

August 13, 2023

## 0.1 Summary on Association Analysis:

### **Preliminaries:**

Association analysis is pivotal in identifying patterns in data, especially in revealing how items relate or associate with each other. Its application spans a wide range, with market-basket analysis standing out as a classic example, primarily detecting items that are frequently purchased together.

### **Frequent Itemset Generation:**

Central to association analysis is the Apriori Principle. This principle posits that if an itemset is frequent, all of its subsets must also be frequent. Conversely, if an itemset is found to be infrequent, all of its supersets will be infrequent as well. This principle simplifies the process of searching for frequent itemsets by systematically narrowing down candidate sets. Techniques such as Candidate Generation and Pruning harness the Apriori principle to significantly reduce the number of candidates considered. The frequency of appearance of an itemset within a dataset, termed Support Counting, is crucial in determining its significance. However, it's essential to recognize the computational complexity associated with these processes. Despite the efficiencies introduced by the Apriori principle, mining can still be resource-intensive, necessitating the use of efficient data structures and algorithms.

### **Rule Generation:**

After generating frequent itemsets, the next step is Rule Generation. Confidence-Based Pruning emerges as a critical aspect here. It employs the confidence metric - a measure of how reliable an inferred rule is, to refine and prune rule sets further. One can look at Congressional Voting Records as a case study. This example underscores how association rules can unravel hidden relationships between seemingly unrelated items.

### **Compact Representation of Frequent Itemsets:**

Considering the vast amounts of data involved, representing frequent itemsets compactly becomes imperative. Maximal Frequent Itemsets refer to those frequent itemsets that lack any frequent supersets. On the other hand, Closed Itemsets are those frequent itemsets for which no superset exists with the same support level.

### **Alternative Methods for Generating Frequent Itemsets:**

While the Apriori algorithm is foundational, there exist alternative methods. These methods aim to mine frequent itemsets without the exhaustive process of candidate generation, promising enhanced efficiency.

### **FP-Growth Algorithm (Advanced):**

Venturing into advanced territory, the FP-Growth Algorithm offers a more efficient approach. Its crux lies in the FP-Tree Representation - a compact data structure capturing frequent items and

their inter-relationships. This algorithm allows for the extraction of frequent itemsets without the need for exhaustive searches.

**Evaluation of Association Patterns:**

Any association pattern needs rigorous evaluation. Objective Measures of Interestingness, such as support, confidence, and lift, serve as criteria to gauge the relevance of an association rule. As patterns get intricate, there's a push towards measures that evaluate beyond mere pairs of binary variables. A peculiar phenomenon to be aware of is Simpson's Paradox, where a trend evident in various data groups disappears when these groups merge.

**Effect of Skewed Support Distribution:**

In data landscapes where certain items significantly dominate in frequency (either by being too common or too rare), there's potential for meaningful associations to be overshadowed. Analytical techniques must be geared to correct for this skew to prevent conclusions that might mislead.