

# Data Mining & Machine Learning (DMML)

↓  
Extract insights  
from data

↓  
Extrapolate from past data  
- Prediction

Data collection → Cleaning

Not focus on this

Patterns in data - Statistics

Analytical - build mathematical relationship

"Equations"

Use computational power to derive relations  
from volumes of data

Two types of problems



Patterns

Unsupervised



Clustering in  
groups



Predictions

Supervised



Use a "model"

Find values of parameters  
from historical data

# Market-Basket Analysis

Shopping baskets

Identify frequently occurring combinations

Items  $I = \{i_1, i_2, \dots, i_n\}$

Basket = Transaction  $t \in I$

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

Threshold : - how frequent

is a subset  $X \subseteq I$  ?

$$\frac{\# t_i \text{ containing } x}{M} = \text{Support}(x)$$

Frequent:  $\text{Support}(x) \geq \text{Threshold}$

↑  
input  
parameter

Given threshold, find all  $x$

s.t.  $\text{Support}(x) \geq \text{threshold}$ .

Given a list of integers

$y_1, y_2, \dots, y_M$

find all integers that appear with

support  $\geq 0.01$

Dictionary

Keys - integers (that appear)

Values - frequency

Additional constraint: All values  $\leq 10^6$

Index  $\equiv$  Input Integer

Array Count  $[0 \dots 10^6]$

Initialize to 0

See  $i \rightarrow$  increment  $\text{count}[i]$

Items:  $I = \{i_1, \dots, i_n\}$

Fix  $x \in I$  Go through  $T = \{t_1, \dots, t_n\}$   
& count frequency of  $x$

To do this for every  $X \subseteq I$   $|I|=N$

Potentially -  $2^N$  subsets

Space - Counters take too much space

Time - Each transaction requires  
 $2^N$  updates to be checked

Calculation Assume each  $t_i \leq 10$  items

$$\begin{aligned} |I| = N &= 10^6 \\ |T| = M &= 10^9 \end{aligned}$$

Threshold is 0.01

Simpler question

How many singleton sets  $\{i_k\}$   
are frequent?

$$|I| = 10^6$$

$$|T| = 10^9$$

$$\text{Thresh} = 0.01$$

$$|t_i| \leq 10$$

How many items in  $T$ ?

See at most  $10^{10}$  items  
across all  $\mathcal{D}_b T$

Frequent item must  
appear  $10^7$  times

Frequent items at most  $\frac{10^{10}}{10^7}$   
 $= 1000$



Suppose  $\{x, y\}$  appears frequently

What can we conclude about  $\{x\}, \{y\}$ ?

If  $\{x\}$  or  $\{y\}$  is not frequent,

$\{x, y\}$  cannot be frequent

Example  $\rightarrow$  1000 potentially frequent  $\{x\}$

$\{x, y, z\}$  frequent  $\rightarrow \{x, y\}, \{y, z\},$   
 $\{x, z\}$  all freq.

## A Priori Principle

For  $X$  to be frequent, every subset of  $X$  must be frequent

## (Layered) A Priori Algorithm

Count frequent sets of size 1

↳ List of candidates of size 2

↳ Count frequent sets of size 2

↳ Candidates of size 3 → count size 3

When do we stop?

↳ Candidate size exceeds max  $k_i$  size

↳ Or the frequent item set count goes to zero at some level

↳ Use case - stop at some small level

Given  $|t_i| \leq K$

Make upto  $K$  passes over  $T$

level 1

level 2

⋮

level  $K$

$\simeq 10^9$  ops/sec

Computational bottleneck

$F_i$  - frequent sets at level  $i$   
to  $C_{i+1}$  - Candidates at level  $i+1$

$$F_i \rightarrow C_{i+1} = \{X \text{ of size } i+1, \\ \text{every } i\text{-subset of} \\ X \text{ is in } F_i\}$$

Instead of  $C_{i+1}$ , can use any  $D_{i+1}$   
that is a superset of  $C_{i+1}$

Strategy Order the items  $i_1 < i_2 < \dots < i_n$

Enumerate any  $t \in I$  in ascending  
ordering

take two transactions in  $F_i$

$$\begin{aligned} t_1 &= \{j_1, j_2, \dots, j_{i-1}, j_i\} \\ t_2 &= \{j_1, j_2, \dots, j_{i-1}, j'_i\} \end{aligned}$$

$\rightarrow t_{12} = \{j_1, j_2, \dots, j_{i-1}, j_i, j'_i\}$

$F_i$

Dictionary  
Order

