732A61/TDDD41 Data Mining - Clustering and Association Analysis
Lecture 6: Association Analysis I

Jose M. Peña
IDA, Linköping University, Sweden

# Outline

- Content
  - Association Rules
  - Frequent Itemsets
  - Apriori Algorithm
  - Exercise
  - Rule Generation Algorithm
  - Exercise
  - Summary

- Literature
  - Course book. Second edition: 5.2.1-2, 5.4. Third edition: 6.2.1-2, 6.4.
  - Agrawal, R. and Srikant, R. Fast Algorithms for Mining Association Rules. In Proc. of the 20th Int. Conf. on Very Large Databases, 1994. Extended version available as IBM Research Report RJ9839, 1994.

# Association Rules

- Assume that we have access to some transactional data, e.g.

| Transaction id | Items bought |
|---|---|
| 1 | A, B, D |
| 2 | A, C, D |
| 3 | A, D, E |
| 4 | B, E, F |
| 5 | B, C, D, E, F |

- We are interested in finding rules of the form

$$Item_1, \ldots, Item_m \rightarrow Item_{m+1}, \ldots, Item_n$$

meaning that if the items in the antecedent are purchased, so are the items in the consequent, e.g.

$$milk, eggs \rightarrow bread, butter$$

- Application: Market basket analysis to support business decisions, e.g.
    - Rules with "butter" in the consequent may help to decide how to boost sales of "butter".
    - Rules with "eggs" in the antecedent may help to determine what happens if "eggs" are sold out.
- Note however that the rules do not convey causality, i.e. forcing the antecedent does not guarantee the consequent.

## Association Rules

▸ We are interested in finding rules of the form

$$X_1, \ldots, X_m \to Y_1, \ldots, Y_n \equiv X \to Y$$

▸ However, not all the rules are equally interesting.

▸ We are interested in finding rules with user-defined minimum support and confidence, where

- ▸ Support = fraction of the transactions which contain $X$ and $Y$ = $p(X, Y)$.
- ▸ Support = how general the rule is.
- ▸ Confidence = fraction of the transactions that contain $X$ which also contain $Y$ = $p(Y|X)$.
- ▸ Confidence = how accurate the rule is.
- ▸ Confidence = $p(Y|X) = p(X, Y)/p(X)$ = support$(X, Y)$ / support$(X)$.

▸ Assume the following transactional data.

| Transaction id | Items bought |
|----------------|--------------|
| 1              | A, B, D      |
| 2              | A, C, D      |
| 3              | A, D, E      |
| 4              | B, E, F      |
| 5              | B, C, D, E, F |

▸ $A \to D$ has support 0.6 and confidence 1.

▸ $D \to A$ has support 0.6 and confidence 0.75.

# Frequent Itemsets

- We are interested in finding rules of the form

$$X_1, \ldots, X_m \to Y_1, \ldots, Y_n \equiv X \to Y$$

  with user-defined minimum support and confidence.

- We define a frequent or large itemset as a set of items that has minimum support.
  - E.g., $\{A, D\}$ is a frequent itemset in the previous example when the minimum support is 0.5.

- We will find the desired rules in two steps:
  1. Find all the frequent itemsets (via the apriori or FP grow algorithm).
  2. Generate all the rules with minimum confidence from the frequent itemsets.

- The first step above will make use of the following apriori property:
  - Every subset of a frequent itemset is frequent.
  - Or, alternatively, every superset of an infrequent itemset is infrequent.
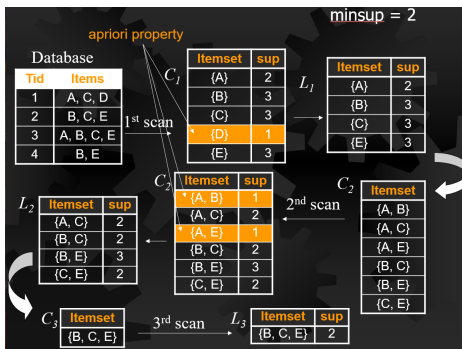
# Apriori Algorithm

**Algorithm**: apriori($D$, *minsup*)
**Input**: A transactional database $D$ and the minimum support *minsup*.
**Output**: All the large itemsets in $D$.

```
1   L₁ = { large 1-itemsets }
2   for (k = 2; L_{k-1} ≠ ∅; k + +) do
3       C_k = apriori-gen(L_{k-1})        // Generate candidate large k-itemsets
4       for all t ∈ D do
5           for all c ∈ C_k such that c ∈ t do
6               c.count + +
7       L_k = {c ∈ C_k | c.count ≥ minsup}
8   return ∪_k L_k
```



minsup = 2

apriori property

Database

| Tid | Items |
|-----|-------|
| 1 | A, C, D |
| 2 | B, C, E |
| 3 | A, B, C, E |
| 4 | B, E |

$C_1$

| Itemset | sup |
|---------|-----|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {D} | 1 |
| {E} | 3 |

$L_1$

| Itemset | sup |
|---------|-----|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {E} | 3 |

$C_2$

| Itemset | sup |
|---------|-----|
| {A, B} | 1 |
| {A, C} | 2 |
| {A, E} | 1 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

$C_2$

| Itemset |
|---------|
| {A, B} |
| {A, C} |
| {A, E} |
| {B, C} |
| {B, E} |
| {C, E} |

$L_2$

| Itemset | sup |
|---------|-----|
| {A, C} | 2 |
| {B, C} | 2 |
| {B, E} | 3 |
| {C, E} | 2 |

$C_3$

| Itemset |
|---------|
| {B, C, E} |

$L_3$

| Itemset | sup |
|---------|-----|
| {B, C, E} | 2 |

1st scan

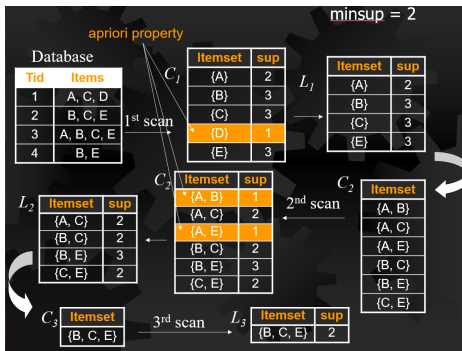2nd scan

3rd scan

# Apriori Algorithm

**Algorithm**: apriori-gen($L_{k-1}$)
**Input**: Large $(k-1)$-itemsets.
**Output**: A superset of $L_k$.

```
1   C_k = ∅                              // Self-join
2   for all I, J ∈ L_{k-1} do
3       if I_1 = J_1, …, I_{k-2} = J_{k-2} and I_{k-1} < J_{k-1} then
4           add {I_1, …, I_{k-1}, J_{k-1}} to C_k
5   for all c ∈ C_k do                   // Prune
6       for all (k-1)-subsets s of c do
7           if s ∉ L_{k-1} then
8               remove c from C_k
9   return C_k
```

# Apriori Algorithm
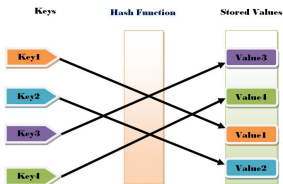
- Self-join step in MySQL:

  insert into $C_k$
  select $I.item_1, \ldots, I.item_{k-1}, J.item_{k-1}$
  from $L_{k-1}$ $I$, $L_{k-1}$ $J$
  where $I.item_1 = J.item_1, \ldots, I.item_{k-2} = J.item_{k-2}, I.item_{k-1} < J.item_{k-1}$

- Self-join step in R:

  $$\text{merge}(L_{k-1}, L_{k-1}, \text{by}=c(L_{k-1}.item_1, \ldots, L_{k-1}.item_{k-2}))$$

  but note that duplicates will be produced because the condition $I.item_{k-1} < J.item_{k-1}$ is not enforced.

- To make the prune step fast, large itemsets are usually stored in a hash table.



- Clever data structures are also typically used for counting the support of the candidates, i.e. lines 4-6 in the apriori algorithm.

# Exercise

- Run the apriori algorithm on the database below with minimum support 0.4, i.e. two transactions.

| Tid | A | B | C | D | E |
|-----|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 0 | 0 |
| 2 | 1 | 1 | 1 | 1 | 1 |
| 3 | 1 | 0 | 1 | 1 | 0 |
| 4 | 1 | 0 | 1 | 1 | 1 |
| 5 | 1 | 1 | 1 | 1 | 0 |

- Show the execution details (i.e. self-join, prune, support counting), not just the large itemsets
$\{A, B, C, D, E,$
$AB, AC, AD, AE, BC, BD, CD, CE, DE,$
$ABC, ABD, ACD, ACE, ADE, BCD, CDE,$
$ABCD, ACDE\}$.

# Apriori Algorithm

```
Algorithm: apriori-gen(L_{k-1})
Input: Large (k − 1)-itemsets.
Output: A superset of L_k.

1   C_k = ∅                                    // Self-join
2   for all I, J ∈ L_{k-1} do
3       if I_1 = J_1, ..., I_{k-2} = J_{k-2} and I_{k-1} < J_{k-1} then
4           add {I_1, ..., I_{k-1}, J_{k-1}} to C_k
5   for all c ∈ C_k do                         // Prune
6       for all (k − 1)-subsets s of c do
7           if s ∉ L_{k-1} then
8               remove c from C_k
9   return C_k
```

- We prove by induction on $k$ that the apriori algorithm is correct.
- Trivial case: The algorithm is correct for $k = 1$.
- Induction hypothesis: Assume that the algorithm is correct up to $k - 1$. We now prove that the algorithm is correct for $k$. It suffices to prove that $L_k \subseteq C_k$.
- Assume to the contrary that $I \in L_k$ but $I \notin C_k$. Then,
    - $\{I_1, \ldots, I_{k-2}, I_{k-1}\} \in L_{k-1}$ follows from $I \in L_k$ by the apriori property and the induction hypothesis.
    - $\{I_1, \ldots, I_{k-2}, I_k\} \in L_{k-1}$ follows from $I \in L_k$ by the apriori property and the induction hypothesis.
    - Then, $I \in C_k$ in line 5, i.e. it is generated by the self-join step.
    - Moreover, every subset of $I$ is large by $I \in L_k$ and the apriori property.
    - Then, $I \in C_k$ in line 9, i.e. it is not removed by the prune step.
    - This contradicts our assumption and, thus, the algorithm is correct for $k$.

## Rule Generation Algorithm

- We want to generate all the rules of the form

$$X \to L \smallsetminus X$$

where $L$ is a large itemset, $X \subseteq L$, and the rule has minimum confidence.

- We will make use of the following apriori property:

  - If $X$ does not result in a rule with minimum confidence for $L$, neither does any subset of $X$.

---

```
1   for all large itemsets l_k with k ≥ 2 do
2       call genrules(l_k, l_k, minconf)

    Algorithm: genrules(l_k, a_m, minconf)
    Input: A large itemset l_k, a set a_m ⊆ l_k, the minimum confidence minconf.
    Output: All the rules of the form a → l_k ∖ a with a ⊆ a_m and confidence equal or above minconf.

1   𝔸 = {(m − 1)-itemsets a_{m−1}|a_{m−1} ⊆ a_m}
2   for all a_{m−1} ∈ 𝔸 do
3       conf = support(l_k) / support(a_{m−1})
4       if conf ≥ minconf then
5           output the rule a_{m−1} → l_k ∖ a_{m−1} with confidence=conf and support=support(l_k)
6           if m − 1 > 1 then call genrules(l_k, a_{m−1}, minconf)
```

---

- A faster algorithm exists.

# Exercise

- Run the genrule algorithm on the database below for the large itemset $\{A, B, C\}$ with minimum confidence 0.8.

| Tid | A | B | C | D | E |
|-----|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 0 | 0 |
| 2 | 1 | 1 | 1 | 1 | 1 |
| 3 | 1 | 0 | 1 | 1 | 0 |
| 4 | 1 | 0 | 1 | 1 | 1 |
| 5 | 1 | 1 | 1 | 1 | 0 |

- Show the execution details (i.e. antecedent generation, recursive calls), not just the rules $\{AB \to C, BC \to A, B \to AC\}$.

# Rule Generation Algorithm

```
1   for all large itemsets l_k with k ≥ 2 do
2       call genrules(l_k, l_k, minconf)

    Algorithm: genrules(l_k, a_m, minconf)
    Input: A large itemset l_k, a set a_m ⊆ l_k, the minimum confidence minconf.
    Output: All the rules of the form a → l_k ∖ a with a ⊆ a_m and confidence equal or above minconf.

1   A = {(m − 1)-itemsets a_{m−1}|a_{m−1} ⊆ a_m}
2   for all a_{m−1} ∈ A do
3       conf = support(l_k) / support(a_{m−1})
4       if conf ≥ minconf then
5           output the rule a_{m−1} → l_k ∖ a_{m−1} with confidence=conf and support=support(l_k)
6           if m − 1 > 1 then call genrules(l_k, a_{m−1}, minconf)
```

- We prove by contradiction that the rule generation algorithm is correct.
- Assume to the contrary that the algorithm missed a rule. Let
  $a_{m−1} → l_k ∖ a_{m−1}$ denote one of the missing rules with the largest antecedent. Then,
    - Note that $l_k$ has minimum support and, thus, it is outputted by the apriori algorithm since this is correct.
    - Then, the algorithm cannot have missed the rule if $m = k$.
    - Moreover if $m < k$, then
      confidence($a_m → l_k ∖ a_m$) = support($l_k$) / support($a_m$)
      ≥ support($l_k$) / support($a_{m−1}$) = confidence($a_{m−1} → l_k ∖ a_{m−1}$)
      ≥ minconf
    - Note that the algorithm did not miss the rule $a_m → l_k ∖ a_m$ because, otherwise, it would contradict our assumption.
    - Then, the algorithm cannot have missed the rule $a_{m−1} → l_k ∖ a_{m−1}$.
    - This contradicts our assumption and, thus, the algorithm is correct.

# Summary

- Mining transactions to find rules of the form

$$Item_1, \ldots, Item_m \to Item_{m+1}, \ldots, Item_n$$

  with user-defined minimum support and confidence.
- Two-step solution:
    1. Find all the large itemsets (via the apriori algorithm).
    2. Generate all the rules with minimum confidence from the large itemsets.
- The two steps above make use of apriori properties.
- Drawbacks of the apriori algorithm:
    - Candidate generate-and-test.
    - Too many candidates to generate, e.g. if there are $10^4$ large 1-itemsets, then more than $10^7$ candidate 2-itemsets.
    - Each candidate implies expensive operations, e.g. pattern matching, subset checking, storing.
- Can candidate generation be avoided ? Yes, use the FP grow algorithm.