

Agenda

- Business Problem
- Apriori Algorithm
- Association Rule
- Metrics for Association Rule:
 - Support
 - Confidence
 - Lift
 - Leverage
 - Conviction

} important

Business Problem

Task
⇒ which products are frequently bought Together

- ① product placement
- ② promotion
- ③ Bundles

id	product

Reliance Store

Products : Milk, bread, butter ---
--- 1000's

$D \in (1, 2, 3, \dots, n)$

2^n & 2^{100}

$C_1 \rightarrow T_1 \rightarrow \{1, 3, 8, 10\}$

$C_2 \rightarrow T_2 \rightarrow \{10, 15, 1, 3\}$

$C_3 \rightarrow T_3 \rightarrow \{1, 1000, 10, 3\}$

..
..
..
..

..

..

Pattern \rightarrow

$\rightarrow \{1, 3, 10\}$

$\rightarrow \{1, 3\}$

$\rightarrow \{3, 10\}$

$\rightarrow \{1, 10\}$

$1 \xrightarrow{\text{or command}} 3, 10$
 $1, 3 \rightarrow 10$

Data Set

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	536365	85123A	WHITE HANGING HEART T-LIGHT HOLDER	6	01/12/10 8:26	2.55	17850.0	United Kingdom
1	536365	71053	WHITE METAL LANTERN	6	01/12/10 8:26	3.39	17850.0	United Kingdom
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8	01/12/10 8:26	2.75	17850.0	United Kingdom
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	6	01/12/10 8:26	3.39	17850.0	United Kingdom
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	6	01/12/10 8:26	3.39	17850.0	United Kingdom



Pre-processing

Invoice No	item1	item2	item3	...	itemN
101	1	0	1	...	1
102	0	0	1	...	1

Sparse Matrix

Task : Find Frequently occurring items

D : Set of all items

T : Set of all Transactions

$T_i \subseteq D$

Apriori Algorithm

	Key	Value	
itemsets	$\{1, 2\}$	10	frequency
	$\{1, 3, 10\}$	20	
	$\{i, j, k, \dots\}$	Count	

$D \ni \{1, 2, 3\}$

itemsets $\{1\}, \{2\}, \{3\}, \{1, 2\}$
 $\{1, 3\}, \{2, 3\}, \{1, 2, 3\}$

$D \ni \{1, 2, 3, \dots, n\}$ 7

itemsets? $2^n - 1 - n$

$100 \ni 2^{100} - 100 - 1$

Brute-force

$\{1, 3\}$

```

i-sets = get-all-sets(D)
for item-set in i-sets:
    for t in transactions:
        if item-set in t:
            Q[item-set] += 1
    
```

T.C. $\mathcal{O}(Q^n \times \text{transaction})$
 $Q^{4000} \times 15,000$

Very Very Slow

Thresholding

$$T = 100$$

Ex:

$$\{2\} \rightarrow 99$$

$$\{2, 3\} \leq 99$$

minimum
Support

itemset $<$ threshold

remove all supersets of
this item-set

Ex 2:

$$\{4, 3\} < T$$

$$\# \{2, 4, 3\} \times \quad \{2, 4, 3, 10, 11\} \times$$

Brute force $\Rightarrow O(2^n \times m)$



reduce this drastically by removing all sets which doesn't satisfy minimum support

* FP Growth (frequent itemset)



TRIES

Market Basket Analysis

Key	Value
$\{1, 2, 3\}$	10
$\{1, 3, 10\}$	20
$\{i, j, k, \dots\}$	Count

(5-6) (8) 2^{44}

2^{19}

Applications of Market Basket Analysis

Bioinformatics

two chemicals $\begin{cases} \rightarrow c_1 \\ \rightarrow c_2 \end{cases}$

frequently occurring position
across chemicals

gene sequences $\begin{cases} \rightarrow m_1 \\ \rightarrow m_2 \end{cases}$

ATTC and AGTC



frequently observed

Medicine $\begin{cases} \rightarrow m_1 \\ \rightarrow m_2 \\ \rightarrow m_3 \end{cases}$ (relationship)

Webpage mining $\begin{cases} \rightarrow w_1 \\ \rightarrow w_2 \\ \rightarrow w_3 \end{cases}$

$\{w_1, w_2, w_3, \dots, w_k\} \Rightarrow$ frequent

Drawback \rightarrow Dataset \rightarrow Millions
of Products

Market Basket Analysis with
huge number products becomes
very costly and time-intensive

Association Rule

Quantify the relationship
discovered in market
Basket Analysis output

$X \Rightarrow \{1, 2, 3\}$

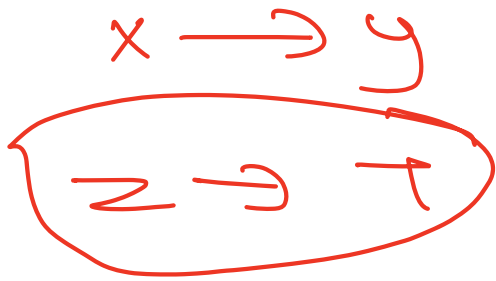
$Y \Rightarrow \{4, 5\}$

$Z \Rightarrow X \cup Y$

$Z \Rightarrow \{1, 2, 3, 4, 5\} \rightarrow$ frequent
set

$X \rightarrow Y$

People who buy X have high likelihood
of buying Y



q_1

q_2

$q_2 > q_1$

Domino's Pizza

Pizza \rightarrow Coke \Rightarrow rule 1

Pizza \rightarrow Garlic Bread \Rightarrow rule 2

rule 1 $>$ rule 2

Coke \rightarrow Pizza = Pizza \rightarrow Coke
 No X

$x \longrightarrow y$
 (Antecedent) \longrightarrow Consequent (then)
 (if)

Bread \longrightarrow Butter 63%
 Ante Conseq

Association Metrics

Support	
Confidence	Conviction
Lift	Leverage

Support

$$\text{Support}(X) \equiv \frac{\# \text{Transactions}(X)}{\# \text{Total Transaction}}$$

$$\# D(X \text{ in } T)$$

milk = 100

Total \approx 1000

$$\text{Support}(\text{milk}) \approx \frac{100}{1000} \approx 0.1\%$$

Confidence

$$\text{Conf}(X \rightarrow Y) \equiv \frac{\#(X \text{ and } Y)}{\#(X)}$$

$$1) \Rightarrow P(Y/X)$$

$$P \rightarrow C \Leftrightarrow \frac{P \text{ and } C}{P}$$

$$C \rightarrow P \Leftrightarrow \frac{C \text{ and } P}{C}$$

$$\text{Conf}(P \rightarrow C) \neq \text{Conf}(C \rightarrow P)$$

Interpretation

of all the times X occurs
how many times Y also occurs

① $\{i+1, i+2\} \Rightarrow$ High frequency
High Support

$i+1 \rightarrow i+2 \Rightarrow$ low confidence

② $\{i+1, i+2\} \Rightarrow$ Low frequency
Low support

$i+1 \rightarrow i+2 \Rightarrow P_{i+1, i+2} \Rightarrow$ High Confidence

Light

Yogurt \rightarrow milk (High Confidence)

Cornflakes \rightarrow milk (High Confidence)

Toothpaste \rightarrow milk ?

Not Sure

* Milk itself is a very frequent item

Transaction $\rightarrow 100$

Milk $\rightarrow 80$

Toothpaste $\rightarrow 14$

* Let's say: 10 of them contain both TP and milk

$$\text{Confidence } (TP \rightarrow \text{milk}) = \frac{10}{14} \approx 0.7$$

$$\text{Lift } (x \rightarrow y)$$

$$\frac{\text{Support}(x \cap y)}{\text{Support}(x) \times \text{Support}(y)}$$

$$\approx \frac{P(x \text{ and } y)}{P(x) \times P(y)} = 1$$

Interpretation

Increase in probability of x and y together compared to when they are independent.

$\text{Lift}(x \rightarrow y) = 1$, x and y are independent

$\text{Lift}(x \rightarrow y) > 1$, likely to be bought together

$\text{Lift}(x \rightarrow y) < 1$, Negatively related

$$\text{Lift}(\text{TP} \rightarrow \text{milk}) = \frac{14/100}{\frac{80}{100} \times \frac{10}{100}} = 0.8$$

$$\text{Lift} \Rightarrow \frac{P(X \text{ and } Y)}{P(X) \times P(Y)}$$

$\rightarrow \text{conf}(x \rightarrow y)$

$$\text{Lift}(X \rightarrow Y) = \frac{\text{Conf}(X \rightarrow Y)}{\text{Support}(Y)}$$

$$\Rightarrow \frac{0.7}{0.8} = 0.875 < 1$$

Leverage

- ① Same as Lift
- ② Easier to interpret

Range:

$$\text{Lift} \rightarrow (0, \infty)$$

$$\text{Leverage} \rightarrow (-1, 1)$$

$$\text{Lev}(x \rightarrow y) = \frac{\text{Support}(x \cap y)}{\text{Support}(x) * \text{Support}(y)}$$

$$\text{Support}(x) * \text{Support}(y)$$

⑨ Implement of the Metric in Python

⑨ Conviction

Conviction

Ratio of Expected frequency that x occurs without y if x and y were independent / Observed frequency of incorrect predictions

$$\text{Conv}(x \rightarrow y) = \frac{1 - s(y)}{1 - c(x \rightarrow y)}$$

Note: A High value means Consequent depends Strongly on the antecedent

Q: Can you calculate Lift of Milk and Bread from given data:

- Total Transactions = 600,000
- Transaction {Bread} = 7500
- Transaction {Milk} = 60000
- Transactions {Milk, Bread} = 6000

In Break

Lift (M \rightarrow B)

$$\frac{\text{Conf}(M \rightarrow B)}{\text{Support}(B)}$$

Lift (B \rightarrow M)

$$\frac{\text{Conf}(B \rightarrow M)}{\text{Support}(M)}$$

Cut ≥ 100

only 10 beg!

Milk + Bread > 10

$\Sigma 23 \Rightarrow$ all customers who
bought