

Association Rule Mining.

Trace the results of using the Apriori algorithm on the grocery store example with support threshold $s=33.34\%$ and confidence threshold $c=60\%$. Show the candidate and frequent itemsets for each database scan. Enumerate all the final frequent itemsets. Also indicate the association rules that are generated and highlight the strong ones, sort them by confidence.

Transaction ID	Items
T1	HotDogs, Buns, Ketchup
T2	HotDogs, Buns
T3	HotDogs, Coke, Chips
T4	Chips, Coke
T5	Chips, Ketchup
T6	HotDogs, Coke, Chips

Solution:

Support threshold $=33.34\% \Rightarrow$ threshold is at least 2 transactions.

Applying Apriori

Pass (k)	Candidate k-itemsets and their support	Frequent k-itemsets
k=1	HotDogs(4), Buns(2), Ketchup(2), Coke(3), Chips(4)	HotDogs, Buns, Ketchup, Coke, Chips
k=2	{HotDogs, Buns}(2), {HotDogs, Ketchup}(1), {HotDogs, Coke}(2), {HotDogs, Chips}(2), {Buns, Ketchup}(1), {Buns, Coke}(0), {Buns, Chips}(0), {Ketchup, Coke}(0), {Ketchup, Chips}(1), {Coke, Chips}(3)}	{HotDogs, Buns}, {HotDogs, Coke}, {HotDogs, Chips}, {Coke, Chips}
k=3	{HotDogs, Coke, Chips}(2)	{HotDogs, Coke, Chips}
k=4	{}	

Note that {HotDogs, Buns, Coke} and {HotDogs, Buns, Chips} are not candidates when $k=3$ because their subsets {Buns, Coke} and {Buns, Chips} are not frequent.

Note also that normally, there is no need to go to $k=4$ since the longest transaction has only 3 items.

All Frequent Itemsets: {HotDogs}, {Buns}, {Ketchup}, {Coke}, {Chips}, {HotDogs, Buns}, {HotDogs, Coke}, {HotDogs, Chips}, {Coke, Chips}, {HotDogs, Coke, Chips}.

Association rules:

{HotDogs, Buns} would generate: HotDogs \rightarrow Buns ($2/6=0.33$, $2/4=0.5$) and
Buns \rightarrow HotDogs ($2/6=0.33$, $2/2=1$);
 {HotDogs, Coke} would generate: HotDogs \rightarrow Coke (0.33 , 0.5) and
Coke \rightarrow HotDogs ($2/6=0.33$, $2/3=0.66$);
 {HotDogs, Chips} would generate: HotDogs \rightarrow Chips (0.33 , 0.5) and
 Chips \rightarrow HotDogs ($2/6=0.33$, $2/4=0.5$);
 {Coke, Chips} would generate: **Coke \rightarrow Chips ($3/6=0.5$, $3/3=1$) and**
Chips \rightarrow Coke ($3/6=0.5$, $3/4=0.75$);
 {HotDogs, Coke, Chips} would generate: HotDogs \rightarrow Coke \wedge Chips ($2/6=0.33$, $2/4=0.5$),
Coke \rightarrow Chips \wedge HotDogs ($2/6=0.33$, $2/3=0.66$),
 Chips \rightarrow Coke \wedge HotDogs ($2/6=0.33$, $2/4=0.5$),
HotDogs \wedge Coke \rightarrow Chips ($2/6=0.33$, $2/2=1$),
HotDogs \wedge Chips \rightarrow Coke ($2/6=0.33$, $2/2=1$) and
Coke \wedge Chips \rightarrow HotDogs ($2/6=0.33$, $2/3=0.66$).

With the confidence threshold set to 60%, the Strong Association Rules are (sorted by confidence):

- | | |
|---|---|
| 1. Coke \rightarrow Chips (0.5, 1) | 5. Chips \rightarrow Coke (0.5, 0.75); |
| 2. Buns \rightarrow HotDogs (0.33, 1); | 6. Coke \rightarrow HotDogs (0.33, 0.66); |
| 3. HotDogs \wedge Coke \rightarrow Chips(0.33, 1) | 7. Coke \rightarrow Chips \wedge HotDogs (0.33, 0.66) |
| 4. HotDogs \wedge Chips \rightarrow Coke(0.33, 1) | 8. Coke \wedge Chips \rightarrow HotDogs(0.33, 0.66). |