

Search and Data Mining: Techniques

Association Rules

Anya Yarygina

Boris Novikov

Introduction

- Basic concepts and road map
- Efficient and scalable frequent itemset mining methods
- Mining various kinds of association rules
- From association mining to correlation analysis
- Constraint-based association mining

Reference

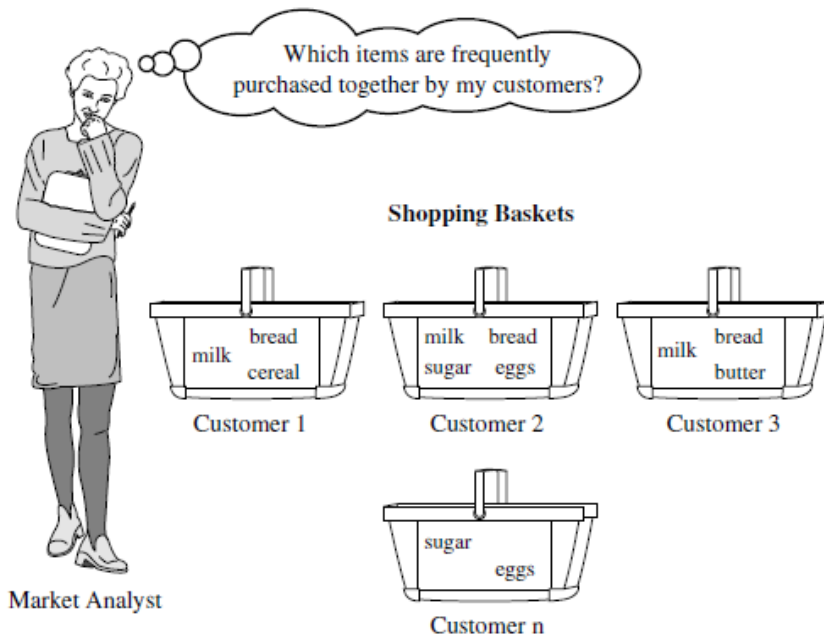
Jiawei Han, Micheline Kamber. Data Mining: Concepts and Techniques, 2nd Edition. The Morgan Kaufmann Series in Data Management Systems, 2006

Frequent patterns

- Frequent patterns
 - Frequent itemset
 - Frequent sequential pattern
 - Frequent structure pattern

Market basket analysis: a motivating example

- Analyze customer buying habits by finding associations between the different items that customers place in their “shopping baskets”



- Basket as boolean vector
- Buying patterns that reflect items that are frequently associated or purchased together
- Association rules
 - computer=>antivirus_sw [support=2%, confidence=60%]
- Measures of rule interestingness
 - Support (usefulness)
 - Confidence (certainty)
- Interestingness criteria
 - Minimum support threshold
 - Minimum confidence threshold

Frequent itemsets, closed itemsets, association rules

- $I = \{I_1, I_2, \dots, I_m\}$
- $D = \{T\}$
- $A \Rightarrow B$
 - $\text{Support}(A \Rightarrow B) = P(AB)$
 - $\text{Confidence}(A \Rightarrow B) = P(B|A)$
- Frequent itemset
- Strong rules
- Closed itemset
- Closed frequent itemset
- Maximal frequent itemset

Association rule mining

- Find all frequent itemsets
- Generate strong association rules from frequent itemsets

Frequent pattern mining: a road map

- Completeness of patterns to be mined
 - Complete set of frequent itemsets
 - Closed frequent itemsets
 - Maximal frequent itemsets
 - Constrained frequent itemsets
 - Approximate frequent itemsets
 - Near-match frequent itemsets
 - Top-k frequent itemsets
- Levels of abstraction involved in rule set
 - Single-level association rules
 - Multilevel association rules
- Number of data dimensions involved in the rule
 - Single-dimensional association rules
 - Multidimensional association rules
- Types of values handled in the rule
 - Boolean association rules
 - Quantitative association rules
- Kinds of rules to be mined
 - Association rules
 - Correlation rules
 - Strong gradient relationships
- Kinds of patterns to be mined
 - Frequent itemset mining
 - Sequential pattern mining
 - Structured pattern mining

Efficient and scalable frequent itemset mining methods

- Finding frequent itemsets
 - With candidate generation
 - Without candidate generation
 - From data in vertical format
 - Closed frequent itemsets
- Generating association rules

Apriori algorithm: finding frequent itemsets using candidate generation

- Level-wise search
- Apriori property
 - All nonempty subsets of a frequent itemset must also be frequent
- Construct L_k from L_{k-1}
 - Join step
 - Join L_{k-1} with itself to construct C_k
 - Members of L_{k-1} are joinable if their first $(k-2)$ items are in common
 - Prune step
 - Prune infrequent itemsets
 - Use apriori property
 - Any subset of candidate should be frequent
 - Subset testing using hash tree of all frequent itemsets

Apriori algorithm: pseudo-code

Input:

- D , a database of transactions;
- min_sup , the minimum support count threshold.

Output: L , frequent itemsets in D .

Method:

```
(1)  $L_1 = \text{find\_frequent\_1-itemsets}(D)$ ;  
(2) for ( $k = 2; L_{k-1} \neq \emptyset; k++$ ) {  
(3)    $C_k = \text{apriori\_gen}(L_{k-1})$ ;  
(4)   for each transaction  $t \in D$  { // scan  $D$  for counts  
(5)      $C_t = \text{subset}(C_k, t)$ ; // get the subsets of  $t$  that are candidates  
(6)     for each candidate  $c \in C_t$   
(7)        $c.\text{count}++$ ;  
(8)   }  
(9)    $L_k = \{c \in C_k \mid c.\text{count} \geq min\_sup\}$   
(10) }  
(11) return  $L = \cup_k L_k$ ;
```

procedure $\text{apriori_gen}(L_{k-1}:\text{frequent } (k-1)\text{-itemsets})$

```
(1) for each itemset  $l_1 \in L_{k-1}$   
(2)   for each itemset  $l_2 \in L_{k-1}$   
(3)     if ( $l_1[1] = l_2[1] \wedge l_1[2] = l_2[2] \wedge \dots \wedge l_1[k-2] = l_2[k-2] \wedge l_1[k-1] < l_2[k-1]$ ) then {  
(4)        $c = l_1 \bowtie l_2$ ; // join step: generate candidates  
(5)       if  $\text{has\_infrequent\_subset}(c, L_{k-1})$  then  
(6)         delete  $c$ ; // prune step: remove unfruitful candidate  
(7)       else add  $c$  to  $C_k$ ;  
(8)     }  
(9) return  $C_k$ ;
```

procedure $\text{has_infrequent_subset}(c:\text{candidate } k\text{-itemset};$

```
   $L_{k-1}:\text{frequent } (k-1)\text{-itemsets}$ ); // use prior knowledge  
(1) for each  $(k-1)$ -subset  $s$  of  $c$   
(2)   if  $s \notin L_{k-1}$  then  
(3)     return TRUE;  
(4) return FALSE;
```

Generating association rules from frequent itemsets

- Strong association rules
 - Minimum support
 - Minimum confidence
 - $\text{Confidence}(A \Rightarrow B) = P(B|A)$
- Generation of association rules
 - For each frequent itemset I generate all nonempty subsets of I
 - For every nonempty subset s of I output the rule “ $s \Rightarrow (I-s)$ ” if confidence is “high”

Improving efficiency of Apriori

- Hash-based technique
 - Hashing items into corresponding buckets
- Transaction reduction
 - Remove transactions from further consideration when it does not contain any frequent k-items
- Partitioning
 - Subdivide the transactions into nonoverlapping partitions
 - Find (local) frequent itemsets within partition
 - Calculate actual support to select global frequent itemsets
- Sampling
 - Mining on subset of given data
- Dynamic itemset counting
 - Adding candidate itemsets at different points during a scan

Mining frequent itemsets without candidate generation

- Frequent-pattern growth
 - Divide-and-conquer strategy
 - Compress the database representing frequent items into a frequent-pattern tree
 - Divide the compressed database into a set of conditional databases, each associated with one frequent item or “pattern fragment”
 - Mine each conditional database separately
- Set of frequent items sorted in the order of descending support count
- FP-tree construction
 - Branch created for each transaction
 - Count of each node in the common prefix is incremented
- Mining FP-tree
 - Each frequent length-1 pattern (an initial suffix pattern)
 - Conditional pattern base (set of prefix paths)
 - Conditional FP-tree
 - Frequent patterns generation

Mining frequent itemsets using vertical data format

- Horizontal data format {TID-itemset}
- Vertical data format {item-TIDset}
- ECLAT algorithm
 - Intersecting the TIDsets of every pair of frequent items
 - Apriori property
 - Support calculation
 - Diffset technique (track only the differences of TIDsets)

Mining closed frequent itemsets

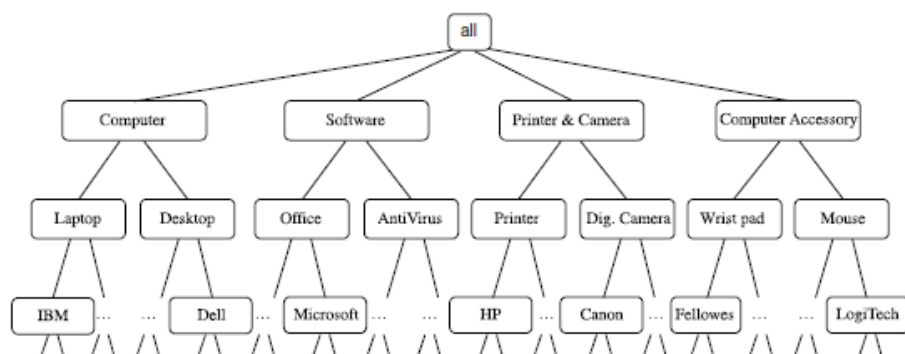
- Naïve approach
- Search for closed itemsets directly during the mining process (pruning techniques)
 - Item merging
 - Sub-itemset pruning
 - Item skipping
- Efficient checking of newly derived frequent itemset
 - Superset checking
 - Subset checking

Mining various kinds of association rules

- Multilevel association rules
- Multidimensional association rules

Mining multilevel association rules

- Transaction database
- Concept hierarchy



- Top-down strategy for support calculation
 - Using uniform minimum support for all levels
 - Using reduced minimum support at lower levels
 - Using item or group-based minimum support

Mining multidimensional association rules

- Database attributes
 - Categorical
 - Quantitative
 - Discretized using predefined concept hierarchies
 - Discretized or “clustered” into bins based on distribution of the data
- Search for frequent predicate sets

Mining multidimensional association rules using static discretization of quantitative attributes

- Quantitative attributes are discretized before mining using predefined concept hierarchies or data discretization techniques
- Categorical attributes may be generalized
- Relational database
 - Frequent itemset mining algorithm
- Data cube

Mining quantitative association rules

- Multidimensional association rules in which the numeric attributes are dynamically discretized during the mining process so as to satisfy some mining criteria (maximizing the confidence or compactness of the rules mined)
- Aquan1,Aquan2=>Acat
- Association rule clustering system
 - Map pairs of quantitative attributed into 2D grid for tuple satisfying a given categorical condition; search for clusters of point for which association rules are generated
 - Binning
 - Equal-width binning
 - Equal-frequency binning
 - Clustering-based binning
 - Finding frequent predicate sets
 - Clustering the association rules (rectangular clusters of rules)

From association mining to correlation analysis

- Strong rules are not necessarily interesting
- $A \Rightarrow B$ [support, confidence, correlation]

– $\text{lift}(A, B) = P(AB) / (P(A)P(B))$

– chi-squared $\chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}}$

– all_conf

$$\text{all_conf}(X) = \frac{\text{sup}(X)}{\text{max_item_sup}(X)} = \frac{\text{sup}(X)}{\max\{\text{sup}(i_j) \mid \forall i_j \in X\}}$$

– cosine

$$\text{cosine}(A, B) = \frac{P(A \cup B)}{\sqrt{P(A) \times P(B)}} = \frac{\text{sup}(A \cup B)}{\sqrt{\text{sup}(A) \times \text{sup}(B)}}$$

Constraint-based association mining

- Knowledge type constraints
 - Association
 - Correlation
- Data constraints
 - Set of task-relevant data
- Dimension/level constraints
 - Desired dimensions
 - Desired levels in the concept hierarchy
- Interestingness constraints
 - Thresholds on statistical measures
- Rule constraints
 - Form of rules
- Metarule-guided mining association rules
- Constraint pushing: mining guided by rule constraints
 - Anti-monotonic
 - Monotonic
 - Succinct
 - Convertible
 - Inconvertible

Outline

- Basic concepts and road map
- Efficient and scalable frequent itemset mining methods
- Mining various kinds of association rules
- From association mining to correlation analysis
- Constraint-based association mining