

1 September 2nd, 2019

1.1 Association Rule Mining

Suppose we have the following dataset:

Customer	Shopping List		
Raymond	Apple	Coke	Coffee
David	Diaper	Coke	
Emily	Milk	Biscuit	
Derek	Coke	Milk	

Definition 1.1. The things on the RHS are **items**.

Definition 1.2. For each customer we have their **history** or **transaction**.

We want to find some **associations** between items. An example of an interesting association might be:

Example 1.3 (Example of an interesting association)

Diapers and Beers are usually bought together.

This association could have different reasons, e.g. people buy both diapers and beer after work usually.

1.2 Applications of Association Rule Mining

Here are some examples of where association rule mining might be used:

- Supermarket - For recommendation
- Web Mining - Google for their autocomplete
- Medical Analysis - Diagnosis from the patient's attributes or finding key attributes linked to illnesses (diabetes and obesity)
- Bioinformatics - Patterns in genomes
- Network Analysis - Associating IP and DoS, e.g. seeing if your packet goes through
- Programming Pattern Finding - e.g. linking segmentation faults and users

1.3 Problem Definition

Consider the following dataset (TID = Transaction ID):

- TID: t_1 , Items: A, D
- TID: t_2 , Items: A, B, D, E

- TID: t_3 , Items: B, C
- TID: t_4 , Items: A, B, C, D, E
- TID: t_5 , Items: B, C, E

In table form this would be:

TID	A	B	C	D	E
t_1	1	0	0	1	0
t_2	1	1	0	1	1
t_3	0	1	1	0	0
t_4	1	1	1	1	1
t_5	0	1	1	0	1

Definition 1.4. A **single item** is a single item (duh).

Example 1.5 (Examples of Single Items)

A, B, C, D , or E

Definition 1.6. An **itemset** is a set of items (again, duh).

Example 1.7 (Examples of Itemsets)

$\{B, C\}, \{A, B, C\}, \{B, C, D\}, \{A\}$

Definition 1.8. An **n -itemset** is a set of n items.

Example 1.9 (Examples of n -itemsets)

From Example 1.7, we have the following:

- $\{B, C\}$ is a 2-itemset
- $\{A, B, C\}$ and $\{B, C, D\}$ are 3-itemsets
- $\{A\}$ is a 1-itemset

Definition 1.10. The **support** (or **frequency**) of an item or an itemset is the number of times it appears in the dataset.

Example 1.11 (Examples of the Support of Items and Itemsets)

From Example 1.7, we have the following:

- The support of A is 3: $A \in t_1, t_2, t_4$
- The support of B is 4: $B \in t_2, t_3, t_4, t_5$
- The support of $\{B, C\}$ is 3: $\{B, C\} \subseteq t_3, t_4, t_5$
- The support of $\{A, B, C\}$ is 3: $\{A, B, C\} \subseteq t_4$

As such, we might try to classify large itemsets or **frequent itemsets** as itemsets with support greater than a threshold, e.g. 3.

Definition 1.12. An n -frequent itemset is an itemset with support n .

Example 1.13

$\{B, C\}$ is a 3-frequent itemset of size 2.

Definition 1.14. An **association rule** is a association between an item/itemset and another.

Definition 1.15. The **support** of an association rule is the number of transaction with both the LHS and RHS of the association rule.

Definition 1.16. The **confidence** of an association rule is the support of the association rule divided by the number of transaction with the LHS of the rule.

Example 1.17

$\{B, C\} \rightarrow E$ is an example of an association rule. It has a support of 3 (t_3, t_4, t_5) and a confidence of $\frac{2}{3}$ (it's true for t_3 and t_4 but not t_5).

In essence, we want to find association rules with:

- Support greater than a threshold e.g. (≥ 3)
- Confidence greater than a threshold e.g. ($\geq 50\%$)

We can do split this into two steps:

1. Find all “large” itemsets (e.g. itemsets with support ≥ 3)
2. Find all “interesting” association rules after Step 1:
 - From all “large” itemsets, find the association rules with confidence $\geq 50\%$.
 - This can be done by taking every pair of elements from Step 1, X and Y , where $X \subset Y$, and checking if $\frac{\text{supp}(Y)}{\text{supp}(X)} \geq 50\%$.
 - If yes, then generate the rule: “ $X \rightarrow Y - X$ ”

Homework

Show that the support of the association rule is still large. This can be easily seen, as X is large, and $Y - X \subseteq Y$, making it large from ??