



Fast algorithm for mining association rules

Presenter: Linda
Discusser :Massih
Instructor : Rachel Pottinger
University of British Columbia



Overview

- Introduction of data Ming (application of mining)
- Association rule mining
- Apriori algorithm
- Comparison of different algorithms
- Performance comparison



Introduction of data mining

• What is of data mining ?

Data mining (DM) is the process of automatically searching large volumes of data for patterns such as association rules.

• A Real-World Example:

a supermarket chain who, through analysis of transactions over a long period of time, found that beer and diapers were often bought together



Introduction of data mining

What are the uses of data mining ?

- *Market segmentation* - Identify common characteristics of customers who buy the same products from your company.
- *Customer churn* - Predict which customers are likely to leave and go to a competitor.
- *Fraud detection* - Identify which transactions are likely to be fraudulent.
- *Direct marketing* - Identify which prospects should be included in a mailing list to obtain the highest response rate.
- *Interactive marketing* - Predict what each individual accessing a Web site is most likely interested in seeing.
- *Market basket analysis* - Understand what products or services are commonly purchased together; e.g., beer and diapers.
- *Trend analysis* - Reveal the difference between a typical customer this month and last.



Introduction of data mining

• How does data mining works ?

Data mining software analyzes **relationships and patterns** in stored transaction data based on open-ended user queries.

• What types of relationships are sought?

- **Classes:** Stored data is used to locate data in predetermined groups.
example: a restaurant chain could mine customer purchase data to determine when customers visit and what they typically order. This information could be used to increase traffic by having daily specials.
 - **Clusters:** Data items are grouped according to logical relationships or consumer preferences.
example, data can be mined to identify market segments or consumer affinities.
- (cont.)



Relationship Types (cont.)

- **Sequential patterns:** Data is mined to anticipate behavior patterns and trends.
example, an outdoor equipment retailer could predict the likelihood of a backpack being purchased based on a consumer's purchase of sleeping bags and hiking shoes.
- **Associations:** Data can be mined to identify associations.
The beer-diaper example is an example of associative mining.

Association rules mining



Itemset

- A set of itemsets is referenced to as itemset
- An itemset containing k items is called k-itemset

Associate rule

- $X \rightarrow Y$ (X and Y are disjoint sets of item set)
- Support of an association rule is the percentage of the relevant data transaction for the which the rule is true.

$$\text{support}(A \rightarrow B) = \frac{\text{IF } A \rightarrow B \text{ \# tuples containig both A and B}}{\text{total_#_of_tuples}}$$

- Confidence of an association is the measure of certainty associated with each discovered pattern.

$$\text{confidence}(A \rightarrow B) = \frac{\text{IF } A \rightarrow B \text{ \# tuples containig both A and B}}{\text{total_#_of_tuples}}$$

Association rule mining



- Find all **frequent** itemsets
- Generate **strong association rules** from the frequent itemsets
- The **problem** of association rule mining is: efficiently find all rules with support > Minsup , confidence > Minconf
- Example

Example



- Itemset $acm = \{a, c, m\}$
- Support of itemset $\text{Sup}(acm) = 3$
- Given $\text{min_sup} = 3$, acm is a frequent patten
- Frequent pattern mining find all frequent patterns in a data base

Transaction database TDB

TID	Items bought
100	F, a, c, d, g, l, m, p
200	A, b, c, f, l, m, o
300	B, f, h, j, o
400	B, c, k, s, p
500	A, f, c, e, l, p, m, n

Problem decomposition



- Discover large itemsets
- Use the largest itemsets to generate the desired rules

Algorithm Apriori



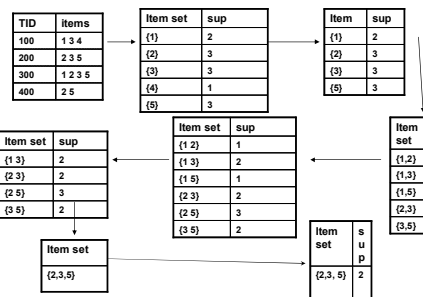
```

F1 = {frequent 1-item sets};
k = 2;
while( Fk-1 is not empty ) {
    Ck = Apriori_generate( Fk-1 );
    for all transactions t in T {
        Subset( Ck, t );
    }
    Fk = { c in Ck s.t. c.count >= minimum_support };
}
Answer = union of all sets Fk;
    
```

Find large items

Find candidate pairs, count them
→ large triplets of items

The Apriori algorithm-- Example



AprioriTid

- Uses the database only once.
- Builds a storage set C^k
 - Members has the form $\langle TID, \{X_k\} \rangle$
 - X_k are potentially large k-items in transaction TID.
 - For $k=1$, C^1 is the database.
- Uses C^k in pass $k+1$.

Example Of AprioriTid

Database		C^1		L_1	
TID	Items	TID	Set-of-Items	Itemset	Support
100	1 3 4	100	{1}, {3}, {4}	{1}	2
200	2 3 5	200	{2}, {3}, {5}	{3}	3
300	1 2 3 5	300	{1}, {2}, {3}, {5}	{3}	3
400	2 5	400	{2}, {5}	{5}	3

C_2		C^2		L_2	
Itemset	TID	Set-of-Items	Itemset	Support	
{1, 2}	100	{1, 3}	{1, 3}	{1, 3}	2
{1, 3}	200	{2, 3}, {3, 5}	{2, 3}	{2, 3}	2
{1, 5}	300	{1, 3}, {1, 5}	{2, 5}	{2, 5}	3
{2, 3}	300	{2, 3}, {2, 5}, {3, 5}	{3, 5}	{3, 5}	2
{2, 5}	400	{2, 5}			

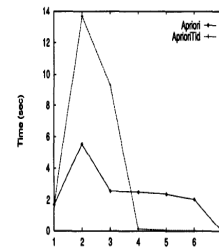
C_3		C^3		L_3	
Itemset	TID	Set-of-Items	Itemset	Support	
{2, 3, 5}	200	{2, 3, 5}			
	300	{2, 3, 5}			

Figure 3: Example

Other algorithms

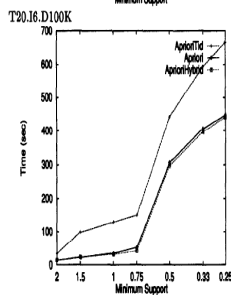
Apriori VS AprioriTid

- In the earlier passes, Apriori does better than AprioriTid.
- AprioriTid beats Apriori in later passes.



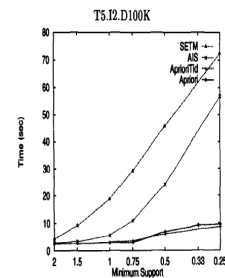
AprioriHybrid

- Uses Apriori in the initial passes and switches to AprioriTid when it expects that the set C_k at the end of the pass will fit in memory



Other algorithms

- AIS
- SETM
- Apriori and AprioriTid algorithms are much better than the AIS algorithm.
- AIS always did considerably better than SETM.



summary



- Association rules are an important tool in analyzing databases.
- We've seen an algorithm which finds all association rules in a database.
- The algorithm has better time results than previous algorithms.
- The algorithm maintains its performance for large databases.

Discussion



- What are challenges of data mining
- Scalability
- Dimensionality
- Complex and heterogeneous data
- Data quality
- Data ownership and distribution
- Streaming data