

association analysis

→ 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

* implication means cooccurrence, not causality

Market-Basket transactions

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

support count (σ) = frequency of occurrence of an itemset → $\sigma(\{milk, bread, diaper\}) = 2$

association rule = $X \rightarrow Y$ ^{item sets} ⇒ $\{milk, diaper\} \rightarrow \{beer\}$

support (s) = fraction of transactions that contain both X and Y ⇒ $s = \frac{\sigma(\{milk, bread, diaper\})}{ITI} = \frac{2}{5}$

confidence (c) = measure of how often Y appear in transactions that contain X ⇒ $c = \frac{\sigma(\{milk, bread, diaper\})}{\sigma(\{milk, bread\})} = \frac{2}{3}$

how to find rules? first find frequent itemset ($\geq \text{minsup}$), then generate rules from these that have high confidence

frequent itemset generation strategies

• reducing number of candidates = for d items, 2^d candidates

→ apriori principle = if an itemset is frequent, then all of its subsets must also be frequent

↳ anti-monotone property = support of an itemset never exceeds the support of its subsets

• if AB is not frequent, ABC , $ABCD$, $ABCDE$ → all its supersets are also not frequent (prune)

$\text{minsup} = 3$

if every subset is considered:

$$\binom{6}{1} + \binom{6}{2} + \binom{6}{3} = 11$$

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

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

Itemset	Count
{Beer, Diaper, Milk}	2
{Beer, Bread, Diaper}	2
{Bread, Diaper, Milk}	2
{Beer, Bread, Milk}	1

with support-based pruning, $6+6+1 = 13$ subsets are considered

→ candidate generation: $F_{k-1} \times F_{k-1}$ method

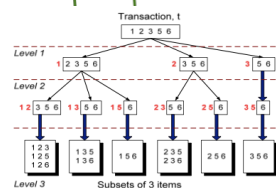
↳ Merge ($\underline{ABC}, \underline{ABD}$) = $ABCD$ merge if first k items are same

↳ do not merge ACD , ABD → must have same prefix

↳ alternate $F_{k-1} \times F_{k-1}$ method

↳ Merge ($\underline{ABC}, \underline{BCD}$) = $ABCD$ merge if last k items of first one are same with the first k items of second one

support counting of candidate items = instead of matching each transaction against every candidate, match with hash buckets



for transactions with length 3

→ {1, 3, 5} supported by transaction {1, 2, 3, 5, 6}?

rule generation = from frequent itemsets, find the ones that have confidence $\geq \text{minconf}$

↳ for $\{A, B, C, D\}$ ⇒ $A \rightarrow BCD$ $B \rightarrow ACD$... $\binom{4}{1} + \binom{4}{2} + \binom{4}{3} = 14$ possible rules $2^4 - 2 = 14$

* confidence does not have anti-monotone property ($A \rightarrow BCD$ and $AB \rightarrow D$ are independent) for different itemsets but have for same itemset

$$c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$$

↳ if this below the minimum confidence, no need to check rest, eliminate, prune

complexity of apriori = increases with decrease of support threshold

maximal frequent itemset = if none of itemset's immediate supersets (one level up) is frequent

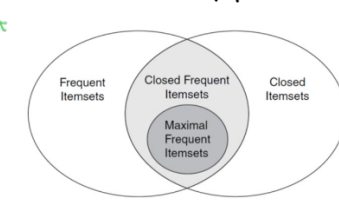
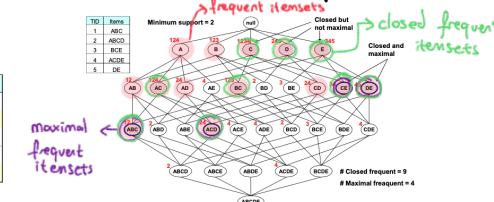
Transaction	A	B	C	D	E	F	G	H	I	J
1										
2										
3										
4										
5										
6										
7										
8										
9										
10										

support threshold =	5	4	3
frequent itemsets =	F	E, F, J, EF	C, D, D...
maximal itemsets =	F	EF, J	CDEF, J

closed itemset = if none of an itemsets immediate supersets has the same support as X

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,B,C,D}
4	{A,B,D}
5	{A,B,C,D}

Itemset	Support
{A}	4
{B}	5
{C}	3
{D}	4
{A,B}	4
{A,C}	2
{A,D}	3
{B,C}	3
{B,D}	4
{C,D}	3



Items	Support (count)	Closed Itemset
{C}	3	✓
{D}	2	✗
{E}	2	✗
{C,D}	2	✗
{C,E}	2	✗
{D,E}	2	✗
{C,D,E}	2	✓

interestiness measure = $X \rightarrow Y$

Contingency table	Y	\bar{Y}
X	f_{11}	f_{10}
\bar{X}	f_{01}	f_{00}
	f_{1+}	f_{0+}
	f_{+1}	f_{+0}
	N	

tea → coffee

	Coffee	$\bar{\text{Coffee}}$	
Tea	150	50	200
$\bar{\text{Tea}}$	650	150	800
	800	200	1000

confidence = $P(\text{coffee} | \text{tea}) = 150/200 = 0.75$

$$P(\text{coffee}) = 800/1000 = 0.8$$

contradiction

means knowing that a person drinks tea reduces the probability that he/she drinks coffee should not reduce

* confidence ($X \rightarrow Y$) > support(Y): otherwise rule will be misleading

$P(X, Y) > P(X) \times P(Y) \rightarrow X \& Y$ positively correlated

$P(X, Y) < P(X) \times P(Y) \rightarrow X \& Y$ negatively correlated

$P(X, Y) = P(X) \times P(Y) \rightarrow X \& Y$ are independent ⇒ confidence($X \rightarrow Y$) = support(Y) ⇒ $P(Y | X) = P(Y)$

lift = $\frac{P(Y | X)}{P(Y)} \rightarrow 1$ if they are independent → used to measure the importance of a rule

they are same if they are independent

interest = $\frac{P(X, Y)}{P(X) \cdot P(Y)} \rightarrow 1$ if they are independent

invariant measures to inverse = cosine, jaccard, confidence

non-invariant measures to inverse = correlation, interest/lift, odds ratio → they don't change

simpson's paradox = observed relationship in data may be influenced by the presence of hidden variables

↳ recovery rate in hospitals, hidden variable = young or old patients

cross support and H-confidence ⇒

conf(caviar → milk) → very high

conf(milk → caviar) → very low

$$0 \leq h\text{conf}(X) \leq r(X) \leq 1$$

h-confidence = min conf of any association rule formed from itemset $X \rightarrow s(X) / \max \{s(x_i)\}$

cross-support = $\min \{s(x_i)\} / \max \{s(x_j)\}$

hyperclique = items in the itemset are strongly correlated → not necessarily frequent itemsets