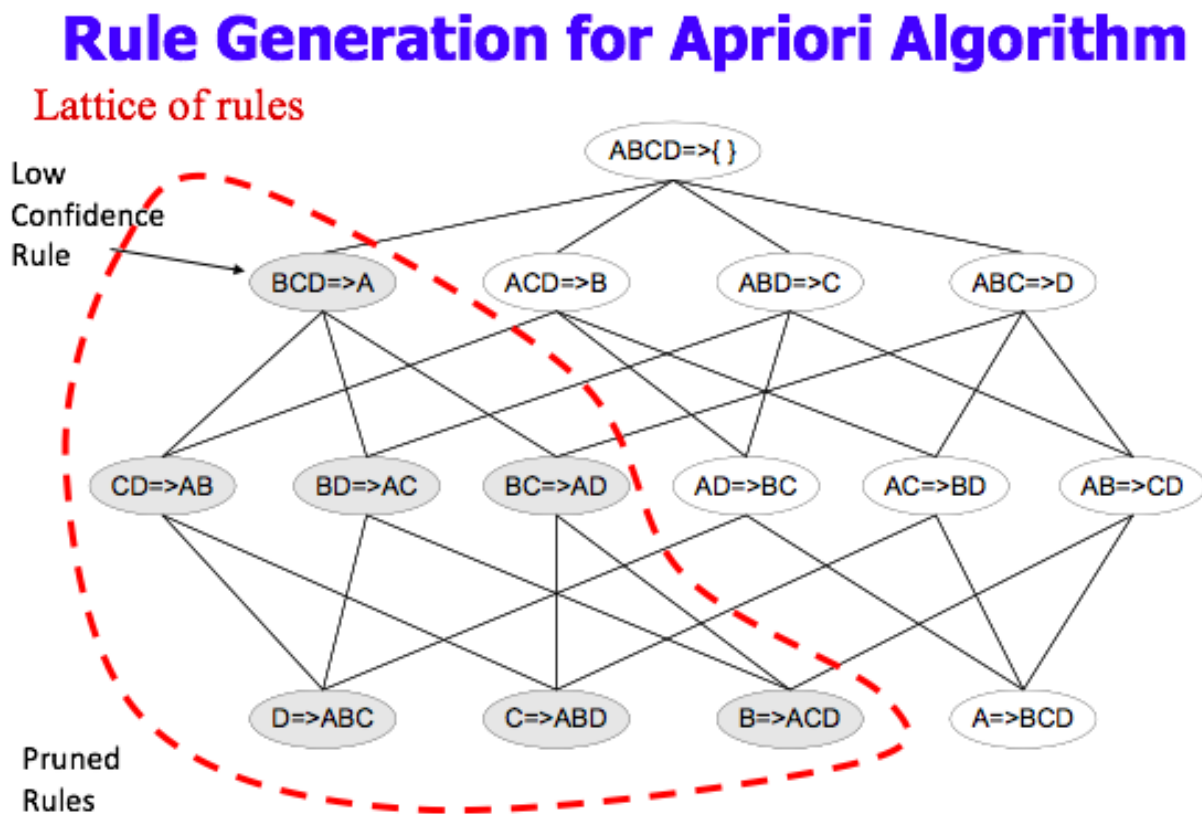


Figure 1: Apriori Pruning Search Space

Although our algorithm is $O(n!)$ where n is the number of attributes, this pruning significantly reduces our runtime. This is because we exploit the fact that a superset of an infrequent set must also be infrequent. In Fig. 1,

$$BCD \Rightarrow A$$

was infrequent therefore all of its supersets are also infrequent so we don't have to explore it.

Results from Part 1

With a support level of 30%, we obtain the values:

196 of length 1 frequent item sets.

5340 of length 2 frequent item sets.

5287 of length 3 frequent item sets.

1518 of length 4 frequent item sets.

438 of length 5 frequent item sets.

88 of length 6 frequent item sets.

11 of length 7 frequent item sets.

1 of length 8 frequent item sets.

With a support level of 40%, we obtain the values:

167 of length 1 frequent item sets.

753 of length 2 frequent item sets.

149 of length 3 frequent item sets.

7 of length 4 frequent item sets.

1 of length 5 frequent item sets.

With a support level of 50%, we obtain the values:

109 of length 1 frequent item sets.

63 of length 2 frequent item sets.

2 of length 3 frequent item sets.

With a support level of 60%, we obtain the values:

31 of length 1 frequent item sets.

2 of length 2 frequent item sets.

With a support level of 70%, we obtain the values:

7 of length 1 frequent item sets.

Results from Part 2

The results from Part 2 were obtained by combining our results from Part 1 and combining it with the counting method described in piazza post 113.

1. RULE HAS ANY OF G6_UP

Result:

2. RULE HAS 1 OF G1_UP

Result:

3. RULE HAS 1 OF (G1_UP, G10_DOWN)

Result:

4. BODY HAS ANY OF G6_UP

Result:

5. BODY HAS NONE OF G72_UP

Result:

6. BODY HAS 1 OF (G1_UP, G10_DOWN)

Result:

7. HEAD HAS ANY OF G6_UP

Result:

8. HEAD HAS NONE OF (G1_UP, G6_UP)

Result:

9. HEAD HAS 1 OF (G6_UP, G8_UP)

Result:

10. RULE HAS 1 OF (G1_UP, G6_UP, G72_UP)

Result:

11. RULE HAS ANY OF (G1_UP, G6_UP, G72_UP)

Result:

12. SIZE OF RULE $i = 3$

Result:

13. SIZE OF BODY $i=2$

Result:

14. SIZE OF HEAD $i=2$

Result:

15. BODY HAS ANY OF G1_UP AND HEAD HAS 1 OF G59_UP

Result:

16. BODY HAS ANY OF G1_UP OR HEAD HAS 1 OF G6_UP

Result:

17. BODY HAS 1 OF G1_UP OR HEAD HAS 2 OF G6_UP

Result:

18. HEAD HAS 1 OF G1_UP AND BODY HAS 0 OF DISEASE

Result:

19. HEAD HAS 1 OF DISEASE OR RULE HAS 1 OF (G72_UP, G96_DOWN)

Result:

20. BODY HAS 1 of (G59_UP, G96_DOWN) AND SIZE OF RULE $i=3$

Result:

Conclusion

In conclusion, BioGalaxy can have reasonably timed queries since having multiple fact tables allows for fewer joins. Our schema also supports many to many relationships, temporal data, and hierarchal querying.

References and Notes

[1] Jing Gao *Association1 Slides*.