

Association Rule Mining

Basics, Apriori

Huiping Cao

Motivation

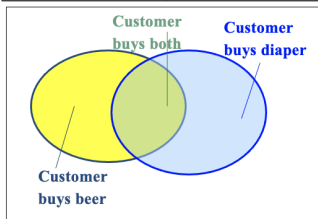
- Market basket analysis
 - Products that are often **bought together** by a customer
- Applications
 - Design store layouts → combined sale
 - Discount for package sale

What Is Frequent pattern?

- **Frequent pattern**: a pattern (a set of items, subsequences, substructures, etc.) that occurs frequently in a data set
- First proposed by Agrawal, Imielinski, and Swami [AIS93] in the context of **frequent itemsets** and **association rule mining**.

Basic Concepts: Frequent Itemsets and Association Rules

Transaction-id	Items Bought
10	A, B, D
20	A, C, D
30	A, D, E
40	B, E, F
50	B, C, D, E, F



- Itemset $X = \{x_1, \dots, x_k\}$
- Find all the rules $X \rightarrow Y$ with minimum **support** and **confidence**
 - **support**, s , probability that a transaction contains $X \cup Y$
 - **confidence**, c , conditional probability that a transaction having X also contains Y

Basic concepts

- $\text{Support}(X \rightarrow Y) = P(X \cup Y) = \frac{\text{count}(X \cup Y)}{\text{total transaction count}}$

- Relative support

- Absolute support: occurrence frequency vfill

- $X \cup Y$: means the transaction contains all the items in X and Y ; does NOT mean that it contains EITHER X or Y

- K -itemset

- Frequent k -itemset L_k

- $\text{Confidence}(X \rightarrow Y) = P(Y|X)$

- $= \frac{\text{support_count}(X \cup Y)}{\text{support_count}(X)}$

Basic Concepts: Frequent Itemsets and Association Rules

Transaction-id	Items Bought
10	A, B, D
20	A, C, D
30	A, D, E
40	B, E, F
50	B, C, D, E, F

- Let $\text{supmin} = 50\%$, $\text{confmin} = 50\%$
 Freq. Itemset.: $\{A:3, B:3, D:4, E:3, AD:3\}$
 Association rules:
 $A \rightarrow D$ (60%, 100%)
 $D \rightarrow A$ (60%, 75%)

$\{A\}: 3$

$\{B\}: 3$

$\{D\}: 4$

$\{E\}: 3$

$\{AD\}: 3$

$A \rightarrow D \iff \{A\} \rightarrow \{D\}$

$\text{count}(A,D): 3$

$\text{total count}: 5$

$\text{Support} = 3/5$

$\text{count}(A,D) = 3$

$\text{count}(A) = 3$

$\text{conf.} = 3/3 = 1.0$

$D \rightarrow A$

$\text{count}(A,D) = 3$

$\text{count}(D) = 4$

$\text{conf.} = 3/4$

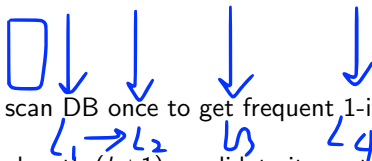
Scalable Methods for Mining Frequent Patterns

- The downward closure property of frequent patterns
- Any subset of a frequent itemset must be frequent
- If {beer, diaper, nuts} is frequent, so is {beer, diaper}
i.e., every transaction having {beer, diaper, nuts} also contains {beer, diaper}
- Scalable mining methods: Three major approaches
 - Apriori (Agrawal & Srikant VLDB'94)
 - Freq. pattern growth (FPgrowth Han, Pei & Yin SIGMOD'00)
 - Vertical data format approach (Charm Zaki & Hsiao SDM'02)

Apriori: A Candidate Generation-and-Test Approach

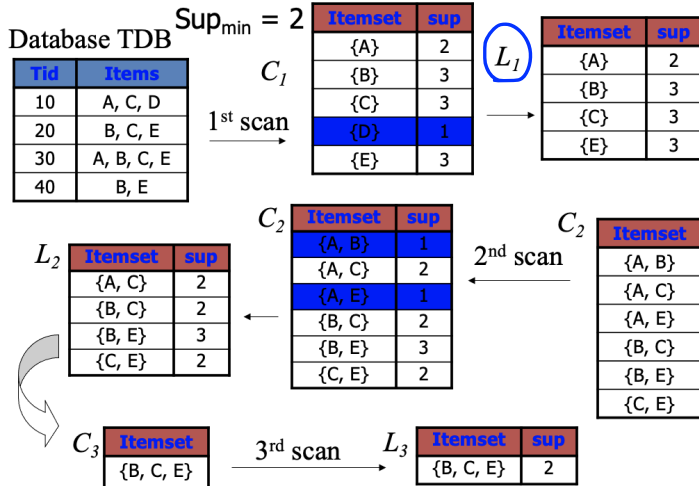
- **Apriori pruning principle**: If there is **any** itemset which is infrequent, its superset should not be generated/tested!
(Agrawal & Srikant VLDB'94, Mannila, et al. KDD' 94)

- **Method**:



- Initially, scan DB once to get frequent 1-itemset
- **Generate** length $(k+1)$ candidate itemsets from length k frequent itemsets
- **Test** the candidates against DB
- Terminate when no frequent or candidate set can be generated

The Apriori Algorithm An Example



Huiping Cao, Association Rule Mining, Slide 11/24

Important Details of Apriori

■ How to generate candidates?

- Step 1: self-joining L_k
- Step 2: pruning

$L_k \rightarrow C_{k+1}$

■ How to count supports of candidates?


■ Example of Candidate-generation

- $L_3 = \{abc, abd, acd, ace, bcd\}$
- Self-joining: $L_3 * L_3$
 $abcd$ from abc and abd
 $acde$ from acd and ace
- Pruning:
 $acde$ is removed because ade is not in L_3
- $C_4 = \{abcd\}$

How to Generate Candidates?

- Suppose the items in L_{k-1} are listed in an order
- Step 1: self-joining L_{k-1} **insert into** C_k
 $\text{select } p.\text{item}_1, p.\text{item}_2, \dots, p.\text{item}_{k-1}, q.\text{item}_{k-1}$
 $\text{from } L_{k-1} \text{ } p, L_{k-1} \text{ } q$
 $\text{where } p.\text{item}_1 = q.\text{item}_1, \dots, p.\text{item}_{k-2} = q.\text{item}_{k-2},$
 $p.\text{item}_{k-1} < q.\text{item}_{k-1}$
- Step 2: pruning
for all itemsets c in C_k do
 - **for all** $(k-1)$ -subsets s of c do
 - If $(s \text{ is not in } L_{k-1})$ then delete c from C_k

How to Count Supports of Candidates?

- Why counting supports of candidates a problem?
 - The total number of candidates can be very huge
 - One transaction may contain many candidates
- Naive Method:
 - Compare each transaction against every candidate itemset
 - Expensive!

- Hash-tree based method:

- Candidate itemsets are stored in a **hash tree**
- **Leaf node** of hash-tree contains a list of candidate itemsets and counts
- **Interior node** contains a hash table
- **Subset function**: finds all the candidates contained in a transaction

Subset function $x \% 3$

```
graph TD; A[ ] --- B[1,4,7]; A --- C[2,5,8]; A --- D[3,6,9];
```



Complexity

- Support Threshold
 - Maximum size of frequent itemsets
- Number of items (dimensionality)
- Number of transactions
- Average transaction width
 - More frequent itemsets

$L_k \cup 21$

Challenges of Frequent Pattern Mining

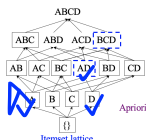
- Challenges
 - Multiple scans of transaction database
 - Huge number of candidates
 - Tedious workload of support counting for candidates
- Improving Apriori: general ideas
 - Reduce passes of transaction database scans
 - Shrink number of candidates
 - Facilitate support counting of candidates

Partition: Scan Database Only Twice

- Any itemset that is potentially frequent in DB must be frequent in at least one of the partitions of DB
 - Scan 1: partition database and find local frequent patterns
 - Scan 2: consolidate global frequent patterns
- A. Savasere, E. Omiecinski, and S. Navathe. An efficient algorithm for mining association in large databases. In VLDB'95

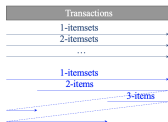
- A **partition** p of the database refers to any subset of the transactions.
- Any two different partitions are **non-overlapping**.
- **Local support of a partition**: the fraction of transactions containing that item in a partition

DIC: Reduce Number of Scans



Apriori

- Once both A and D are determined frequent, the counting of AD begins
- Once all length-2 subsets of BCD are determined frequent, the counting of BCD begins



Itemset lattice
S. Brin, R. Motwani, J. Ullman,
and S. Y. Chong. Dynamic itemset
counting and implication rules for
market basket data. In
SIGMOD '97

DIC

Bottleneck of Frequent-pattern Mining

- Multiple database scans are costly
- Mining long patterns needs many passes of scanning and generates lots of candidates
 - To find frequent itemset i_1, i_2, \dots, i_{100}
of scans: 100 # of candidates: $2^{100} - 1 = 1.27 * 10^{30}$
- **Bottleneck:** candidate-generation-and-test
- Can we avoid candidate generation?

- Chapter 5: Introduction to Data Mining (2nd Edition) by Pang-Ning Tan, Michael Steinbach, Anuj Karpatne, and Vipin Kumar
- Implementation of Apriori algorithm with Python 2.7 and 3.3 - 3.5: <https://pypi.org/project/apyori/>

12