



Association Rules Discovery Part one: Apriori

Market basket problem

- Given:
 - Large number of items : bread, milk, banana, cereal, ...
 - Customer fill their basket with a subset of these items
 - Available recordings of the content of the baskets
- Goal:
 - Derive “What items do people buy together?”
- Purpose:
 - Marketers use the information to position items, and control the way a typical customer traverse the store
 - Targeted advertisement during online shopping/ browsing

Related Problems

- Similar problems
 - Information gathering from text documents
 - Baskets : documents
 - Items : words
 - Goal: documents share groups of words may indicate the resemblance of their content → information retrieval and intelligence gathering
 - Finding mirror sites
 - Baskets : documents
 - Items : sentences
 - Goal : web page containing groups of the same sentences may indicate they are mirror sites on the web

Goal of Market-basket analysis

- Association rule discovery
- Causality analysis

What is Association Rule ?

- Assume :
 - $J = \{i_1, i_2, \dots, i_m\}$ is a set of items;
 - D , be a set of database transactions where each transaction T is a set of items, such that T is a subset of J ;
- $T_1 : i_2, i_5, i_6, i_{12}, i_9, i_{30}, i_{55}, \dots$
 $T_2 : i_1, i_2, i_5, i_{38}, i_{39}, i_{100}, i_{121}, \dots$
 $T_3 : i_4, i_6, i_{23}, i_{29}, i_{59}, i_{44}, \dots$
 \dots
 $T_N : i_6, i_9, i_{35}, i_{26}, i_{40}, \dots$

Association Rule

- An **association rule** is an implication of the form

$$X \Rightarrow Y,$$
 where X and Y are subsets of J , and X and Y do not share common item.

X and Y are sets of items.

A transaction T is said to **contain** X (or Y) if and only if X is a subset of T .

Support and Confidence

- Support of the rule:
The rule $X \rightarrow Y$ holds in the transaction set D with **support** s , where s is the percentage of transactions in D that contain both X and Y. (Computed as $P(X \text{ and } Y)$)
- Confidence of the rule:
The rule $X \rightarrow Y$ has **confidence** c in the transaction set D if c is the percentage of transactions in D containing X that also contain Y. (Computed as $P(Y|X)$)

Middle Tennessee State University

Example Association Rules

- Examples:
 1. Computer \rightarrow financial_management_software
[support = 2%, confidence = 60%]
 2. Diaper \rightarrow beer
[support = 3.5%, confidence = 45%]
 3. Milk, butter \rightarrow bread
[support = 6%, confidence = 60%]

Middle Tennessee State University

Practice problem

- Given transaction database, D:
T1 bread, milk, banana, cereal, apple, sugar, flour, butter
T2 cereal, pear, sugar, salt, egg, flour, milk
T3 bread, milk, potato, onion, apple
T4 potato chip, orange juice, coke, ice cream
T5 coke, potato chip, sugar, flour, milk

Assume that : **milk \rightarrow apple** is an association rule discovered:
 - What is the support for this rule?
 - What is the confidence of this rule ?
 - how about **{sugar, flour} \rightarrow egg** ?

Middle Tennessee State University

Causality Analysis

- Causality analysis:
 - Does the presence of X actually causes Y to be bought?
 - Test method:
 - Question: does diaper causes beer to be bought? Or does beer causes diaper to be bought?
 - Approach 1: lower the price of diaper, raise the price of beer
 - Approach 2: lower the price of beer, raise the price of diaper
 - Result

Middle Tennessee State University

Terminologies

- **itemset** : a set of items
- **k-itemset** : an itemset that contains k items
- **occurrence frequency of an item**: the number of transactions that contain the itemset
- **frequent k-itemset**: a k-itemset whose occurrence frequency is greater than or equal to a pre-defined minimum support count
- **minimum support (count)**
- **minimum confidence**
- **strong association rules** : association rules that satisfy both the minimum support and the minimum confidence threshold

Middle Tennessee State University

Apriori based association rule discovery methods

- Finding frequent item sets
 - Frequent item set : set of items appearing in at least fraction s of the baskets
 - Why do we need to find frequent item sets?
 - What do we mean by frequent?
 - Frequent item sets can be found efficiently using the **Apriori** property
- What is the **Apriori** property?
Apriori property : If a set of items S is frequent, then every subset of S is also frequent

Middle Tennessee State University

Apriori property

- How does Apriori property help in finding the frequent item sets efficiently?
 - Proceeds level wise, start from frequent single items, then find the frequent pairs, the frequent triples, ...

Middle Tennessee State University

The Basic Process

- The algorithm proceeds level-wise:
 - Given minimum support count s , in the first pass find the 1-itemsets (sets with single items) that appear in at least s number of baskets. (L_1)
 - Pairs of items in L_1 become the candidate pairs C_2 for the second pass. The pairs in C_2 whose count reaches s are the frequent pairs, L_2 .

Middle Tennessee State University

The basic process (cont.)

- The candidate triples, C_3 are those sets such that all of the 2-itemsets are frequent. E.g., for $\{A, B, C\}$ to be candidate 3-itemset, $\{A, B\}$, $\{B, C\}$, and $\{A, C\}$ should all be frequent 2-itemset. Count the occurrences of triples in C_3 , those with a count of at least s are the frequent triples, L_3 .
- Proceed as this until the i th frequent item set becomes empty.
 - L_i : frequent item set of size i
 - C_i : candidate item set of size i

Middle Tennessee State University

The Apriori Algorithm

Identify frequent 1-itemset by scanning the database
For each k value (starting with $k=2$, ends when L_{k-1} becomes empty):

1. Applying candidate generation to generate all possible k -itemsets based on the frequent 1-itemsets
 - the join step :
 - L_k is found by joining L_{k-1} with itself
 - Apriori assumes that all items within a transaction or itemset are sorted in lexicographic order.
 - Two items may be joined if they share the first $k-2$ items (Why?)

Middle Tennessee State University

The Apriori Algorithm

- the pruning step:
 - C_k is a superset of L_{k-1} that is, its members may or may not be frequent, but all of the frequent k -itemsets are included in C_k .
 - 2. Eliminate candidate k -itemsets that have support smaller than the minimum support count
- Return the union of all L_k

Middle Tennessee State University

Algorithm Apriori

```

L1 = find_frequent_1-itemsets (D);
for (k=2; Lk-1 != φ; k++) {
    Ck = apriori_gen (Lk-1);
    for each transaction t in D {
        Ct = subset(Ck, t);
        for each candidate c in Ct
            c.count ++;
    }
    Lk = {c in Ck | c.count >= minimum_support_count}
    L = L union Lk;
}
return L;
```

Middle Tennessee State University

Algorithm Apriori (cont.)

```

procedure apriori_gen( $L_{k-1}$ )    // frequent (k-1)-itemset
  for each itemset  $l_1$  in  $L_{k-1}$ 
    for each itemset  $l_2$  in  $L_{k-1}$  {
      if ( $l_1[1] = l_2[1] \wedge (l_1[2]=l_2[2]) \wedge \dots \wedge (l_1[k-2]=l_2[k-2]) \wedge (l_1[k-1] < l_2[k-1])$ )
      {
         $c = l_1$  join  $l_2$ ; // join step
        if has_infrequent_subset( $c, L_{k-1}$ ) then
          delete  $c$ ; // pruning step
        else
          add  $c$  to  $C_k$ ;
      }
    }
  return  $C_k$ ;

```

Middle Tennessee State University

Algorithm Apriori (cont.)

```

procedure has_infrequent_subset ( $C_k, L_{k-1}$ )
  for each (k-1)-subset  $s$  of  $c$ 
    if  $s$  not in  $L_{k-1}$  then
      return TRUE;
  return FALSE;

```

Middle Tennessee State University

Practice Question

TID	List of items	
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	

Given the transaction data,
find all frequent item sets
having minimum support
count = 2

Middle Tennessee State University

Step Two: Generate strong association rules from the frequent itemsets

- Computation:**
 $\text{confidence}(X \rightarrow Y) = p(Y|X)$
 $= \frac{\text{support_count}(X \wedge Y)}{\text{support_count}(X)}$
- $\text{support_count}(X \wedge Y)$: number of transactions containing the itemsets X and Y , and
 $\text{support_count}(X)$ is the number of transactions containing the itemset X .

Middle Tennessee State University

Generate strong association rules

- The basic idea:** given a frequent itemset I , generate all strong associate rules based on I
- Direct approach:**
 - 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 $\text{support_count}(I) / \text{support_count}(s) \geq \text{minimum confidence count}$

Middle Tennessee State University

Properties of the confidence measure

Given : a is a subset of I , a' is a subset of a

$\text{support}(I) \leq \text{support}(a) \leq \text{support}(a')$

$\text{confidence}(a \rightarrow (I - a)) \geq \text{confidence}(a' \rightarrow (I - a'))$

Why?

Middle Tennessee State University

Properties of the confidence measure

Similarly,

confidence of $((l - a) \rightarrow a) \leq \text{confidence of } ((l - a') \rightarrow a')$

Therefore,

If $((l - a') \rightarrow a')$ is not a strong rule, then none of the rules of the form $((l - a) \rightarrow a)$ can be strong, where a' is a subset of a

Middle Tennessee State University

Rule Generation

• The apriori approach for rule generation

– Basis: If $((l - a') \rightarrow a')$ is not a strong rule, then none of the rules of the form $((l - a) \rightarrow a)$ can be strong, (a' is a subset of a)

– Approach :

- Start by generating rules that have a single consequent
- Increase size of consequent to 2 by the “ap_gen” function, based on only the successful single consequents
- Continue to increase the number of consequents, (equivalently, decreasing the size of the antecedent)...., one item at a time, until the antecedent becomes empty

Middle Tennessee State University

Example

TID	List of items	Given the transaction data, find all association rules having minimum support count = 2, and minimum confidence of 70%.
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	

Middle Tennessee State University

Example

TID	List of items	Given the transaction data, find all association rules having minimum support count = 2, and minimum confidence of 70%. We already derived frequent itemsets for this TD as: {I1}, {I2}, {I3}, {I4}, {I5}, {I1, I2}, {I1, I3}, {I1, I5}, {I2, I3}, {I2, I4}, {I2, I5}, {I1, I2, I3}, {I1, I2, I5}
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	

Suppose $I_k = \{I1, I2, I5\}$

Middle Tennessee State University

Rule generation algorithm

- 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 = $\bigcup_k L_k$;

Middle Tennessee State University

Rule generation algorithm (cont.)

$H_1 = \{\text{consequences of rules from } I_k \text{ with one item in the consequent}\}$;

Given I_k :

for each item i in I_k ,

- (1) Treat it as consequent,
- (2) form rule $I_k - i \rightarrow i$
- (3) compute confidence = $\text{support}(I_k) / \text{support}(I_k - i)$
if confidence > min_conf
output rule $I_k - i \rightarrow i$
add i to H_1

Middle Tennessee State University

Rule generation algorithm (Cont.)

```
ap-genrules ( $l_k$ : large k-itemset,  $H_m$ : set of m-item consequents)
if ( $k > m+1$ ) then begin
     $H_{m+1} = \text{apriori-gen}(H_m)$ ;
    forall  $h_{m+1}$  belongs to  $H_{m+1}$  do
        Conf = support( $l_k$ )/support( $l_k - h_{m+1}$ );
        if (Conf  $\geq$  minimum_confidence) then
            Output the rule ( $l_k - h_{m+1}$ )  $\rightarrow h_{m+1}$  with
            confidence = Conf, support = support( $l_k$ );
        else
            Delete  $h_{m+1}$  from  $H_{m+1}$ 
    Call ap-genrules( $l_k$ ,  $H_{m+1}$ );
```

Middle Tennessee State University

Discussion

- Why use apriori-gen to create the consequents?
 - Join: form the consequent part of the rule, with successively larger sizes
 - Prune: eliminate no-hope candidates
- Why is it necessary to delete h_{m+1} from H_{m+1} ?

Middle Tennessee State University

Implementation of Apriori

- Hash table
- Hash tree
- Hash tree is used for two steps of Apriori
 - Test whether “subset s of candidate itemset is in L_{k-1} ”
 - Test whether “a candidate itemset C_k is in a transaction t”
- Time complexity of the algorithm

Middle Tennessee State University

Subset(candidate itemset, L_{k-1})

- Goal : check “is there any (k-1) subset of c that is not in L_{k-1} ?
- Approach:
 - All items in L_{k-1} are stored in a hash tree
 - For each non-empty (k-1) subsets of a k-itemset, checking whether a (k-1) subset is in L_{k-1} takes $O(1)$

Middle Tennessee State University

Subset (c_k , t)

- Assumption:
 - Items in transactions are ordered
 - Items in candidate set are ordered
 - Candidate c_k are put in a hash tree
- Approach:
 - At root level, hash on every item in the transaction,
 - At level i,
 - if it is an interior node, hash on every item following the ith item,
 - if it is a leaf node, check if the candidate c is in the list
 - if yes, update the counter for that candidate

Middle Tennessee State University

Improve the efficiency of Apriori

- AprioriTid for frequent item set generation
- Data partitioning
- Data sampling
- Other approaches

Middle Tennessee State University

AprioriTid

- Objectives: as the size of the frequent itemsets increases,
 - reduce the length of each transaction
 - reduce the number of transactions necessary to check for support
- Method:
 - Database D is not used for counting support after the first pass.
 - The set $\bar{C}_k < \text{Tid}, \{X_k\} >$ is used for counting support afterwards, where $\{X_k\}$ is a potentially large k-itemset present in the transaction with identifier Tid.

Middle Tennessee State University

Practice Problem

TID	items
T1	I1, I2, I3, I5
T2	I2, I4
T3	I2, I6
T4	I1, I2, I4, I5
T5	I1, I2
T6	I1, I2, I3, I5
T7	I1, I2, I3

Middle Tennessee State University

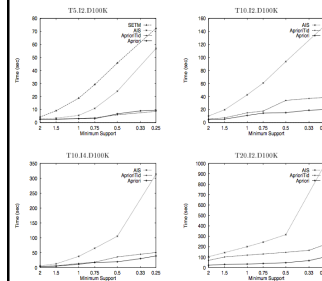
Algorithm AprioriTid

```

1)  $L_1 = \{\text{large 1-itemsets}\};$ 
2)  $\bar{C}_1 = \text{database } D;$ 
3) for (  $k = 2; L_{k-1} \neq \emptyset; k++$  ) do begin
4)    $C_k = \text{apriori-gen}(L_{k-1});$  // New candidates
5)    $\bar{C}_k = \emptyset;$ 
6)   forall entries  $t \in \bar{C}_{k-1}$  do begin
7)     // determine candidate itemsets in  $C_k$  contained
       // in the transaction with identifier t.TID
        $C_t = \{c \in C_k \mid (c - c[k]) \in t.\text{set-of-itemsets} \wedge$ 
          $(c - c[k-1]) \in t.\text{set-of-itemsets}\};$ 
8)     forall candidates  $c \in C_t$  do
9)       c.count++;
10)    if ( $C_t \neq \emptyset$ ) then  $\bar{C}_k += \langle t.\text{TID}, C_t \rangle;$ 
11)  end
12)   $L_k = \{c \in C_k \mid c.\text{count} \geq \text{minsup}\}$ 
13) end
14) Answer =  $\bigcup_k L_k;$ 
    
```

Middle Tennessee State University

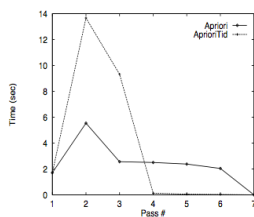
Apriori vs. AprioriTid Performance



- AprioriTid replaces a pass over the original dataset by a pass over the set \bar{C}_k .
 - AprioriTid is very effective in later passes when the size of \bar{C}_k becomes small compared to the size of the database.
- AprioriTid beats Apriori when its \bar{C}_k sets can fit in memory.
- When \bar{C}_k does not fit in memory, there is a jump in the execution time for AprioriTid.

Middle Tennessee State University

Algorithm AprioriHybrid



Apriori vs. AprioriTid (T10.I4.D100k, minsup = 0.75%)

Middle Tennessee State University

Use Apriori for the initial passes, and switch to AprioriTid when it expects that the set \bar{C}_k at the end of the pass will fit in memory.

Improve the efficiency of Apriori

- AprioriTid for frequent item set generation
- Data partitioning
- Data sampling
- Other approaches

Middle Tennessee State University

Data Partitioning

- Phase I :
 - subdivide the transactions of D into N non-overlapping partitions. Minimum_support_count for each partition = min_support(%) * the number of transactions in that partition
 - find all frequent itemsets within each partition → **local frequent itemsets**
 - form **global frequent itemset** by collecting the **local frequent itemsets**
 - Questions:
 - *Do all itemsets in the global frequent itemset necessarily be a frequent itemset in the overall database D?*
 - *If an itemset is frequent given database D, does it have to be in the global frequent itemset?*

Middle Tennessee State University

Data Partitioning

- Phase II :
 - scan database to determine the actual support of each candidate itemset in the global frequent itemset. Delete those having small support count

Middle Tennessee State University

Sampling

- Pick a random sample of the given data D, and search for frequent itemsets in S instead of D.
 - *Do all candidate itemsets from S necessarily be a frequent itemset in the overall database D?*
 - *If an itemset is frequent given database D, does it have to appear in the candidate frequent itemset?*
- Method:
 - lower the support count requirement when finding frequent itemsets in S, L^S
 - the original data D is used to find the actual support count for candidate itemsets in L^S

Middle Tennessee State University