

Archit Jain

AT

Crib sheet

Con Association ($X \rightarrow Y$) (2 step)

Apriori algorithm

$\{Diaper \rightarrow Beer\}$ \rightarrow (Co-occurrence)
(Inventory management) (not causality)

(Frequent Itemset)

Before \rightarrow rules

* Support - fraction of transaction
contains X & Y

1) Itemset - collection of items $\{Milk, Bread\}$

2) Support Count \rightarrow Count $\rightarrow \sigma\{Milk, Bread\} = 2$

3) Support \rightarrow fraction of $\frac{\sigma\{Milk, Bread\}}{Total} = \frac{2}{5}$

4) Frequent-Itemset \rightarrow

* Confidence Measure, how often items in X Support appear in transaction X
Support is greater or min Support
Support \geq minsup threshold

* Prune \rightarrow delete rules. (Computational \rightarrow Brute force)

Two step approach: 1) frequent itemset generation (Support \geq minsup)

(Computationally expensive) 2) rule generation (Confidence \geq min Conf)

$\sim O(NMw)$ (width)
Transaction \rightarrow frequent itemset

$$\left[\text{Conf} = \frac{\sigma(Y)}{\sigma(X)} \right]$$

Minsupport issue \rightarrow too high \rightarrow miss itemset with rare items
too low \rightarrow computationally expensive.

Drawback Confidence.

Hides negative correlation

Statistical Independence $P(S \cap B) = P(S) \times P(B)$ = independent

$P(S \cap B) > P(S) \times P(B)$ = Positive correlation

$P(S \cap B) < P(S) \times P(B)$ = Negative correlation

$$\text{Lift} = \frac{\text{Conf}(X \rightarrow Y)}{\text{Sup}(Y)} = \frac{P(X, Y)}{P(X) \cdot P(Y)} = \begin{matrix} 1 & \text{independent} \\ > 1 & \text{positive} \\ < 1 & \text{negative} \end{matrix}$$

Interest = when positive

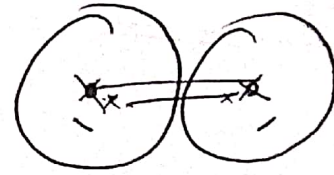
	Coffee	
Tea	15	20
Tea	90	
		$P(\text{Coffee}) = 0.75$
		$\frac{P(\text{Coffee})}{P(\text{Tea})} = \frac{15}{20}$

(K-means)

(2)

fix number of iterations

Stopping criteria - iteration
Same Centroid.



Initial Centroid random,

6 points \rightarrow

Euclidean distance:

Manhattan \rightarrow Summation of absolute value

Cosine Simi :

Convergence

[K means \rightarrow Important to initial Centroids.

doesn't give global, just local solution.

Min (SSE = Sum of Squared error)



$$SSE = \sum_{i=1}^K \sum_{x \in C_i} \text{dist}^2(m_i, x)$$

minimize \rightarrow distance from Center

choose less SSE

min(SSE) increase $\rightarrow (K) \uparrow$

good cluster with smaller K have smaller (SSE)

