

Ve572 Lecture 10

Manuel and Jing

UM-SJTU Joint Institute

June 19, 2018

- **Association analysis**, a.k.a **Association rule mining/learning** is a method for discovering interesting relations between variables in large databases.
- It is often used by retailers, again your cat and dog, which has a very large database on transactions, known as **market basket data**, e.g.

TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

- Retailers are interested in analysing the data to learn about the purchasing behaviour of their customers, and make recommendations to their customers.

Q: Can you guess the difference between clustering and association analysis?

- Association analysis is often used in
 1. Bioinformatics
 2. Medical diagnosis
 3. Finance
 4. Government monitoring
- We discuss association analysis using market basket data. Let

$$\mathcal{I} = \{z_1, z_2, \dots, z_p\}$$

be the set of all items in a market basket data, and

$$\mathcal{T} = \{t_1, t_2, \dots, t_n\}$$

be the set of all transactions, where each transaction

$$t_i$$

has a unique transaction ID, and contains a subset of items from \mathcal{I} .

- Imagine we have a very big market basket data, that is, very large p and n in

$$\mathcal{I} = \{z_1, z_2, \dots, z_p\} \quad \text{and} \quad \mathcal{T} = \{t_1, t_2, \dots, t_n\}$$

Q: How would you make recommendations? That is, how to expand the idea

	Gillette Razors	Christian Louboutin Shoes	Air Jordan Shoes
obs. 1	1	0	1
obs. 2	0	1	0
obs. 1	1	0	1

to slightly more complicated data but vastly big

TID	Bread	Milk	Diapers	Beer	Egg	Cola
1	1	1	0	0	0	0
2	1	0	1	1	1	0
3	0	1	1	1	0	1
4	1	1	1	1	0	0
5	1	1	1	0	0	1

where 1 denotes at least 1 of that item was brought in the i th transaction.

- In association analysis, a subset of \mathcal{I} , $\mathcal{X} \subset \mathcal{I}$, is known as an **itemset**, e.g.

$$\mathcal{A} = \{ \text{Milk, Diapers, Beer} \}$$

- The **support** of an itemset \mathcal{X} is defined as

$$\text{supp}(\mathcal{X}) = \frac{|\{t_i \mid \mathcal{X} \subset t_i, t_i \in \mathcal{T}\}|}{|\mathcal{T}|}$$

where $|\cdot|$ denotes the number of elements in a set, the **size** of the set.

- For instance, the support of \mathcal{A} is 0.4,

TID	Bread	Milk	Diapers	Beer	Egg	Cola
1	1	1	0	0	0	0
2	1	0	1	1	1	0
3	0	1	1	1	0	1
4	1	1	1	1	0	0
5	1	1	1	0	0	1

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

$$\mathcal{X} \implies \mathcal{Y}$$

where \mathcal{X} and \mathcal{Y} are itemsets of \mathcal{I} , and are disjoint

$$\mathcal{X} \cap \mathcal{Y} = \emptyset$$

- The **confidence** of a rule is defined as

$$\text{conf}(\mathcal{X} \implies \mathcal{Y}) = \frac{\text{supp}(\mathcal{X} \cup \mathcal{Y})}{\text{supp}(\mathcal{X})}$$

Q: What is the confidence of the rule $\{\text{Milk, Diapers}\} \implies \{\text{Beer}\}$?

TID	Bread	Milk	Diapers	Beer	Egg	Cola
1	1	1	0	0	0	0
2	1	0	1	1	1	0
3	0	1	1	1	0	1
4	1	1	1	1	0	0
5	1	1	1	0	0	1

- The support of an itemset \mathcal{X} is an indication of how frequently \mathcal{X} in \mathcal{T} .

Q: What can we interpret the confidence of a rule

$$\text{conf}(\mathcal{X} \implies \mathcal{Y}) = \frac{\text{supp}(\mathcal{X} \cup \mathcal{Y})}{\text{supp}(\mathcal{X})}$$

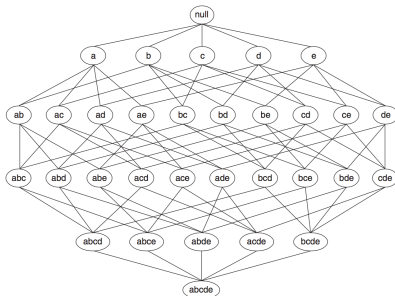
- Of course, we wish to discover rules with high support and high confidence.
- Similar to hierarchical clustering in the sense that the number of rules satisfy

$$\text{supp}(\mathcal{X} \cup \mathcal{Y}) \geq s_{min} \quad \text{and} \quad \text{conf}(\mathcal{X} \implies \mathcal{Y}) \geq c_{min}$$

decreases as s_{min} and c_{min} increases, just like the number of clusters and d.

- A choice then can be made regarding the recommendations for customers once we have a set of rules with the highest support and confidence.
- Association rule mining is about finding all rules achieve certain s_{min} or c_{min}
- A brutal force approach is NOT a good idea since it grows really quickly.

- We will have a lot of rules with only 6 items,



- In practice, when apply the support requirement, a lot of rules are out, e.g.

$$\begin{aligned} \{ \text{Beer, Diapers} \} &\implies \{ \text{Milk} \}, \{ \text{Milk} \} \implies \{ \text{Beer, Diapers} \}, \\ \{ \text{Diapers, Milk} \} &\implies \{ \text{Beer} \}, \{ \text{Beer} \} \implies \{ \text{Diapers, Milk} \}, \\ \{ \text{Beer, Milk} \} &\implies \{ \text{Diapers} \}, \{ \text{Diapers} \} \implies \{ \text{Beer, Milk} \}, \end{aligned}$$

all involve the same $\text{supp}(\mathcal{X} \cup \mathcal{Y})$, where $\mathcal{X} \cup \mathcal{Y} = \{ \text{Milk, Diapers, Beer} \}$

- Hence a common strategy used by many association rule mining algorithms is to decompose the problem into the following sequential steps:

1. Frequent Itemset Generation

$$\text{supp}(\mathcal{X} \cup \mathcal{Y}) \geq s_{min}$$

2. Strong Rule Generation

$$\text{conf}(\mathcal{X} \implies \mathcal{Y}) \geq c_{min}$$

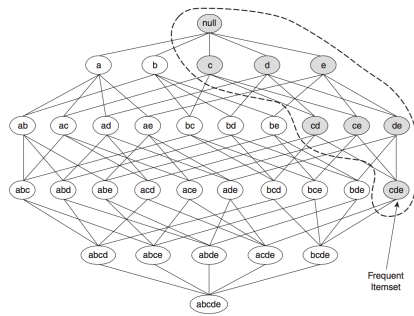
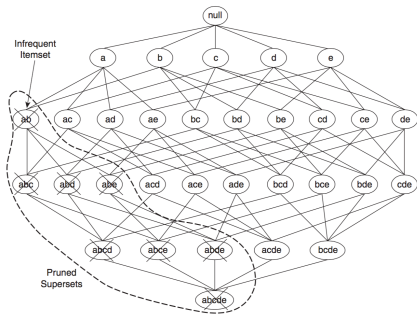
- Notice the definition of support

$$\text{supp}(\mathcal{X}) = \frac{|\{t_i \mid \mathcal{X} \subset t_i, t_i \in \mathcal{T}\}|}{|\mathcal{T}|}$$

ensures the following key inequality is true

$$\text{supp}(\mathcal{X} \cup \mathcal{Y}) \leq \text{supp}(\mathcal{X})$$

- This means if $\{a, b\}$ is infrequent, then all the supersets are infrequent,



and if $\{c, d, e\}$ is frequent, then all the subsets are frequent.

1. This leads to the sub-steps in the Frequent Itemset Generation step in the

Apriori algorithm

- Roughly, this step consists of two iterative sub-steps:

- (a) Candidate generation
- (b) Candidate pruning

- To give a high-level illustration, consider the following small dataset again

TID	Bread	Milk	Diapers	Beer	Egg	Cola
1	1	1	0	0	0	0
2	1	0	1	1	1	0
3	0	1	1	1	0	1
4	1	1	1	1	0	0
5	1	1	1	0	0	1

- Note we want to compute/check the support, however $|\mathcal{T}|$ is a constant and

$$\sigma(\mathcal{A}) = |\mathcal{T}| \cdot \text{supp}(\mathcal{A}) = |\{t_i \mid \mathcal{A} \subset t_i, t_i \in \mathcal{T}\}| = \text{counts}$$

thus it is just a matter of converting this table in wide form into a long form.

(a) Candidate generation:

Generate all itemset \mathcal{A} of size 1, that is $|\mathcal{A}| = 1$,

Itemset	Count
Bread	4
Milk	4
Diapers	4
Beer	3
Egg	1
Cola	2

- Suppose $s_{min} = 60\%$, which is equivalent to a minimum count equal to 3.

(b) Candidate pruning:

Remove itemsets with a low count, that is, we left with

Itemset	Count
Bread	4
Milk	4
Diapers	4
Beer	3

- Next iteration, we repeat sub-step (a) and (b), but with itemset of \mathcal{A} of size

$$|\mathcal{A}| = 2$$

- Invoking the key inequality,

$$\text{supp}(\mathcal{X} \cup \mathcal{Y}) \leq \text{supp}(\mathcal{X})$$

the relevant itemset are the ones that do not involve Egg or Cola

TID	Bread	Milk	Diapers	Beer	Egg	Cola
1	1	1	0	0	0	0
2	1	0	1	1	1	0
3	0	1	1	1	0	1
4	1	1	1	1	0	0
5	1	1	1	0	0	1

- Thus, we have $\{ \text{Bread, Milk} \}, \{ \text{Bread, Diapers} \}, \{ \text{Bread, Beer} \},$
 $\{ \text{Milk, Diapers} \}, \{ \text{Milk, Beer} \},$
 $\{ \text{Diapers, Beer} \},$

(a) Candidate generation:

Itemset	Count
Bread, Milk	3
Bread, Diapers	3
Bread, Beer	2
Milk, Diapers	3
Milk, Beer	2
Diapers, Beer	3

(b) Candidate pruning:

Itemset	Count
Bread, Milk	3
Bread, Diapers	3
Milk, Diapers	3
Diapers, Beer	3

- Similarly, in the next iteration, we have

Itemset	Count
Bread, Milk, Diapers	3

- Notice every itemset of size 3 involving Beer is a superset of an infrequent itemset in the previous two iterations.

TID	Bread	Milk	Diapers	Beer	Egg	Cola	Itemset	Count
							Bread, Milk	3
1	1	1	0	0	0	0	Bread, Diapers	3
2	1	0	1	1	1	0	Bread, Beer	2
3	0	1	1	1	0	1	Milk, Diapers	3
4	1	1	1	1	0	0	Milk, Beer	2
5	1	1	1	0	0	1	Diapers, Beer	3

- In general, the candidate generation sub-step in Apriori algorithm merges a pair of itemsets of size $(k - 1)$ only if their first $(k - 2)$ items are identical.

- Let

$$\mathcal{A} = \{a_1, a_2, \dots, a_{k-1}\} \quad \text{and} \quad \mathcal{B} = \{b_1, b_2, \dots, b_{k-1}\}$$

be a pair of frequent itemsets generated in the $(k-1)$ th iteration.

- In the k th iteration, \mathcal{A} and \mathcal{B} are merged

$$\mathcal{A} \cup \mathcal{B}$$

to form an itemset of size k if the following conditions are satisfied

$$\begin{aligned} a_i &= b_i & \text{for } i &= 1, 2, \dots, k-2 \\ a_{k-1} &= b_{k-1} \end{aligned}$$

- Iterating until no such pair of frequent itemsets exists, which is the end of

1. Frequent Itemset Generation

- (a) Candidate generation
- (b) Candidate pruning

- Now consider rule generation given a frequent itemset $\mathcal{A} = \mathcal{X} \cup \mathcal{Y}$ of size k

$$\text{conf}(\mathcal{X} \implies \mathcal{Y}) \geq c_{min}$$

- For example, we have the following as one of the frequent itemsets

$$\mathcal{A} = \{ \text{Bread, Milk, Diapers} \}$$

- There are six candidate association rules that can be generated from \mathcal{A} :

$$\{ \text{Bread, Milk} \} \implies \{ \text{Diapers} \}$$

$$\{ \text{Bread, Diapers} \} \implies \{ \text{Milk} \}$$

$$\{ \text{Milk, Diapers} \} \implies \{ \text{Bread} \}$$

$$\{ \text{Bread} \} \implies \{ \text{Milk, Diapers} \}$$

$$\{ \text{Milk} \} \implies \{ \text{Bread, Diapers} \}$$

$$\{ \text{Diapers} \} \implies \{ \text{Bread, Milk} \}$$

Q: How can we efficiently generate and prune rules given a frequent itemset?

- Note if a rule $\mathcal{C} \implies \mathcal{A} - \mathcal{C}$, where $\mathcal{C} \subset \mathcal{A}$, is low confidence rule, that is,

$$\text{conf}(\mathcal{C} \implies \mathcal{A} - \mathcal{C}) < c_{min}$$

then any rule $\mathcal{C}^* \implies \mathcal{A} - \mathcal{C}^*$, where $\mathcal{C}^* \subset \mathcal{C}$, is also low confident,

$$\text{conf}(\mathcal{C}^* \implies \mathcal{A} - \mathcal{C}^*) < c_{min}$$

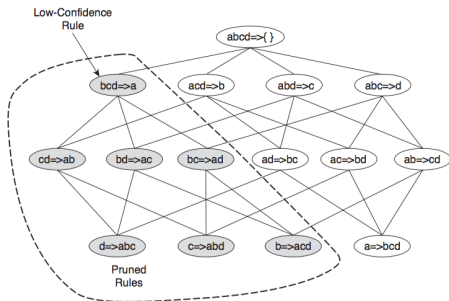
- This is clear from the definition of

$$\begin{aligned}\text{conf}(\mathcal{C} \implies \mathcal{A} - \mathcal{C}) &= \frac{\text{supp}(\mathcal{C} \cup (\mathcal{A} - \mathcal{C}))}{\text{supp}(\mathcal{C})} = \frac{\text{supp}(\mathcal{A})}{\text{supp}(\mathcal{C})} \\ \text{conf}(\mathcal{C}^* \implies \mathcal{A} - \mathcal{C}^*) &= \frac{\text{supp}(\mathcal{C}^* \cup (\mathcal{A} - \mathcal{C}^*))}{\text{supp}(\mathcal{C}^*)} = \frac{\text{supp}(\mathcal{A})}{\text{supp}(\mathcal{C}^*)}\end{aligned}$$

Since $\mathcal{C}^* \subset \mathcal{C}$, we have

$$\text{supp}(\mathcal{C}^*) \geq \text{supp}(\mathcal{C})$$

- Thus we should start with the largest possible \mathcal{C} for a given frequent itemset,



and precede to the next level only if both rules satisfy

$$\begin{aligned} \text{conf}(\mathcal{C}_1 \Rightarrow \mathcal{A} - \mathcal{C}_1) < c_{min} \\ \text{conf}(\mathcal{C}_2 \Rightarrow \mathcal{A} - \mathcal{C}_2) < c_{min} \end{aligned} \implies \text{conf}(\mathcal{C}_1 \cap \mathcal{C}_2 \Rightarrow \mathcal{A} - \mathcal{C}_1 \cap \mathcal{C}_2)$$

- Rules from 2 itemsets are not related, so are generated and pruned separately