

ASSOCIATION RULE

Everything in the world is repeated. Francis Bacon

Data Mining

Outline

2

- Basic concepts
- Frequent itemset mining
- Association Rule Mining

Data Mining © Mikhail Zymbler

Motivation



- What are the products that often bought together?
- □ Is beer usually bought with diapers?
- Is the milk bought when both toasts and honey are bought together?

Basic definitions 4 □ Itemset – a collection of one or TID Items more items. 1 bread, milk bread, coffee, eggs, sugar 2 □ *k*-*itemset* – an itemset that contains 3 milk, coffee, coke, sugar k items. bread, coffee, milk, sugar 4 □ *Support count* – frequency of bread, coke , milk, sugar occurrence of an itemset. ■ P({bread,milk,sugar})=2 □ Support – fraction of transactions that contain an itemset. ■ sup({bread,milk,sugar})=0.4

□ *Frequent itemset* – an itemset whose support is greater than or equal to a *minsup* threshold.



Basic approaches

- □ Brute force
- Apriori

6

- Vertical data format
- □ FP-growth



Frequent itemset generation strategies

- Reduce the number of candidates (*M*)
 Complete search: *M*=2^d
 Use pruning techniques to reduce *M*
- Reduce the number of transactions (N)
 Reduce size of N as the size of itemset increases
- Used by vertical-based mining algorithms
- Reduce the number of comparisons (N · M)
 Use efficient data structures to store the candidates or transactions

 \blacksquare No need to match every candidate against every transaction

Data Mining © Mikhail Zymbler

Reducing number of candidates

□ Anti-monotone property of support

■ Support of an itemset never exceeds the support of its subsets (e.g. *sup*({bread, milk, chips}))≤*sup*({bread, milk}))

 $\forall X, Y : (X \subseteq Y) \Longrightarrow \sup(X) \ge \sup(Y)$

□ Apriori principle

- □ If an itemset is frequent, then all of its subsets must also be frequent
- Converse assertion: if some subset is infrequent then whole set is infrequent









Apriori algorithm

- □ Generate frequent 1-itemsets
- \square Let k=1
- Repeat until no new frequent itemsets are identified
 Generate candidate (k+1)-itemsets using frequent k-itemsets
 - Prune candidate itemsets containing infrequent *k*-subsets
 - Count the support of each candidate by scanning the DB
 - Eliminate candidates that are infrequent, leaving only those that are frequent

Apriori algorithm

```
\begin{array}{l} C_k: \mbox{ candidate itemset of size } k\\ L_k: \mbox{ frequent itemset of size } k\\ L_1 = \{\mbox{ frequent items}\};\\ \mbox{ for } (k=1; \ L_k \ != \ \varnothing; \ k++) \ \mbox{ do begin}\\ C_{k+1} = \mbox{ candidates generated from } L_k;\\ \mbox{ for each transaction } t \ \mbox{ in database } \mbox{ do begin}\\ \mbox{ increment the count of all candidates in } C_{k+1}\\ \mbox{ that are contained in } t\\ L_{k+1} = \mbox{ candidates in } C_{k+1} \ \mbox{ with } minsup\\ \mbox{ end}\\ \mbox{ return } \cup_k \ L_k; \end{array}
```

Data Mining © Mikhail Zymbler



How to generate candidates

```
□ Suppose the items in L_{k-1} are listed in an order
```

□ Step 1: self-joining L_{k-1}

Step 2: pruning

forall itemsets c in C_k do
forall (k-1)-subsets s of c do
 if (s is not in L_{k-1}) then delete c from C_k

Factors affecting complexity

- Choice of minimum support threshold
- I lowering support threshold results in more frequent itemsets; this may increase number of candidates and max length of frequent itemsets
- Dimensionality (number of items) of the data set
 more space is needed to store support count of each item
 if number of frequent items also increases, both computation and I/O costs may also increase
- □ Size of database
- since Apriori makes multiple passes, run time of algorithm may increase with number of transactions
- Average transaction width
 - transaction width increases with denser data sets
 - this may increase max length of frequent itemsets and traversals of hash tree (number of subsets in a transaction increases with its width)

Data Mining © Mikhail Zymbler

Improving the efficiency of Apriori

- □ Reduce the number of comparisons
- □ Reduce the number of transactions
- □ Partitioning the data to find candidate itemsets
- □ Sampling: mine on a subset of the given data





Reduce the number of transactions

A transaction that doesn't contain any frequent kitemset can not contain any frequent j-itemset, for any j>k. So such transaction can be marked or removed from subsequent scans of the database for j-itemsets.

Data Mining © Mikhail Zymbler





Sampling: mine on a subset of the given data

- □ Pick random sample *S* of the given data *D* (S should fit in the available memory)
- Search for frequent itemsets in S instead in D.
 Lower *minsup* to reduce the number of missed frequent itemsets.
- \Box Find the L_s , frequent itemsets in S.
- □ The rest of the database can be used to compute the actual frequencies of the itemset in L_s .
- □ If L_s doesn't contain all the frequent itemsets in D, then a second pass will be needed.

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

22

of candidates: $({}_{100}{}^1) + ({}_{100}{}^2) + \dots + ({}_{100}{}^{100}) = 2^{100} \cdot 1 = 1.27*10^{30}!$

Data Mining © Mikhail Zymbler

Bottleneck is candidate-generation-and-test

Vertical data format TID Items В С D Е A 2 3 2 1 A,B,E 1 1 1 2 B,C,D 4 2 4 3 3 C,E 5 5 4 5 6 6 7 4 A,C,D 7 8 9 5 A,B,C,D 8 9 6 A,E 8 10 7 A,B 9 8 A,B,C 9 A,C,D 10 B tid-lists Data Mining © Mikhail Zymbler







Using RDBMS and SQL

- Store data in relational database, implement mining algorithms in SQL
- Pros
 - No need to export data for external mining utilities and import mining results into database
 - No memory limits
- Cons
 - **RDBMS** will lose the competition with mining utilities in case of data fits in memory
 - SQL has less flexibility than high-level programming languages and may produce cumbersome source code

Data Mining © Mikhail Zymbler



Candidate generation in SQL

insert into Cand select * from TDB where item in (select item from TDB group by item having count(*) >= minsup);

Finding frequent itemsets in SQL

```
-- Frequent 2-itemsets
select A.item, B.item, count(A.tid) as
sup_count
from Cand A, Cand B
where A.tid=B.tid and A.item<B.item
group by A.item, B.item
having count(A.tid) >= minsup;
```

-- Frequent itemsets for k>2
-- How to implement k:=k+1 ?

Data Mining © Mikhail Zymbler

Association rules

□ Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction

Market-Basket transactions Examples of Association Rules

TID	Items	$sugar \rightarrow coffee$,
1	bread, milk	${milk, bread} \rightarrow {eggs, cola},$
2	bread, coffee, eggs, sugar	$\{\text{coffee, bread}\} \rightarrow \{\text{milk}\},\$
3	milk, coffee, cola, sugar	$\{\text{mirk}\} \rightarrow \{\text{coffee, bread}\},\$
4	bread, coffee, milk, sugar	· ···
5	bread, cola, milk, sugar	
		Data Mining C Mikhail Zymbler

Association rule

- □ Association rule is an implication expression of the form $X \rightarrow Y$ ("if X then Y"), where X and Y are itemsets
- Rule support

■ is a fraction of transactions that contain both X and Y ■ $sup(X \rightarrow Y)=P(X,Y)/|TID|$

□ Rule confidence

• shows how often items in Y appear in transactions that contain X

 $\Box conf(X \rightarrow Y) = P(X,Y)/P(X)$

Support and confidence of the					
ass	sociation rule				
31					
TID	Items				
1	bread, milk				
2	bread, coffee, eggs, sugar				
3	milk, coffee, cola, sugar				
4	bread, coffee, milk, sugar				
5	bread, cola, milk, sugar				
 [5] [Pictal; end; finit; sugar [6] {milk, sugar} → {coffee} [6] sup=P(milk, sugar, coffee)/[TID]=2/5 [6] conf=P(milk, sugar, coffee)/P(milk, sugar)=2/3 [6] {coffee} → {milk, sugar} [6] sup=2/5 [6] conf=P(coffee, milk, sugar)/P(coffee)=2/4 					
		Data Mining	© Mikhail Zymbler		





Association rule mining task

Given a set of transactions T, the goal of association rule mining is to find all the rules having

- support \geq *minsup* threshold
- confidence \geq *minconf* threshold
- Brute-force approach:
 - List all possible association rules
 - Compute the support and confidence for each rule
 - Prune rules that fail the *minsup* and *minconf* thresholds
 - Computationally prohibitive!

Mining association rules

34				
TID	Items	Rule	sup	conf
1	bread, milk	{milk,sugar}→coffee	0.4	0.67
2	bread, coffee, eggs, sugar	{milk,coffee}→sugar	0.4	1.0
3	milk, coffee, cola, sugar	$\{sugar, coffee\} \rightarrow milk$	0.4	0.67
4	bread, coffee, milk, sugar	$coffee \rightarrow \{milk, sugar\}$	0.4	0.67
5	bread cola milk sugar	$sugar \rightarrow \{milk, coffee\}$	0.4	0.5
5	bread, cold, milly sugar	$milk \rightarrow \{sugar, coffee\}$	0.4	0.5

Observations

- All the above rules are binary partitions of the same itemset: {milk, sugar, coffee}
- Rules originating from the same itemset have identical support but can have different confidence

Data Mining C Mikhail Zymbler

Mining association rules

- □ Two-step approach:
 - 1. Frequent Itemset Generation
 - Generate all itemsets whose support grater than *minsup*
 - 2. Rule Generation
 Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset
- Frequent itemset generation is still computationally expensive