

Foundations of Data Science & Analytics

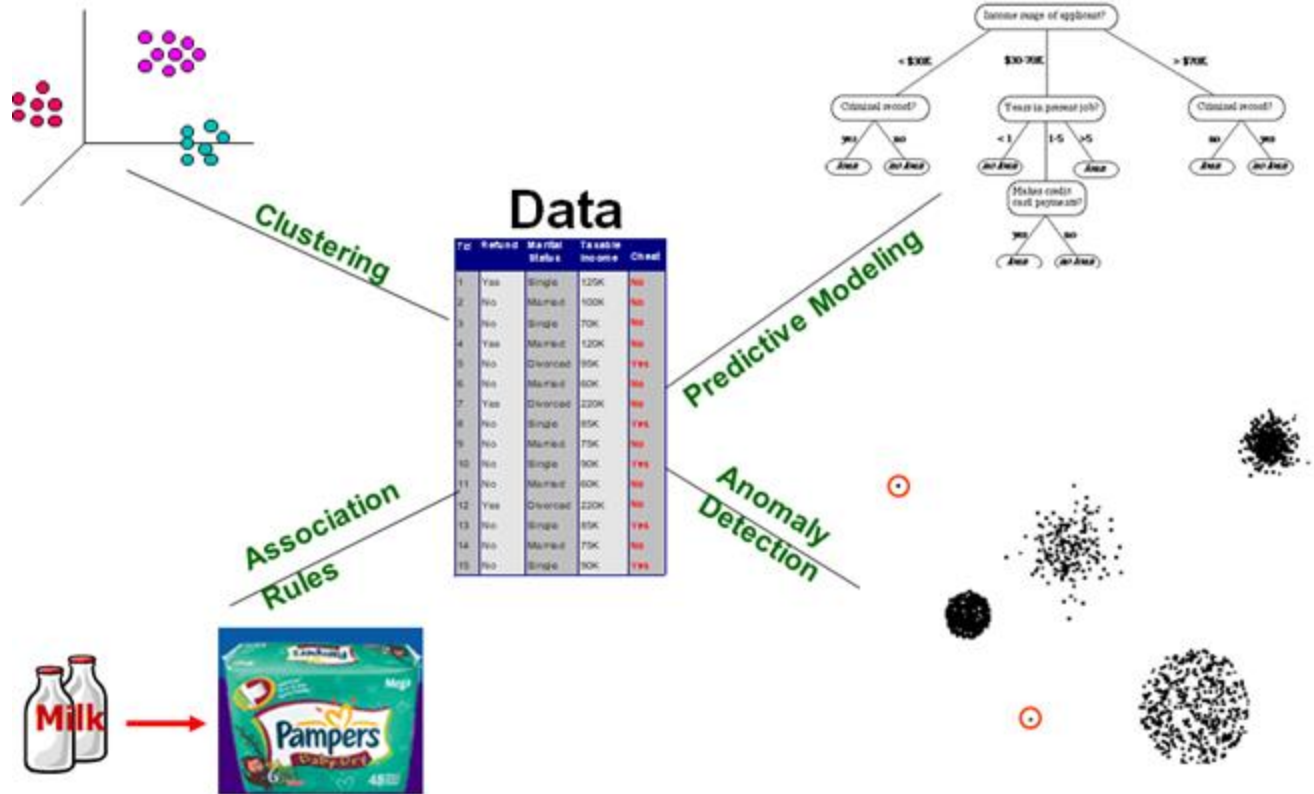
Association Rule Mining

Ezgi Siir Kibris

[Introduction to Data Mining, 2nd Edition](#)

by
Tan, Steinbach, Karpatne, Kumar

Tasks



01/17/2018

Introduction to Data Mining, 2nd Edition

11

Association Rule Mining

Market-Basket transactions

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

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.

$\{\text{Diaper}\} \rightarrow \{\text{Beer}\},$
 $\{\text{Milk, Bread}\} \rightarrow \{\text{Diaper}\},$
 $\{\text{Beer, Bread}\} \rightarrow \{\text{Milk, Diaper}\}$

co-occurrence, not **causality**

Frequent Itemset

Market-Basket transactions

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

- **Itemset**

- A collection of one or more items
 - {Milk, Bread, Diaper}
- k-itemset
 - An itemset that contains k items

- **Support count**

- Frequency of occurrence of an itemset
 - $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$

- **Support**

- Fraction of transactions that contain an itemset
 - $s(\{\text{Milk, Bread, Diaper}\}) = 2/5$

- **Frequent Itemset**

- An itemset whose support is greater than or equal to a *minsup* threshold

Association Rules

Market-Basket transactions

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

- **Association Rule**
 - An implication expression of the form $X \rightarrow Y$, where X and Y are itemsets
 - $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$
- **Rule Evaluation Metrics**
 - **Support (s)**
 - Fraction of transactions that contain both X and Y
 - **Confidence (c)**
 - Measures how often items in Y appear in transactions that contain X

Association Rules

Market-Basket transactions

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example:

$$\{\text{Milk, Diaper}\} \Rightarrow \{\text{Beer}\}$$

$$s = \frac{\sigma(\text{Milk, Diaper, Beer})}{|T|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

Association Rule Mining

Market-Basket transactions

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

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

- support \geq *minsup* threshold
- confidence \geq *minconf* threshold

How?

Market-Basket transactions

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Two-step approach:

- **Frequent Itemset Generation**
 - Generate all itemsets whose support $\geq \text{minsup}$
- **Rule Generation**
 - Generate high confidence ($\geq \text{minconf}$) rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset

How?

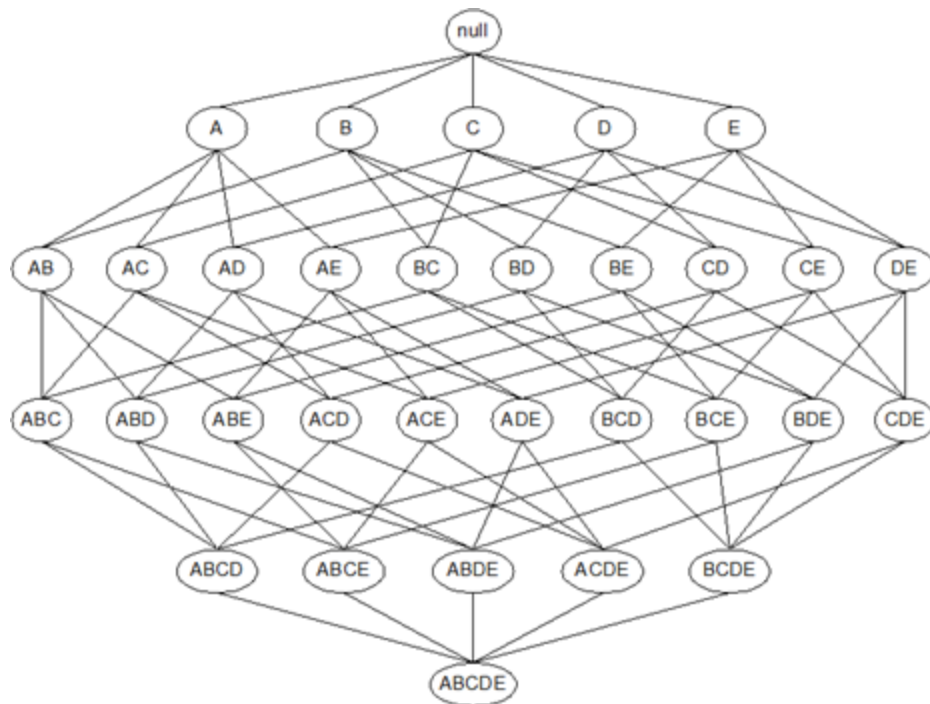
Market-Basket transactions

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Two-step approach:

- **Frequent Itemset Generation**
 - Generate all itemsets whose support $\geq \text{minsup}$
- **Rule Generation**
 - Generate high confidence ($\geq \text{minconf}$) rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset

Frequent Itemset Generation



Given d items, there are 2^d possible candidate itemsets

Frequent Itemset Generation

Apriori principle:

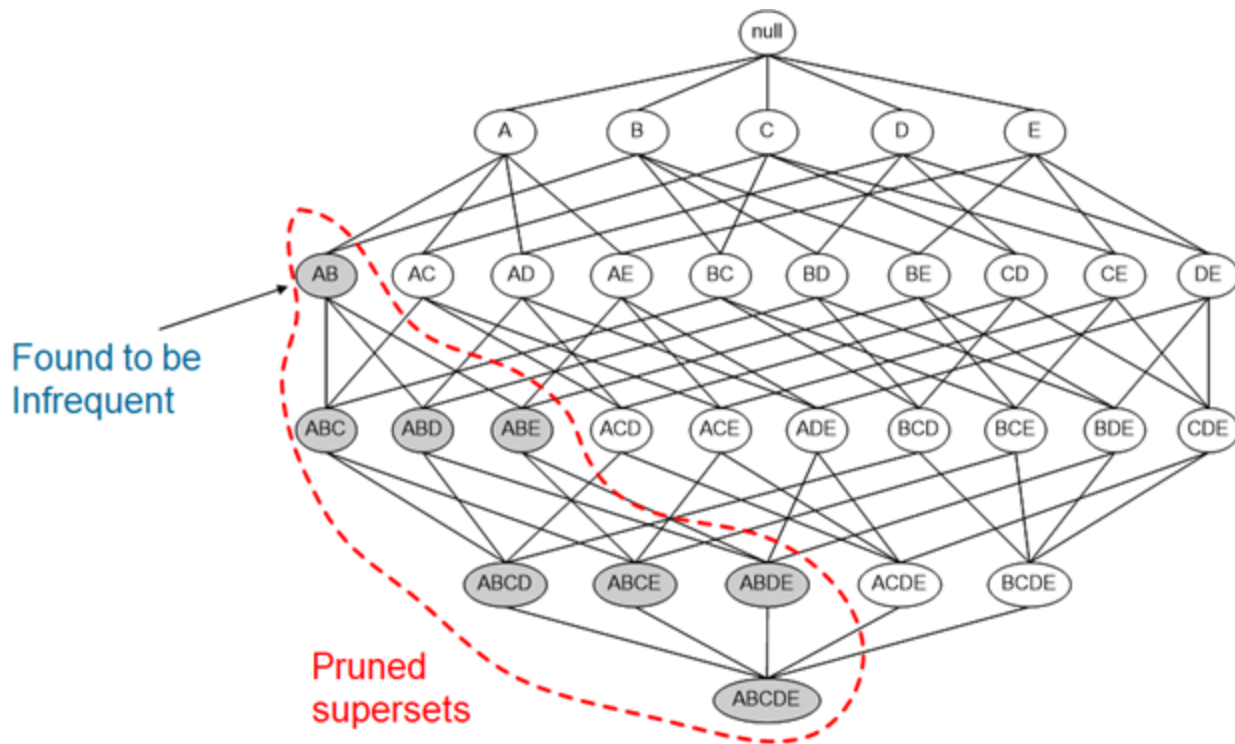
- If an itemset is frequent, then all of its subsets must also be frequent

Apriori principle holds due to the following property of the support measure:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

- Support of an itemset never exceeds the support of its subsets
- This is known as the ***anti-monotone*** property of support

Apriori Principle



Apriori Algorithm

- F_k : frequent k-itemsets
- L_k : candidate k-itemsets

□ Algorithm

- Let $k=1$
- Generate $F_1 = \{\text{frequent 1-itemsets}\}$
- Repeat until F_k is empty
 - ◆ **Candidate Generation**: Generate L_{k+1} from F_k
 - ◆ **Candidate Pruning**: Prune candidate itemsets in L_{k+1} containing subsets of length k that are infrequent
 - ◆ **Support Counting**: Count the support of each candidate in L_{k+1} by scanning the DB
 - ◆ **Candidate Elimination**: Eliminate candidates in L_{k+1} that are infrequent, leaving only those that are frequent $\Rightarrow F_{k+1}$

Generate 1-Itemset

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Beer, Bread, Diaper, Eggs
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Bread, Coke, Diaper, Milk

Minimum Support = 3

Count and Elimination

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Beer, Bread, Diaper, Eggs
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Bread, Coke, Diaper, Milk



Items (1-itemsets)

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Minimum Support = 3

Generate 2-Itemset and prune

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset
{Bread, Milk}
{Bread, Beer}
{Bread, Diaper}
{Beer, Milk}
{Diaper, Milk}
{Beer, Diaper}

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

Count and Elimination

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Beer, Bread}	2
{Bread,Diaper}	3
{Beer,Milk}	2
{Diaper,Milk}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

Generate 3-Itemset and prune

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,
 ${}^6C_1 + {}^6C_2 + {}^6C_3$
 $6 + 15 + 20 = 41$
 With support-based pruning,
 $6 + 6 + 4 = 16$



Triplets (3-itemsets)

Itemset
{ Beer, Diaper, Milk }
{ Beer,Bread,Diaper }
{Bread, Diaper, Milk }
{ Beer, Bread, Milk }

Count and Elimination

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,
 ${}^6C_1 + {}^6C_2 + {}^6C_3$
 $6 + 15 + 20 = 41$
 With support-based pruning,
 $6 + 6 + 4 = 16$
 $6 + 6 + 1 = 13$



Triplets (3-itemsets)

Itemset	Count
{ Beer, Diaper, Milk}	
{ Beer,Bread, Diaper}	
{Bread, Diaper, Milk}	2
{Beer, Bread, Milk}	

More Efficient Candidate Generation

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,

$${}^6C_1 + {}^6C_2 + {}^6C_3 \\ 6 + 15 + 20 = 41$$

With support-based pruning,

$$6 + 6 + 1 = 13$$



Triplets (3-itemsets)

Itemset	Count
{Bread, Diaper, Milk}	2

How?

Market-Basket transactions

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Two-step approach:

- **Frequent Itemset Generation**
 - Generate all itemsets whose support $\geq \text{minsup}$
- **Rule Generation**
 - Generate high confidence ($\geq \text{minconf}$) rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset

Rule Generation

If $\{A,B,C,D\}$ is a frequent itemset, candidate rules:

$ABC \rightarrow D,$	$ABD \rightarrow C,$	$ACD \rightarrow B,$	$BCD \rightarrow A,$
$A \rightarrow BCD,$	$B \rightarrow ACD,$	$C \rightarrow ABD,$	$D \rightarrow ABC$
$AB \rightarrow CD,$	$AC \rightarrow BD,$	$AD \rightarrow BC,$	$BC \rightarrow AD,$
$BD \rightarrow AC,$	$CD \rightarrow AB,$		

If $|L| = k$, then there are $2^k - 2$ candidate association rules (ignoring $L \rightarrow \phi$ and $\phi \rightarrow L$)

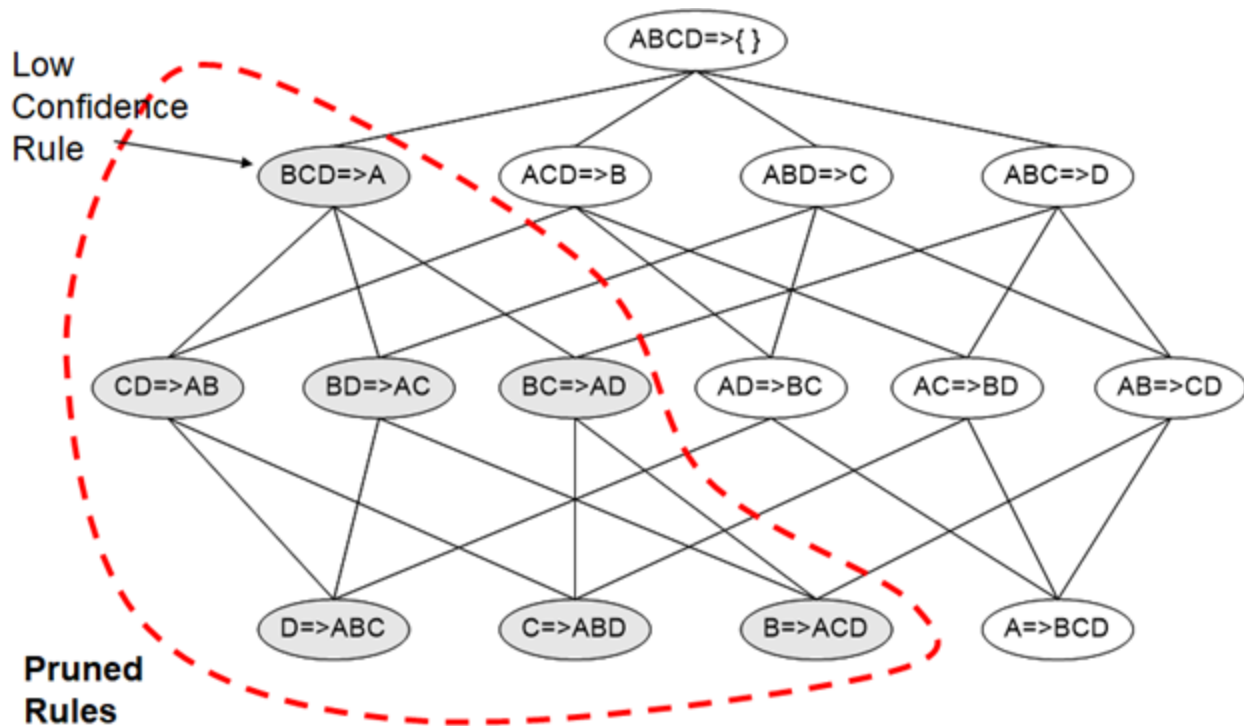
Rule Generation

- Confidence of rules generated from the same itemset has an anti-monotone property
 - E.g., Suppose $\{A,B,C,D\}$ is a frequent 4-itemset:

$$\begin{aligned}c(ABC \rightarrow D) &= \sigma(ABCD) / \sigma(ABC) \\&\geq c(AB \rightarrow CD) = \sigma(ABCD) / \sigma(AB) \\&\geq c(A \rightarrow BCD) = \sigma(ABCD) / \sigma(A)\end{aligned}$$

- Confidence is *anti-monotone* w.r.t. number of items on the RHS of the rule

Rule Generation



Rule Generation

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)



Triplets (3-itemsets)

Itemset	Count
{Bread, Diaper, Milk}	2

$Minconf = 0.8$

$c(\text{Bread} \rightarrow \text{Milk}) = 3 / 4$

$c(\text{Milk} \rightarrow \text{Bread}) = 3 / 4$

$c(\text{Bread} \rightarrow \text{Diaper}) = 3 / 4$

$c(\text{Diaper} \rightarrow \text{Bread}) = 3 / 4$

$c(\text{Milk} \rightarrow \text{Diaper}) = 3 / 4$

$c(\text{Diaper} \rightarrow \text{Milk}) = 3 / 4$

$c(\text{Beer} \rightarrow \text{Diaper}) = 3 / 3$

$c(\text{Diaper} \rightarrow \text{Beer}) = 3 / 4$

Apply Rules

Rules: {**Beer** \rightarrow **Diaper**}

- When a customer buys **Beer**, suggest **Diaper** also.
- Put **Beer** and **Diaper** close.