

17 Jan 2019

Frequent Itemsets - Market Basket Analysis

$$I = \{i_1, i_2, \dots, i_N\}$$

$$T = \{t_1, t_2, \dots, t_M\}$$

Each $t_i \subseteq I$ is an item set

Given $X \subseteq I$, support of X is fraction of T containing X

$$\text{sup}(X) = \frac{|\{t_i \mid X \subseteq t_i\}|}{M}$$

Frequent : $\text{Sup}(X) \geq \text{threshold}$
min-sup

Examples

1. Retail - shopping baskets
2. Text - items are words, baskets are documents
3. Subjects in which schoolchildren fail

Frequent itemsets - possibilities are combinatorially large, but sparsity helps

Apriori Principle

If X is frequent, every $Y \subseteq X$ must also be frequent.

Corollary

If Y is not frequent, no $X \supseteq Y$

is frequent — Prune space for counting

Level by level calculation

F_1 - frequent sets of size 1 - count

C_2 - candidate frequent sets of size 2

↓ $F_1 \times F_1 \setminus Id \setminus \text{symmetric pairs}$

count $\rightarrow F_2$

C_3 - all $\{x, y, z\}$ s.t. $\{x, y\}, \{y, z\}, \{x, z\}$ in F_2

$$F_{k-1} \rightarrow C_k$$

Enumerate all k -subsets of I

Check if each $k-1$ -subset is in F_{k-1}

Combinatorially infeasible

Lexicographically order each subset

Assume I is ordered

$$i_1 < i_2 < \dots < i_m$$

Given $\{i_1, i_2, \dots, i_{k-2}, i_{k-1}\} \in F_{k-1}$
 $\{i_1, i_2, \dots, i_{k-2}, i'_{k-1}\}$



$\{i_1, i_2, \dots, i_{k-2}, i_{k-1}, i'_{k-1}\}$

Option 1 Use this more permissive

list of candidates $\underline{\underline{C'_k \supseteq C_k}}$

Check this

Suppose $\{j_1, j_2, \dots, j_k\}$ genuinely
in C_k

- Every $k-1$ subset is in F_{k-1}

$\{j_1, \dots, j_{k-2}, j_{k-1}\} \in F_{k-1}$
 $\{j_1, \dots, j_{k-2}, j_k\} \in F_{k-1}$

generates this candidate in C_k

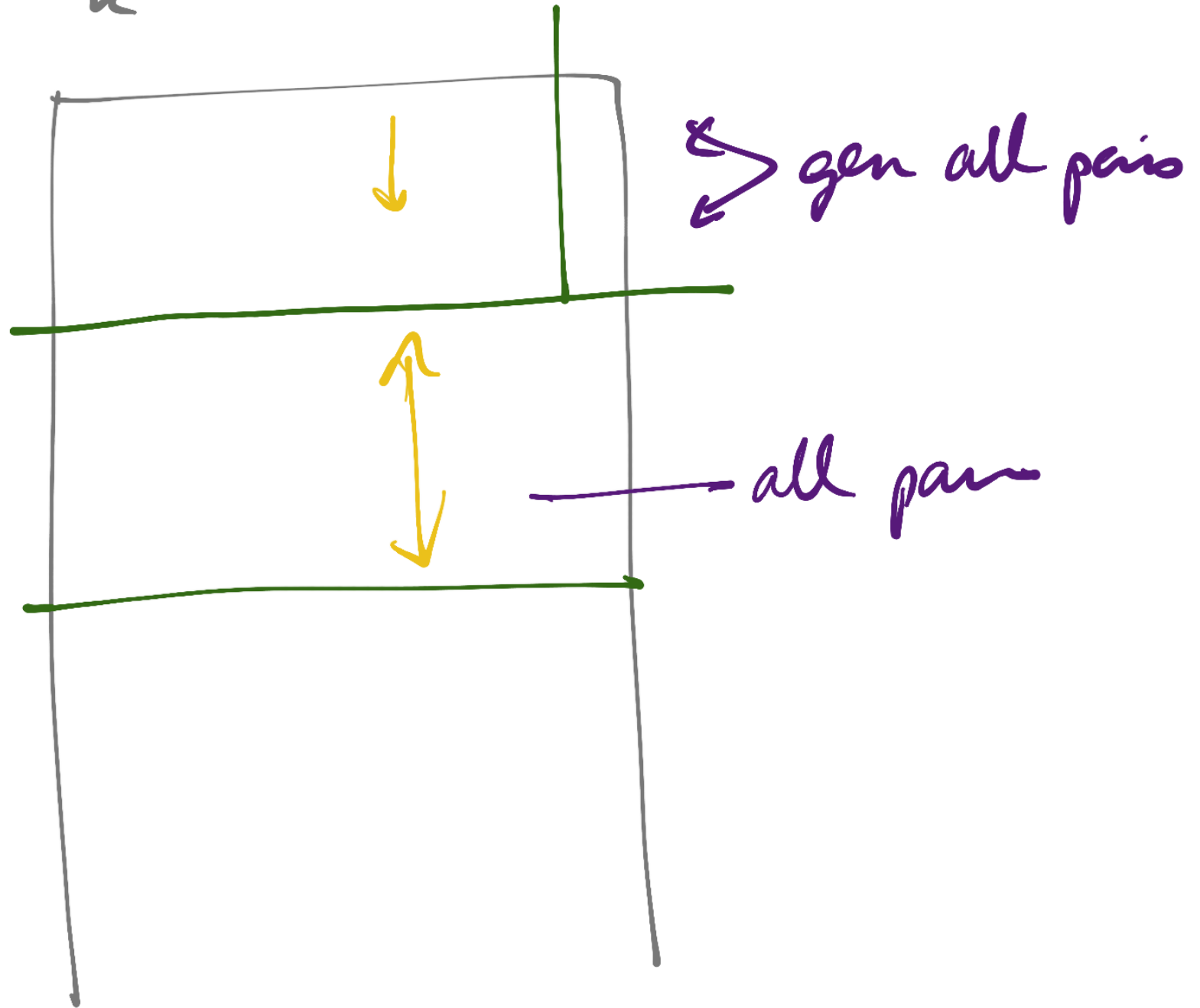
Option 2 Prune C_k' to C_k

$$\{f_1, f_2, \dots, f_k\} \in C_k'$$

Extra work, exact C_k

$$F_{k-1} \rightarrow C'_k$$

F_{k-1}
 in
 lex
 order



Termination

1. All t_i have size $\leq S$

level S is largest possible

2. Some F_{k-1} generates empty C_k

or F_k is empty

3. Fix some small k that is useful

This problem has a correct answer

Uses of frequent item sets

Association Rules

$$X \rightarrow Y, \quad X, Y \subseteq I \\ X \cap Y = \emptyset$$

"People who bought X also bought Y"

Two parameters to make a rule worth considering

- Support — how often is it seen?

$\text{Sup}(X \rightarrow Y)$ is just $\text{sup}(X \cup Y)$

min-Sup

- confidence = $\frac{\text{sup}(X \cup Y)}{\text{sup}(X)}$

min-conf

From frequent itemsets (A-priori)
to association rules

$F = F_1 \cup F_2 \cup \dots \cup F_s$ — candidates
for $X \cup Y$
for a rule
 $X \rightarrow Y$

Given $Z \in F$

does it break up as $X \cup Y$

s.t. $X \rightarrow Y$ has min-conf?

Now, consider every partition of
 Z as $A, Z \setminus A$

Given $A, B = Z \setminus A$

$$\text{conf}(A \rightarrow B) = \frac{\text{sup}(A \cup B = Z)}{\text{sup}(A)}$$

both are in F_i, F_j

for some i, j

Counts are available
already

A priori again!

$A \rightarrow \{a, b\}$ is a rule with
 $\text{conf} > \text{min-conf}$

What about $A \cup \{a\} \rightarrow \{b\}$?

$$\frac{\sup(A \cup \{a, b\})}{\sup(A)} \overset{=}{\longleftarrow} \frac{\sup((A \cup \{a\}) \cup \{b\})}{\sup(A \cup \{a\})}$$

As a fraction $\text{LHS} \leq \text{RHS}$

In general $A \rightarrow B$ is "good" $\Rightarrow A \cup \{a\} \rightarrow B \setminus \{a\}$ is also "good"

Conversely

If $A \rightarrow B$ is not good

$A \setminus \{a\} \rightarrow B \cup \{a\}$ is also not good

Decomposition of $X \subseteq F$ as A, B

level 1 $X \setminus \{a\} \rightarrow \{a\}$

level 2 $X \setminus \{a, b\} \rightarrow \{a, b\}$ if both
 $X \setminus \{a\} \rightarrow \{a\}$ &
 $X \setminus \{b\} \rightarrow \{b\}$ are
validated in level 1

Special Case

Basket as rows in a table

A column that is a category

| | Words | | | | Topic |
|-------|-------|-------|-----|-------|-------|
| | w_1 | w_2 | ... | w_N | |
| d_1 | | | | | Sport |
| d_2 | | | | | Art |
| d_3 | | | | | Art |

Rules of the form: $X \subseteq \text{Words} \rightarrow \text{Topic}$

Association rules \Rightarrow Topic classifier

Association rules can be used for
classification = "supervised learning"

One problem to be addressed -
items have different "natural" frequency

Split I into disjoint buckets

- Will only find relationships within
a category

"Manually" assign $\text{min-sup}(i)$ separately
for each $i \in I$

$X \subseteq I$ $\{i_1, i_2, \dots, i_k\}$

When is X frequent? Each i_j has
different $\text{min-sup}()$

Natural choice: $\text{min-sup}(X) = \min_j \text{min-sup}(i_j)$

Lose a-priori property!

$\{i_1, i_2, i_3\} \rightarrow \text{min-sup}$
0.3 0.1 0.4 is 0.1

$\{i_1, i_3\}$ occurs with sup 0.25
not freq - min-sup is 0.3

Problem only occurs when smallest min-sup
item is dropped

Modified level wise calculation

Earlier - each $X \subseteq I$ ordered acc to I

Instead order by min-sup

$$\{i_1, i_2, \dots, i_k\}$$

$$\text{min-sup}(i_1) \leq \text{min-sup}(i_2) \leq \dots \leq \text{min-sup}(i_k)$$

F_1 ✓

$$C_2 = F_1 \times F_1?$$

$$(i, j) \in C_2$$



should meet $\text{minsup}(i)$

$$C_2 = \left\{ (j, k) \mid \begin{array}{l} j \in F_1 \\ k > j, \quad k\text{-count} \\ \geq \min\text{-sup}(j) \end{array} \right\}$$

Finish next time