

Class 11

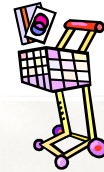
Market Basket Analysis



Learning Objectives

- Market Basket Analysis
 - Introduction and usage
- Association rules
 - Definition
 - Apriori
 - Usage issues

Market Basket Analysis



- Retail – each customer purchases different set of products, different quantities, different times
- MBA uses this information to:
 - Identify who customers are (not by name)
 - Understand why they make certain purchases
 - Gain insight about its merchandise (products):
 - Fast and slow movers
 - **Products which are purchased together**
 - Products which might benefit from promotion
 - Take action:
 - Store layouts
 - Which products to put on specials, promote, coupons...
- Combining all of this with a customer loyalty card it becomes even more valuable

3

Nappies and beer



<http://www.daedalus.es/en/data-mining/nappies-and-beer/>

4

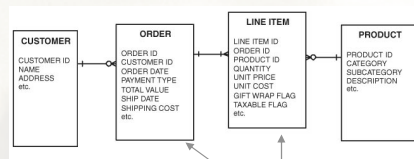
Market Basket Analysis Drill-Down

- MBA is a set of techniques, Association Rules being most common, that focus on point-of-sale (p-o-s) transaction data
- 3 types of market basket data (p-o-s data)
 - Customers
 - Orders (basic purchase data)
 - Items (merchandise/services purchased)

5

Typical Data Structure (Relational Database)

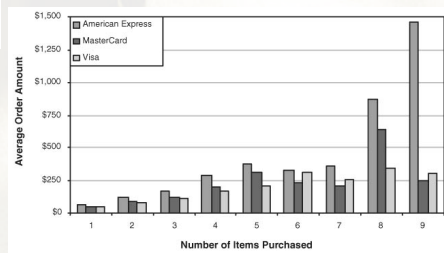
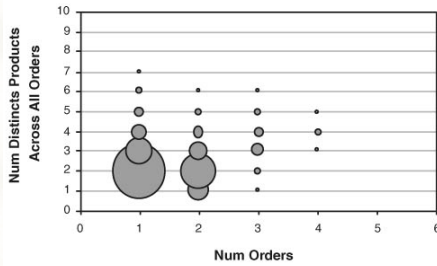
- Lots of questions can be answered
 - Avg # of orders/customer
 - Avg # unique items/order
 - Avg # of items/order
 - For a product
 - What % of customers have purchased
 - Avg # orders/customer include it
 - Avg quantity of it purchased/order
 - Etc...
- Visualization is extremely helpful...next slide



Transaction Data

6

Sales Order Characteristics



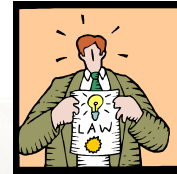
7

Sales Order Characteristics

- Did the order use gift wrap?
- Billing address same as Shipping address?
- Did purchaser accept/decline a cross-sell?
- What is the most common item found on a one-item order?
- What is the most common item found on a multi-item order?
- What is the most common item for repeat customer purchases?
- How has ordering of an item changed over time?
- How does the ordering of an item vary geographically?

8

Association Rules



- Association rule types:
 - Actionable Rules – contain high-quality, actionable information
 - Wal-Mart customers who purchase Barbie dolls have a 60% likelihood of also purchasing one of three types of candy bars [*Forbes*, Sept 8, 1997]
 - Trivial Rules – information already well-known by those familiar with the business
 - Customers who purchase maintenance agreements are very likely to purchase large appliances
 - Inexplicable Rules – no explanation and do not suggest action
 - When a new hardware store opens, one of the most commonly sold items is toilet bowl cleaners
- Trivial and Inexplicable Rules occur most often

9

How Good is an Association Rule?

Customer	Items Purchased
1	OJ, soda
2	Milk, OJ, window cleaner
3	OJ, detergent
4	OJ, detergent, soda
5	Window cleaner, soda

← POS Transactions

↙ Co-occurrence of Products

	OJ	Window cleaner	Milk	Soda	Detergent
OJ	4	1	1	2	2
Window cleaner	1	2	1	1	0
Milk	1	1	1	0	0
Soda	2	1	0	3	1
Detergent	2	0	0	1	2

10

How Good is an Association Rule?

	OJ	Window cleaner	Milk	Soda	Detergent
OJ	4	1	1	2	2
Window cleaner	1	2	1	1	0
Milk	1	1	1	0	0
Soda	2	1	0	3	1
Detergent	2	0	0	1	2

Simple patterns:

1. OJ and soda are more likely purchased together than any other two items
2. Detergent is never purchased with milk or window cleaner
3. Milk is never purchased with soda or detergent

11

How Good is an Association Rule?

Customer	Items Purchased
1	OJ, soda
2	Milk, OJ, window cleaner
3	OJ, detergent
4	OJ, detergent, soda
5	Window cleaner, soda

← POS Transactions

- What is the confidence for this rule:
 - If a customer purchases soda, then customer also purchases OJ
 - 2 out of 3 soda purchases also include OJ, so 67%
- What about the confidence of this rule reversed?
 - 2 out of 4 OJ purchases also include soda, so 50%
- **Confidence** = Ratio of the number of transactions with all the items to the number of transactions with just the “if” items

12

Multi-item association rules

Co-occurrence can occur in 3, 4, or more dimensions...

1. Generate co-occurrence matrix for single items... "if OJ then soda"
2. Generate co-occurrence matrix for two items... "if OJ and Milk then soda"
3. Generate co-occurrence matrix for three items... "if OJ and Milk and Window Cleaner" then soda
4. Etc...

13

Association Rule Mining

1. Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction

Market-Basket transactions

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example of Association Rules

{Diaper} → {Beer},
{Milk, Bread} → {Eggs, Coke},
{Beer, Bread} → {Milk},

Implication means co-occurrence,
not causality!

14

Definition: Frequent Itemset

1 Itemset

- A collection of one or more items
 - Example: {Milk, Bread, Diaper}
- k-itemset
 - An itemset that contains k items

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

1 Support count (σ)

- Frequency of occurrence of an itemset
- E.g. $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$

1 Support

- Fraction of transactions that contain an itemset
- E.g. $s(\{\text{Milk, Bread, Diaper}\}) = 2/5$

1 Frequent Itemset

- An itemset whose support is greater than or equal to a *minsup* threshold

15

Definition: Association Rule

1 Association Rule

- An implication expression of the form $X \rightarrow Y$, where X and Y are itemsets
- Example: $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

1 Rule Evaluation Metrics

- Support (s)
 - Fraction of transactions that contain both X and Y
- Confidence (c)
 - Measures how often items in Y appear in transactions that contain X

Example:

$\{\text{Milk, Diaper}\} \Rightarrow \text{Beer}$

$$s = \frac{\sigma(\text{Milk, Diaper, Beer})}{|T|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

16

Association Rule Mining Task

- 1 Given a set of transactions T, the goal of association rule mining is to find all rules having
 - support \geq *minsup* threshold
 - confidence \geq *minconf* threshold
- 1 Brute-force approach:
 - List all possible association rules
 - Compute the support and confidence for each rule
 - Prune rules that fail the *minsup* and *minconf* thresholds

⇒ **Computationally prohibitive!**

17

Mining Association Rules

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example of Rules:

{Milk,Diaper} \rightarrow {Beer} (s=0.4, c=0.67)
{Milk,Beer} \rightarrow {Diaper} (s=0.4, c=1.0)
{Diaper,Beer} \rightarrow {Milk} (s=0.4, c=0.67)
{Beer} \rightarrow {Milk,Diaper} (s=0.4, c=0.67)
{Diaper} \rightarrow {Milk,Beer} (s=0.4, c=0.5)
{Milk} \rightarrow {Diaper,Beer} (s=0.4, c=0.5)

Observations:

- All the above rules are binary partitions of the same itemset:
{Milk, Diaper, Beer}
- Rules originating from the same itemset have identical support but can have different confidence
- Thus, we may decouple the support and confidence requirements

18

Mining Association Rules

1 Two-step approach:

1. Frequent Itemset Generation

- Generate all itemsets whose support \geq minsup

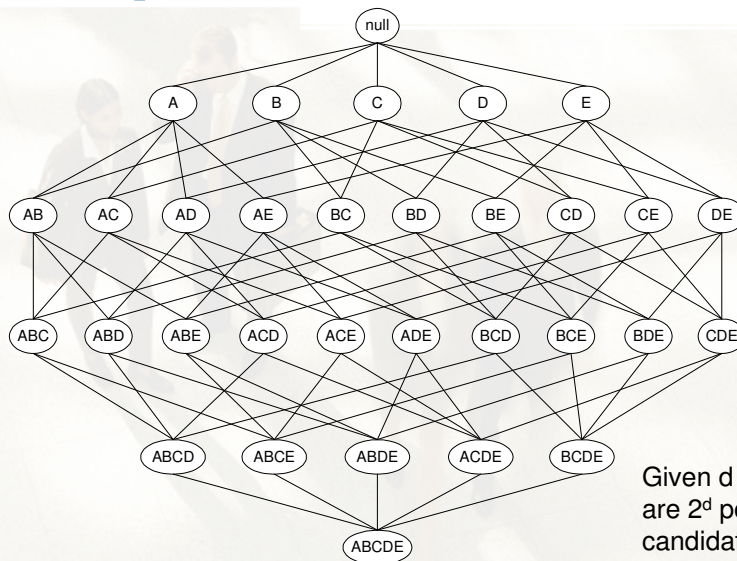
2. Rule Generation

- Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset

1 Frequent itemset generation is still computationally expensive

19

Frequent Itemset Generation



20

Reducing Number of Candidates

1 Apriori principle:

- If an itemset is frequent, then all of its subsets must also be frequent

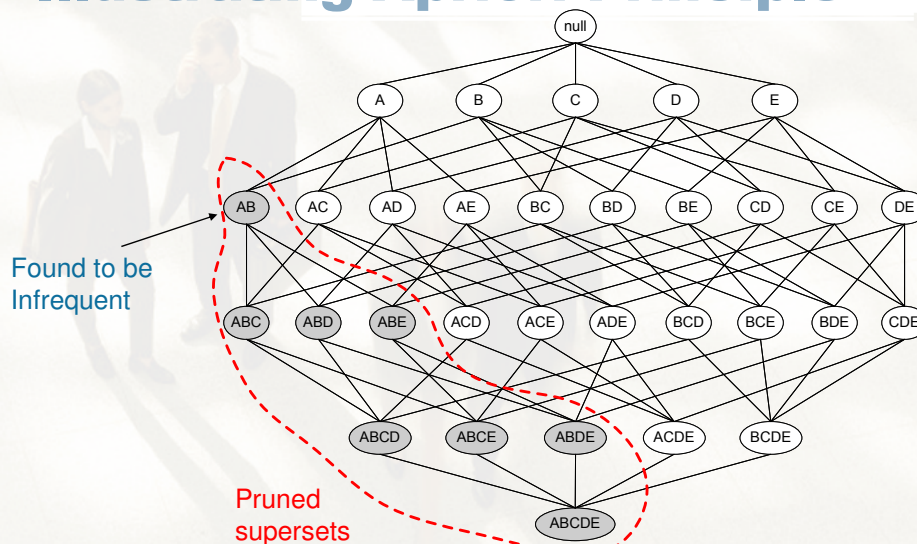
1 Apriori principle holds due to the following property of the support measure:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

- Support of an itemset never exceeds the support of its subsets
- This is known as the **anti-monotone** property of support

21

Illustrating Apriori Principle



22

Apriori Algorithm

1 Method:

- Let $k=1$
- Generate frequent itemsets of length 1
- Repeat until no new frequent itemsets are identified
 - u Generate length $(k+1)$ candidate itemsets from length k frequent itemsets
 - u Prune candidate itemsets containing subsets of length k that are infrequent
 - u Count the support of each candidate by scanning the DB
 - u Eliminate candidates that are infrequent, leaving only those that are frequent

23

Factors Affecting Complexity

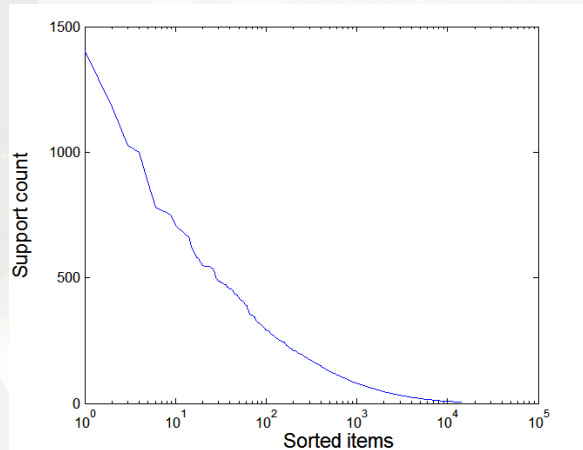
- 1 Choice of minimum support threshold
 - lowering support threshold results in more frequent itemsets
 - this may increase number of candidates and max length of frequent itemsets
- 1 Dimensionality (number of items) of the data set
 - more space is needed to store support count of each item
 - if number of frequent items also increases, both computation and I/O costs may also increase
- 1 Size of database
 - since Apriori makes multiple passes, run time of algorithm may increase with number of transactions
- 1 Average transaction width
 - transaction width increases with denser data sets
 - This may increase max length of frequent itemsets and traversals of hash tree (number of subsets in a transaction increases with its width)

24

Effect of Support Distribution

- 1 Many real data sets have skewed support distribution

Support distribution of a retail data set



25

Effect of Support Distribution

- 1 How to set the appropriate *minsup* threshold?
 - If *minsup* is set too high, we could miss itemsets involving interesting rare items (e.g., expensive products)
 - If *minsup* is set too low, it is computationally expensive and the number of itemsets is very large
- 1 Using a single minimum support threshold may not be effective

26

Multiple Minimum Support

1 How to apply multiple minimum supports?

- $MS(i)$: minimum support for item i
- e.g.: $MS(\text{Milk})=5\%$, $MS(\text{Coke}) = 3\%$,
 $MS(\text{Broccoli})=0.1\%$, $MS(\text{Salmon})=0.5\%$
- $MS(\{\text{Milk}, \text{Broccoli}\}) = \min (MS(\text{Milk}), MS(\text{Broccoli}))$
 $= 0.1\%$
- Challenge: Support is no longer anti-monotone
 - Suppose: $\text{Support}(\text{Milk}, \text{Coke}) = 1.5\%$ and
 $\text{Support}(\text{Milk}, \text{Coke}, \text{Broccoli}) = 0.5\%$
 - $\{\text{Milk}, \text{Coke}\}$ is infrequent but $\{\text{Milk}, \text{Coke}, \text{Broccoli}\}$ is frequent

27

Pattern Evaluation

1 Association rule algorithms tend to produce too many rules

- many of them are uninteresting or redundant
- Redundant if $\{A, B, C\} \rightarrow \{D\}$ and $\{A, B\} \rightarrow \{D\}$
have same support & confidence

1 Interestingness measures can be used to prune/rank the derived patterns

1 In the original formulation of association rules, support & confidence are the only measures used

28

Computing Interestingness Measure

- Given a rule $X \rightarrow Y$, information needed to compute rule interestingness can be obtained from a contingency table

Contingency table for $X \rightarrow Y$

	Y	\bar{Y}	
X	f_{11}	f_{10}	f_{1+}
\bar{X}	f_{01}	f_{00}	f_{0+}
	f_{+1}	f_{+0}	T

f_{11} : support of X and Y
 f_{10} : support of X and \bar{Y}
 f_{01} : support of \bar{X} and Y
 f_{00} : support of \bar{X} and \bar{Y}

Used to define various measures

- support, confidence, lift, Gini, J-measure, etc.

29

Drawback of Confidence

	Coffee	$\bar{\text{Coffee}}$	
Tea	15	5	20
$\bar{\text{Tea}}$	75	5	80
	90	10	100

Association Rule: Tea \rightarrow Coffee

Confidence = $P(\text{Coffee}|\text{Tea}) = 0.75$

but $P(\text{Coffee}) = 0.9$

\Rightarrow Although confidence is high, rule is misleading

$\Rightarrow P(\text{Coffee}|\bar{\text{Tea}}) = 0.9375$

30

Statistical Independence

1 Population of 1000 students

- 600 students know how to swim (S)
- 700 students know how to bike (B)
- 420 students know how to swim and bike (S,B)

- $P(S \wedge B) = 420/1000 = 0.42$
- $P(S) \times P(B) = 0.6 \times 0.7 = 0.42$

- $P(S \wedge B) = P(S) \times P(B) \Rightarrow$ Statistical independence
- $P(S \wedge B) > P(S) \times P(B) \Rightarrow$ Positively correlated
- $P(S \wedge B) < P(S) \times P(B) \Rightarrow$ Negatively correlated

31

Statistical-based Measures

1 Measures that take into account statistical dependence

$$\text{Lift} = \frac{P(Y | X)}{P(Y)}$$

$$\text{Interest} = \frac{P(X, Y)}{P(X)P(Y)}$$

$$PS = P(X, Y) - P(X)P(Y)$$

$$\phi\text{-coefficient} = \frac{P(X, Y) - P(X)P(Y)}{\sqrt{P(X)[1 - P(X)]P(Y)[1 - P(Y)]}}$$

32

Example: Lift/Interest

	Coffee	<u>Coffee</u>	
Tea	15	5	20
<u>Tea</u>	75	5	80
	90	10	100

Association Rule: Tea → Coffee

Confidence= P(Coffee|Tea) = 0.75

but P(Coffee) = 0.9

⇒ Lift = 0.75/0.9= 0.8333 (< 1, therefore is negatively associated)

33

Drawback of Lift & Interest

	Y	<u>Y</u>	
X	10	0	10
<u>X</u>	0	90	90
	10	90	100

$$Lift = \frac{0.1}{(0.1)(0.1)} = 10$$

	Y	<u>Y</u>	
X	90	0	90
<u>X</u>	0	10	10
	90	10	100

$$Lift = \frac{0.9}{(0.9)(0.9)} = 1.11$$

Statistical independence:

If P(X,Y)=P(X)P(Y) => Lift = 1

34