

9. Data mining

- definition
- basic concepts
- applications
- challenges

Definition - Data Mining

- extraction of [unknown] patterns from data for actionability
- combines methods from:
 - databases
 - machine learning
 - visualization
- involves large datasets
- consists of:
 - stating the [business] question
 - data collection and (instance) selection
 - preprocessing
 - transformation
 - model building
 - interpretation/evaluation/deployment

Model building

- Supervised
 - (mainly classification)
 - Also ranking, estimation
- Unsupervised
 - Associations
 - Clustering

Associations

Given:

$I = \{i_1, \dots, i_m\}$ set of items

D set of transactions (a database), each transaction is a set of items $T \subset 2^I$

Association rule: $X \Rightarrow Y$, $X \subset I$, $Y \subset I$, $X \cap Y = \emptyset$

confidence c : ratio of # transactions that contain both X and Y to # of *all* transaction that contain X

support s : ratio of # of transactions that contain both X and Y to # of transactions in D

Itemset is *frequent* if its support $> \theta$

An *association rule* $A \Rightarrow B$ is a conditional implication among itemsets A and B , where $A \subset I$, $B \subset I$ and $A \cap B = \emptyset$.

Support of an association rule = $P(A \cup B)$. The *confidence* of an association rule $r: A \Rightarrow B$ is the conditional probability that a transaction contains B , given that it contains A . Confidence = $P(B|A)$

The support of rule r is defined as: $sup(r) = sup(A \cup B)$. The confidence of rule r can be expressed as $conf(r) = sup(A \cup B) / sup(A)$.

Formally, let $A \subset 2^I$; $sup(A) = |\{t: t \in D, A \subset t\}| / |D|$, if $R = A \Rightarrow B$ then $sup(R) = sup(A \cup B)$, $conf(R) = sup(A \cup B) / sup(A)$

Itemsets and association rules

- Itemset = set of items
- k-itemset = set of k items
- Finding association rules in databases:
 - Find all frequent (or large) itemsets (those with support $> \min_s$)
 - Generate rules that satisfy minimum confidence

Example

- Computer store
- Customers buying computers and financial software
- What does the rule mean:

computer → *financial_mgmt_software*

[support = 2%, conf = 60%]

Associations - mining

Given D , generate all assoc rules with $c, s >$
thresholds \min_c, \min_s

(items are ordered, e.g. by barcode)

Idea:

find all itemsets that have transaction
support $> \min_s$: **large itemsets**

Associations - mining

to do that: start with indiv. items with large support

in ea next step, k ,

- use itemsets from step $k-1$, generate new itemset C_k ,
- count support of C_k (by counting the candidates which are contained in any t),
- prune the ones that are not large

Apriori property

- All [non-empty] subsets of a frequent itemset must be frequent
- Based on the fact that an itemset i that is NOT frequent has support $< \min_s$
- But inserting an additional item A in i will not increase the support of $i \cup A$

Associations - mining

```
1)  $L_1 = \{\text{large 1-itemsets}\};$ 
2) for (  $k = 2; L_{k-1} \neq \emptyset; k++$  ) do begin
3)    $C_k = \text{apriori-gen}(L_{k-1});$  // New candidates
4)   forall transactions  $t \in \mathcal{D}$  do begin
5)      $C_t = \text{subset}(C_k, t);$  // Candidates contained in  $t$ 
6)     forall candidates  $c \in C_t$  do
7)        $c.\text{count}++;$ 
8)     end
9)    $L_k = \{c \in C_k \mid c.\text{count} \geq \text{minsup}\}$ 
10) end
11) Answer =  $\bigcup_k L_k;$ 
```

subset(C_k, t) denotes those itemsets that are contained in transaction t

Candidate generation

$C_k = \text{apriori-gen}(L_{k-1})$

```
insert into  $C_k$ 
select  $p.\text{item}_1, p.\text{item}_2, \dots, p.\text{item}_{k-1}, q.\text{item}_{k-1}$ 
from  $L_{k-1} p, L_{k-1} q$ 
where  $p.\text{item}_1 = q.\text{item}_1, \dots, p.\text{item}_{k-2} = q.\text{item}_{k-2},$ 
       $p.\text{item}_{k-1} < q.\text{item}_{k-1};$ 
```

Next, in the *prune* step, we delete all itemsets $c \in C_k$ such that some $(k-1)$ -subset of c is not in L_{k-1} :

```
forall itemsets  $c \in C_k$  do
  forall  $(k-1)$ -subsets  $s$  of  $c$  do
    if  $(s \notin L_{k-1})$  then
      delete  $c$  from  $C_k;$ 
```

Select from $k-1$ -frequent itemsets two overlapping subsets, add the differences

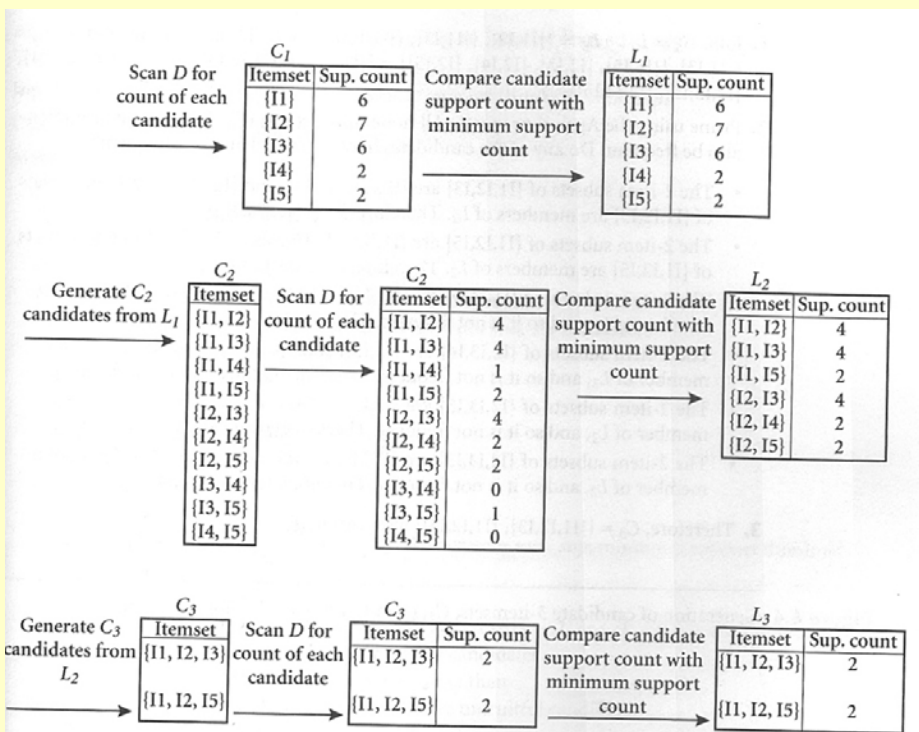
Example

TID	List of item_IDs
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

From Han,
Kamber, "Data
Mining", p. 232

$$I = \{I1, \dots, I5\}$$

$$\min_s = 2$$



Firstly, $C_3 = \{\{I1, I2, I3\}, \{I1, I2, I5\}, \{I1, I3, I5\}, \{I2, I3, I4\}, \{I2, I3, I5\}, \{I2, I4, I5\}\}$

Only $\{I1, I2, I3\}, \{I1, I2, I5\}$ are left

$C_4 = \{I1, I2, I3, I5\}$ is attempted but pruned, $C_4 = \emptyset$ terminates the algorithm

From itemsets to association rules

- For ea. frequent itemset l generate all the partitions of l into $s, l-s$
- Attempt a rule $s \rightarrow l-s$ iff
 $support_count(l)/support_count(s) > min_c$
- e.g. for $min_c = 0.5$, what rules do we get?
[$conf(r) = sup(A \cup B)/sup(A)$]