

Introduction to Association Rule Mining

Introduction to Association Rules

- Substantial part of every E-Commerce and Supermarket

Frequently Bought Together



Total price: \$92.20

Add all three to Cart

Add all three to List



☒ This item: Hadoop: The Definitive Guide by Tom White Paperback \$33.43

☒ Learning Spark: Lightning-Fast Big Data Analysis by Holden Karau Paperback \$29.01

☒ Advanced Analytics with Spark: Patterns for Learning from Data at Scale by Sandy Ryza Paperback \$29.76

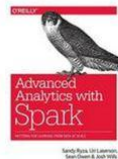
Customers Who Bought This Item Also Bought



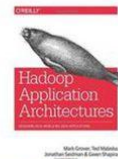
Learning Spark:
Lightning-Fast Big Data
Analysis



Big Data: Principles and
best practices of scalable
realtime data systems



Advanced Analytics with
Spark: Patterns for
Learning from Data at ...



Hadoop Application
Architectures
> Mark Grover



ZooKeeper: Distributed
Process Coordination
> Flavio Junqueira

Introduction to Association Rules

- Supermarket & Market Basket Analysis



Association Rule Mining: Formal Definition

- Commonly used for shopping behavior analysis
 - aims at detection of rules such as
"If a customer purchases baby food then he also buys diapers in 70% of the cases"

Baby Food



Diapers

























What is Association Rule Mining?

- Motivation: finding regularities in data
 - What products are often purchased together in a supermarket?
 - What are the subsequent purchases after buying a PC?
 - Is there a serious chance of Covid-19 patients to suffer a heart attack?
 - Do people who read news about Trump are also interested in NFL?

Market Basket Analysis























- Retail organizations e.g. Supermarkets collect and store massive amounts of sales data called Basket Data
- A record consist of
 - Transaction date/ID
 - Items bought

Transaction 1	   
Transaction 2	  
Transaction 3	 
Transaction 4	 
Transaction 5	   
Transaction 6	  
Transaction 7	 
Transaction 8	 

Basic Terminology & Brute Force Method for mining association rules

Market Basket Analysis

- 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























Transaction 1	   
Transaction 2	  
Transaction 3	 
Transaction 4	 
Transaction 5	   
Transaction 6	  
Transaction 7	 
Transaction 8	 

Example of Association Rules

{Baby Food} \square {Diapers}

{Milk, Rice} \square {Beer}























Support

Transaction 1	   
Transaction 2	  
Transaction 3	 
Transaction 4	 
Transaction 5	   
Transaction 6	  
Transaction 7	 
Transaction 8	 

$$\text{Support} \{ \text{apple} \} = \frac{4}{8}$$

$$\text{Support} \{ \text{apple}, \text{beer mug} \} = \frac{3}{8}$$























Confidence

Transaction 1	   
Transaction 2	  
Transaction 3	 
Transaction 4	 
Transaction 5	   
Transaction 6	  
Transaction 7	 
Transaction 8	 

$$\text{Confidence} \{ \text{apple} \rightarrow \text{beer mug} \} = \frac{\text{Support} \{ \text{apple}, \text{beer mug} \}}{\text{Support} \{ \text{apple} \}} = \frac{3}{4}$$

$$\text{Confidence} \{ \text{beer mug} \rightarrow \text{apple} \} = \frac{\text{Support} \{ \text{apple}, \text{beer mug} \}}{\text{Support} \{ \text{beer mug} \}} = \frac{3}{6}$$

Lift

Transaction 1	   
Transaction 2	  
Transaction 3	 
Transaction 4	 
Transaction 5	   
Transaction 6	  
Transaction 7	 
Transaction 8	 

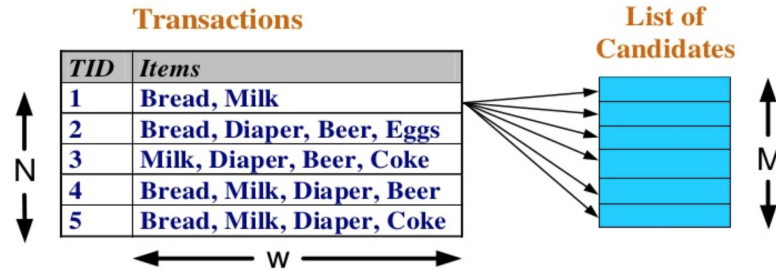
$$\text{Lift} \{ \text{apple} \rightarrow \text{beer} \} = \frac{\text{Support} \{ \text{apple}, \text{beer} \}}{\text{Support} \{ \text{apple} \} \times \text{Support} \{ \text{beer} \}} = \frac{3}{1 \times 2} = 1$$

Terminologies (1)

- **Itemset**
 - A collection of one or more items (Example: {Milk, Bread, Diaper})
- **k-itemset**
 - a set of k items.
 - E.g. {beer, cheese, eggs} is a 3-itemset
 - {cheese} is a 1-itemset
 - {honey, ice-cream} is a 2-itemset
- **Frequent/Large Itemset (L_k)**
 - An itemset whose support is greater than or equal to a minsup threshold
- **Candidate Itemsets**
 - a set of *candidate* large k -itemsets.

Brute Force Method

- List all possible association rules
- Compute the support and confidence for each rule
- Prune rules that fail minimum support & minimum confidence thresholds
- Computationally expensive



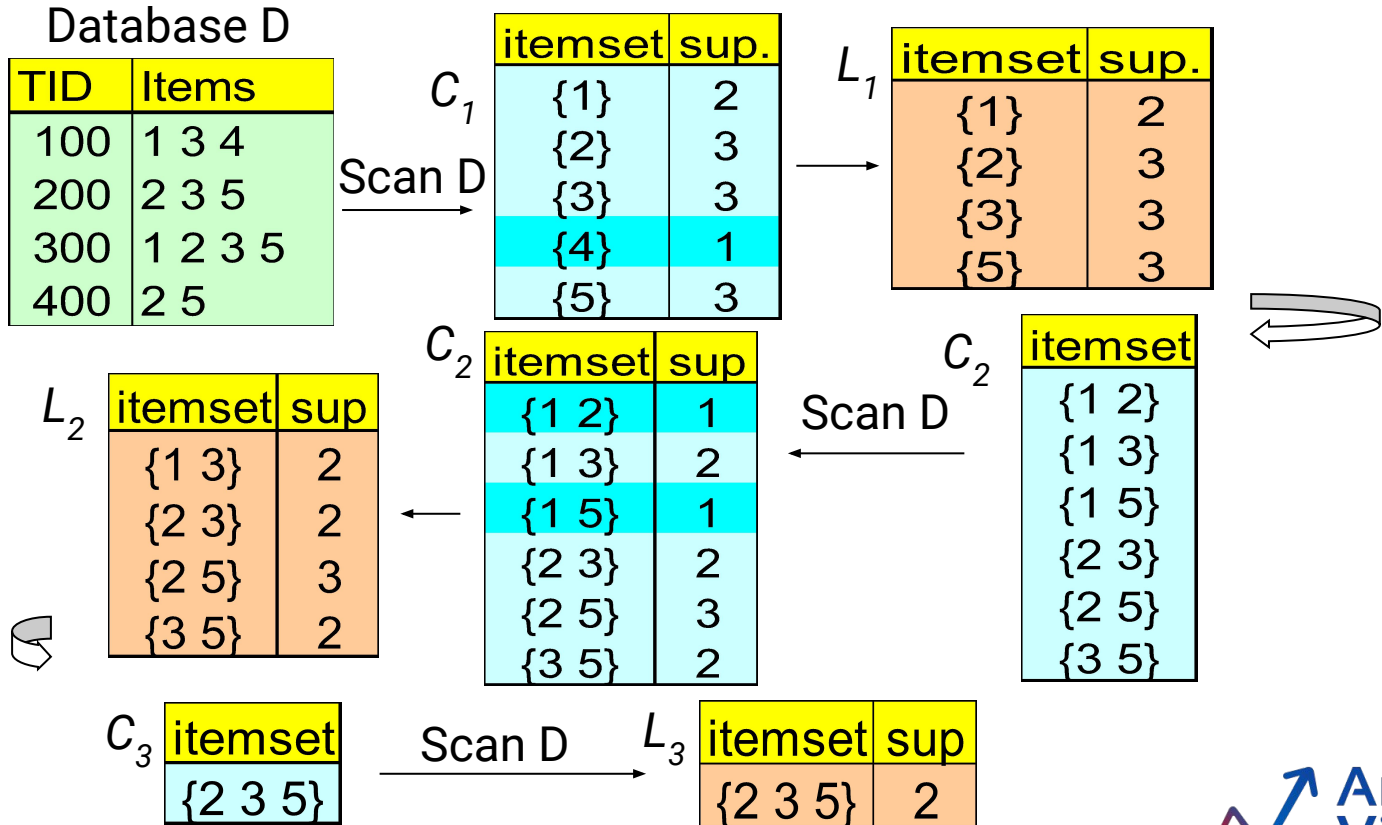
Basket Data

Apriori Algorithm for mining Association Rules

Apriori Algorithm

- There are many association rule mining algorithms
- Most Popular: Apriori Algorithm
 - Identifies the frequent individual items in the database
 - Extends them to larger and larger item sets if those itemsets appear sufficiently often in the database

Apriori Algorithm - Example



Generating Candidate Itemsets C₄

- Suppose these are the only 3-itemsets all have >10% support:
 {1, 2, 3}
 {1, 5, 7}
 {5, 6, 8}
 {5, 6, 11}
 {16, 17, 18}
- How do we generate candidate 4-itemsets that *might* have 10% support?

Generating Candidate Itemsets C₄

- Suppose these are the only 3-itemsets all have >10% support:
 {1, 2, 3}
 {1, 5, 7}
 {5, 6, 8}
 {5, 6, 11}
 {16, 17, 18}

Brute Force:

- Note all the items involved: {1, 2, 3, 5, 6, 7, 8, 11, 16, 17, 18}
- Generate all subsets of 4 of these:
 {1,2,3,5}, {1,2,3,6}, {1,2,3,7}, {1,2,3,8}, {1,2,3,11}, {1,2,3,16} etc ...
 there are 330 possible subsets in this case!

Generating Candidate Itemsets C₄

- Suppose these are the only 3-itemsets all have >10% support:
 - $\{1, 2, 3\}$
 - $\{1, 5, 7\}$
 - $\{5, 6, 8\}$
 - $\{5, 6, 11\}$
 - $\{16, 17, 18\}$
- We can easily see that $\{1, 2, 3, 5\}$ couldn't have 10% support – because $\{1, 2, 5\}$ is *not* one of our 3-itemsets
- Same goes for several other of these subsets

Apriori Trick

{1, 2, 3}

{1, 5, 7}

{5, 6, 8}

{5, 6, 11}

{16, 17, 18}

- Enforce that subsets are always arranged in an order (or similar), as they are already on the left
- **Only** generate $k+1$ -itemset candidates from k -itemsets that differ in the last item.
- So, in this case, the only candidate 4-itemset would be:

{5, 6, 8, 11}

Apriori Trick

This trick

- Guarantees to capture the itemsets that have enough support
- Will still generate some candidates that don't have enough support, so we still have to check them in the 'pruning' step,
- So for example we need to check if {5, 6, 8,11} has support greater than 10% or not
- If it does, algorithm will stop here as there is just 1 large itemset and no possibility of a 5-large itemset

Recommendation based on Association Rule Mining

- Simplest approach
 - Transform 5-point ratings into binary ratings (1 = above user average)
- Mine rules such as
 - Item1 \rightarrow Item5
 - support (2/4), confidence (2/2) (without Alice)
- Make recommendations for Alice (basic method)
 - Determine "relevant" rules based on Alice's transactions (the above rule will be relevant as Alice bought Item1)
 - Determine items not already bought by Alice
 - Sort the items based on the rules' confidence values

	Item1	Item2	Item3	Item4	Item5
Alice	1	0	0	0	?
User1	1	0	1	0	1
User2	1	0	1	0	1
User3	0	0	0	1	1
User4	0	1	1	0	0

Association Rule Mining: Formal Definition

- Commonly used for shopping behavior analysis
 - aims at detection of rules such as
"If a customer purchases baby food then he also buys diapers in 70% of the cases"
- Association rule mining algorithms
 - can detect rules of the form $X \rightarrow Y$ (e.g., beer \rightarrow diapers) from a set of sales transactions $D = \{t_1, t_2, \dots, t_n\}$
 - Here X is called antecedent & Y is called consequent & X,Y have no items in common
 - Each transaction from D will have information regarding the set of items bought together
 - measure of quality: support, confidence
 - used e.g. as a threshold to cut off unimportant rules