

This week

1. Intro to this week
2. Association Rule Analysis – Notes
3. Demo 1
4. Python Library – How to install
5. Demo 2 – Explanation -> Exercise to complete

- Assessment-A Q&A

- Assessment-B - Discussion

Follow this schedule

Do Not jump ahead

Follow me

Step by Step

Do Not jump ahead

TU 257 – Fundamentals of Data Science

Data Analytics

L9 – Association Rule Analysis

Brendan Tierney

Agenda

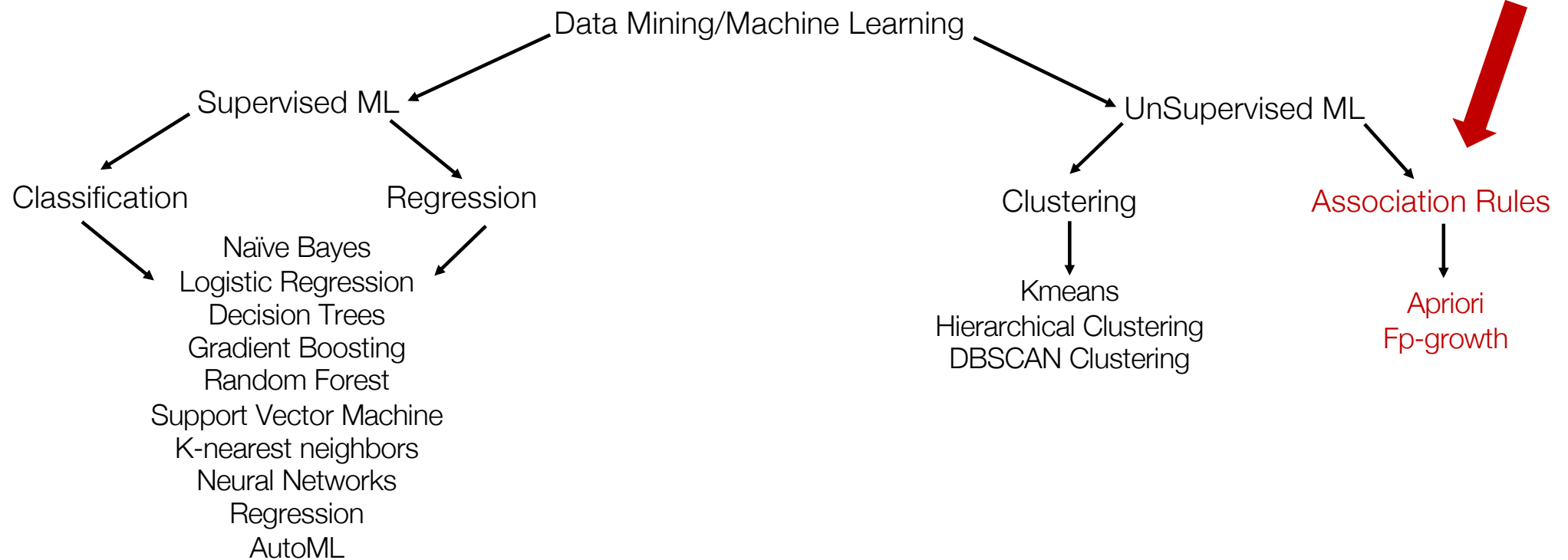
- Application Areas
- Market Basket Analysis
- A Simple Example
- A Big Search Space Problem
- Frequent Item Sets
- Apriori Algorithm
- How the Apriori Algorithms Example
- Evaluating & Filtering Rules (Support, Confidence & Lift)
- Data Privacy Issues



Machine Learning



- the use and development of computer systems that are able to learn and adapt without following explicit instructions, by using algorithms and statistical models to analyse and draw inferences from patterns in data.





theory **VS** reality

- Most Data Analytics etc can be done in a few lines of code
- Don't worry about the Theory, we might touch upon some of it, but it isn't necessary to know in-depth
- You'll never have to write an algorithm from scratch

Domain Knowledge



- Built up over time
- Experience
- Can take many years or decades
- In-depth knowledge of what is happening and why
- Can be based on Gut Feelings

Domain Knowledge



Domain Knowledge



- Built up over time
- Experience
- Can take many years or decades
- In-depth knowledge of what is happening and why
- Can be based on Gut Feelings

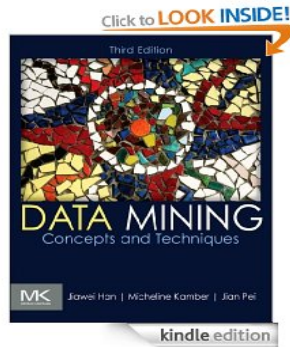
vs

- But what about the DATA
- What is the DATA telling you

Exercise

- Think about 2 items you regularly buy in your local shop, supermarket, convenience store, etc?
 - What are they?
 - What can you tell about their placement in the shop?
 - Why are they placed like that?

Start reading *Data Mining: Concepts and Techniques: Concepts and Techni...* on your Kindle **in under a minute**. Don't have a Kindle? [Get your Kindle](#)



Data Mining: Concepts and Techniques: Concepts and Techniques (The Morgan Kaufmann Series in Data Man

[Jiawei Han](#) (Author), [Micheline Kamber](#) (Author), [Jian Pei](#) (Author)

★☆☆☆☆ (1 customer review)

Print List Price: ~~£45.99~~

Kindle Price: **£28.42** includes VAT* & free wireless delivery via **Amazon Whispernet**

You Save: **£17.57 (38%)**

* Unlike print books, digital books are subject to VAT.

- Length: 626 pages (Contains Real Page Numbers)
- Don't have a Kindle? [Get your Kindle here](#) or start reading now with a free [Kindle Reading App](#).

Formats	Amazon Price	New from	Used from
Kindle Edition	£28.42	--	--
Hardcover	£38.41	£21.17	£32.31

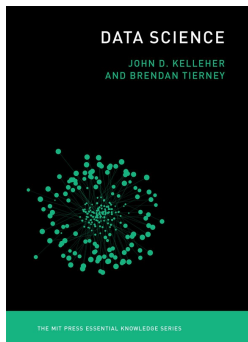
[Share your own customer images](#)

Trade-In Store

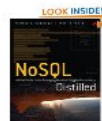
Amazon.co.uk Trade-In Store

Did you know you can trade in your unwanted old books for an Amazon.co.uk Gift Card? Plus, until February 23, 2014, get an extra £5 Promotion. [Books Trade-In Store](#) for more details. [Learn more](#).

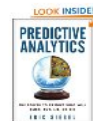
Customers Who Bought This Item Also Bought



Data Mining Techniques: For Marketing, Sales, and Customer Relationship Management
 > [Jiawei Han](#), [Micheline Kamber](#), [Jian Pei](#)
 Kindle Edition
 £20.39



NoSQL Distilled: A Brief Guide to the Emerging World of Distributed Database Systems
 > [Pramod J. Sadalage](#), [Martin Fowler](#)
 Kindle Edition
 ★★★★★ (7)
 £12.99



Predictive Analytics: The Power to Predict Who Will Buy What
 > [Eric Siegel](#)
 Kindle Edition
 ★★★★★ (9)
 £11.45



Seven Databases in Seven Weeks: A Guide to Modern Databases and How to Use Them
 > [Eric Redmond](#)
 Kindle Edition
 ★★★★★ (3)
 £18.81



Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython
 > [Wes McKinney](#)
 Kindle Edition
 ★★★★★ (5)
 £12.67



The Data Warehouse Toolkit: The Complete Guide to Dimensional Modeling
 > [Ralph Kimball](#)
 Kindle Edition
 ★★★★★ (14)
 £27.81



DON'T FORGET TO USE YOUR COUPONS
VALID FOR A LIMITED TIME ONLY!

<p>25P off when you buy any any Harriott soup</p> <p>AINSLEY HARRIOTT</p> <p>Use your Clubcard with this coupon Valid until 13/03/11</p>	<p>£1 off when you spend £5 or more on celebration products Includes cards, wrap, party bags, balloons & table decorations</p> <p>TESCO <i>Every little helps</i></p> <p>Use your Clubcard with this coupon Valid until 13/03/11</p>
<p>25 extra points when you buy any any Harriott soup</p> <p>25 extra points when you spend £1 or more on confectionery</p>	<p>25 extra points when you spend £1 or more on fresh savoury snack items Includes pies, pastries, quiche & scotch eggs</p>
<p>5P off when you buy any Finest catballs 180g</p> <p>Finest★</p> <p>Use your Clubcard with this coupon Valid until 13/03/11</p>	<p>£2.50 off when you spend £10 or more on any kids' clothing Includes back to school</p> <p>TESCO <i>Every little helps</i></p> <p>Use your Clubcard with this coupon Valid until 13/03/11</p>
<p>25 extra points when you buy any any condiments</p>	<p>50 bonus Airmiles when you exchange a minimum of £15 in Clubcard Vouchers for Airmiles in one transaction at tesco.com/rewards. Exchange by 31/03/11. Visit www.airmiles.co.uk/rewards for more information See overleaf for Terms & Conditions</p> <p>AIR MILES MAKE YOUR MONEY FLY</p> <p>Valid until 31/03/11</p>

25 extra points
when you buy any
any
confectionery

25 extra points
when you spend
£1 or more on
fresh savoury snack items
Includes pies, pastries, quiche & scotch eggs

25 extra points
when you buy any
any
wine
Includes sparkling & fortified wines

50P off
when you buy any
Spam
Includes new Spam chopped pork
& ham with bacon 200g

SPAM

25 extra points
when you buy any
any
condiments

£2.00

£1.50

£1.50

In-store • Online • F

MumsMonkey2012

What Is Association Mining?

- Association rule mining:

- Finding frequent patterns, associations, correlations, or causal structures among sets of items or objects in transaction databases, relational databases, and other information repositories

Frequent Pattern: A pattern (set of items, sequence, etc.) that occurs frequently in a database

Examples of Application Areas


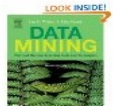





- Beer and Nappies Example

- Others include

- What are the subsequent purchases after buying a PC?
- What kinds of DNA are sensitive to this new drug?
- Can we automatically classify web documents?
- Basket data analysis, cross-marketing, catalog design, sale campaign analysis
- Web log (click stream) analysis, DNA sequence analysis, etc

Customers Who Bought This Item Also Bought

Page

 Data Mining Techniques: For Marketing, Sales, ... Gordon S. Linoff ★★★★☆ (11) Paperback £22.09	 Data Mining: Practical Machine Learning Tools ... I. H. Witten ★★★★☆ (3) Paperback	 Information Theory, Inference and Learning ... David J. C. MacKay ★★★★★ (5) Hardcover £35.20	 Data Mining with Microsoft SQL Server 2008 Jamie MacLennan ★★★★★ (1) Paperback £22.09	 Mining the Social Web: Analyzing Data from ... Matthew A. Russell ★★★★☆ (2) Paperback £20.14	 Handbook of Statistical Analysis and Data Mining Applications Nisbet Hardcover £42.74	 Pattern Recognition and Machine Learning ... Christopher M. Bishop ★★★★☆ (8) Hardcover £56.31
--	---	--	---	--	--	---

Another Example

“A bank’s marketing department is interested in examining associations between various retail banking services used by customers. They would like to determine both typical and atypical service combinations”

- The **BANK** data set has over 32,000 rows coming from 8,000 customers. Each row of the data set represents a customer-service combination. Therefore, a single customer can have multiple rows in the data set, each row representing one of the products he or she owns. The median number of products per customer is three

Name	Model Role	Measurement Level	Description
ACCOUNT	ID	Nominal	Account Number
SERVICE	Target	Nominal	Type of Service
VISIT	Sequence	Ordinal	Order of Product Purchase

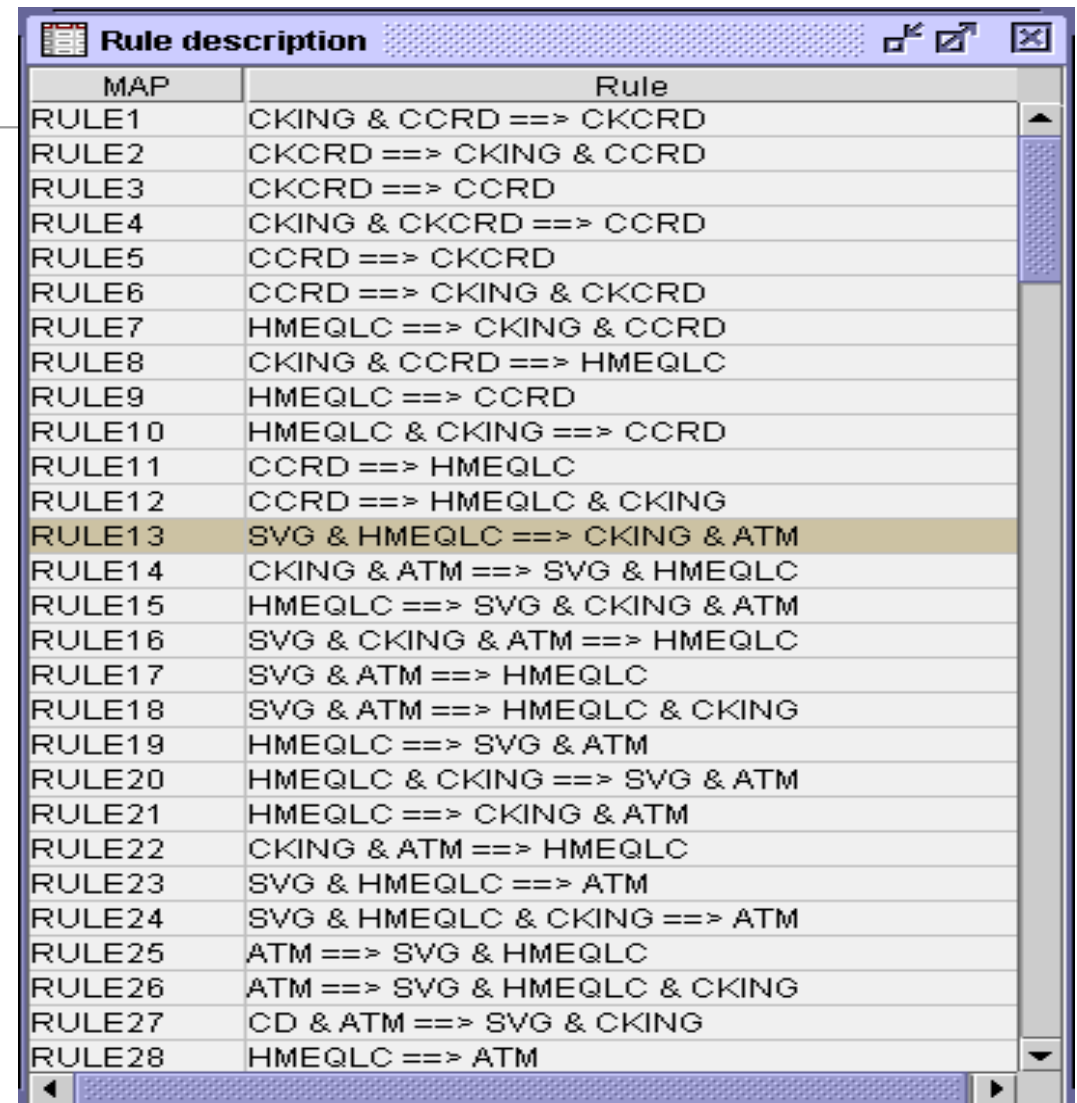
Case Study

- The 13 products are represented in the data set as follows:
 - ATM automated teller machine debit card
 - AUTO automobile installment loan
 - CCRD credit card
 - CD certificate of deposit
 - CKCRD check/debit card
 - CKING checking account
 - HMEQLC home equity line of credit
 - IRA individual retirement account
 - MMDA money market deposit account
 - MTG mortgage
 - PLOAN personal/consumer installment loan
 - SVG saving account
 - TRUST personal trust account

Case Study (cont...)

- Rules generated by analysis

- ATM automated teller machine debit card
- AUTO automobile installment loan
- CCRD credit card
- CD certificate of deposit
- CKCRD check/debit card
- CKING checking account
- HMEQLC home equity line of credit
- IRA individual retirement account
- MMDA money market deposit account
- MTG mortgage
- PLOAN personal/consumer installment loan
- SVG saving account
- TRUST personal trust account



MAP	Rule
RULE1	CKING & CCRD ==> CKCRD
RULE2	CKCRD ==> CKING & CCRD
RULE3	CKCRD ==> CCRD
RULE4	CKING & CKCRD ==> CCRD
RULE5	CCRD ==> CKCRD
RULE6	CCRD ==> CKING & CKCRD
RULE7	HMEQLC ==> CKING & CCRD
RULE8	CKING & CCRD ==> HMEQLC
RULE9	HMEQLC ==> CCRD
RULE10	HMEQLC & CKING ==> CCRD
RULE11	CCRD ==> HMEQLC
RULE12	CCRD ==> HMEQLC & CKING
RULE13	SVG & HMEQLC ==> CKING & ATM
RULE14	CKING & ATM ==> SVG & HMEQLC
RULE15	HMEQLC ==> SVG & CKING & ATM
RULE16	SVG & CKING & ATM ==> HMEQLC
RULE17	SVG & ATM ==> HMEQLC
RULE18	SVG & ATM ==> HMEQLC & CKING
RULE19	HMEQLC ==> SVG & ATM
RULE20	HMEQLC & CKING ==> SVG & ATM
RULE21	HMEQLC ==> CKING & ATM
RULE22	CKING & ATM ==> HMEQLC
RULE23	SVG & HMEQLC ==> ATM
RULE24	SVG & HMEQLC & CKING ==> ATM
RULE25	ATM ==> SVG & HMEQLC
RULE26	ATM ==> SVG & HMEQLC & CKING
RULE27	CD & ATM ==> SVG & CKING
RULE28	HMEQLC ==> ATM

What are the most interesting findings ?

Case Study (cont...)

- The most interesting findings from the analysis included:
 - The strongest rule is **checking**, and **credit card** implies **check card**.
 - This is **not surprising** given that many check cards include credit card logos
 - It appears that customers with auto loans typically have checking and savings accounts (and are ATM users),
 - but do not utilize other services (at least with sufficient support and confidence to be included in the presented analysis)
- Patterns in the data are discovered. Your job is to put a meaning on the patterns.
 - This isn't always possible

Association Rule: Basic Concepts

- Given: (1) **database of transactions**, (2) **each transaction is a list of items** (purchased by a customer in a visit)
- **Find:** **all rules that correlate the presence of one set of items with that of another set of items**
 - E.g., *98% of people who purchase tires and auto accessories also get automotive services done*
- Applications
 - *Maintenance Agreement* (What the store should do to boost Maintenance Agreement sales)
 - *Home Electronics* (What other products should the store stocks up?)
 - Attached mailing in direct marketing
 - Detecting “ping-pong”ing of patients, faulty “collisions”

Market Basket Analysis

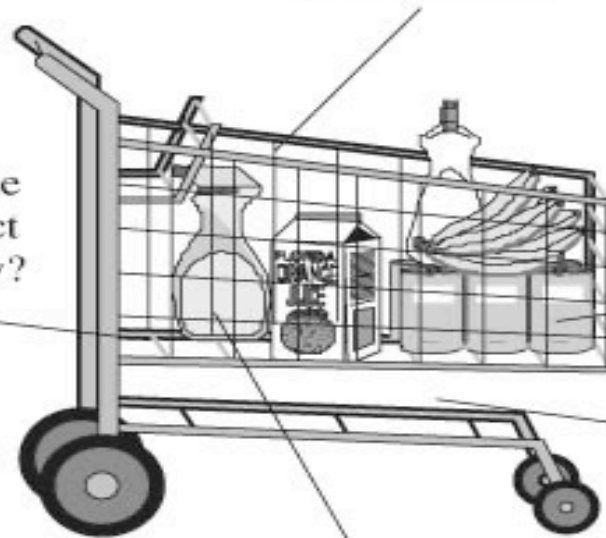
- **Market basket analysis** is a typical example of frequent itemset mining
- Customers **buying habits** are divined by finding associations between different items that customers place in their “shopping baskets”
- This information **can be used** to develop marketing strategies
- One basket tells you about what one customer purchased at one time
- A **loyalty cards** make it possible to tie together purchases by a single customer (or household) over time



In this shopping basket, the shopper purchased a quart of orange juice, some bananas, dish detergent, some window cleaner, and a six pack of soda.

How do the demographics of the neighborhood affect what customers buy?

Is soda typically purchased with bananas? Does the brand of soda make a difference?



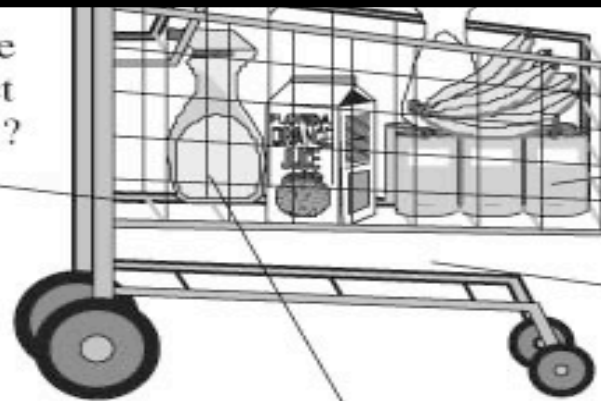
Are window cleaning products purchased when detergent and orange juice are bought together?

What should be in the basket but is not?

In this shopping basket, the shopper purchased a quart of orange juice, some bananas, dish detergent, some window cleaner, and a six

How have these rules / patterns changed over the past few years ?

How do the demographics of the neighborhood affect what customers buy?



Are window cleaning products purchased with bananas? Does the brand of soda make a difference?

Are window cleaning products purchased when detergent and orange juice are bought together?

What should be in the basket but is not?

- **more than just the contents of shopping carts**

- It is also about what customers do not purchase, **and why**.

- If customers purchase baking powder, but no flour, **what are they baking?**

- If customers purchase a mobile phone, but no case, are you **missing an opportunity?**

- **Are they a drug dealer?**

- It is also about key drivers of purchases; for example, the gourmet mustard that seems to lie on a shelf collecting dust until a customer buys that particular brand of special gourmet mustard in a shopping excursion that includes hundreds of dollars' worth of other products. Would eliminating the mustard (to replace it with a better-selling item) threaten the entire customer relationship?

Data has a Value

In the USA

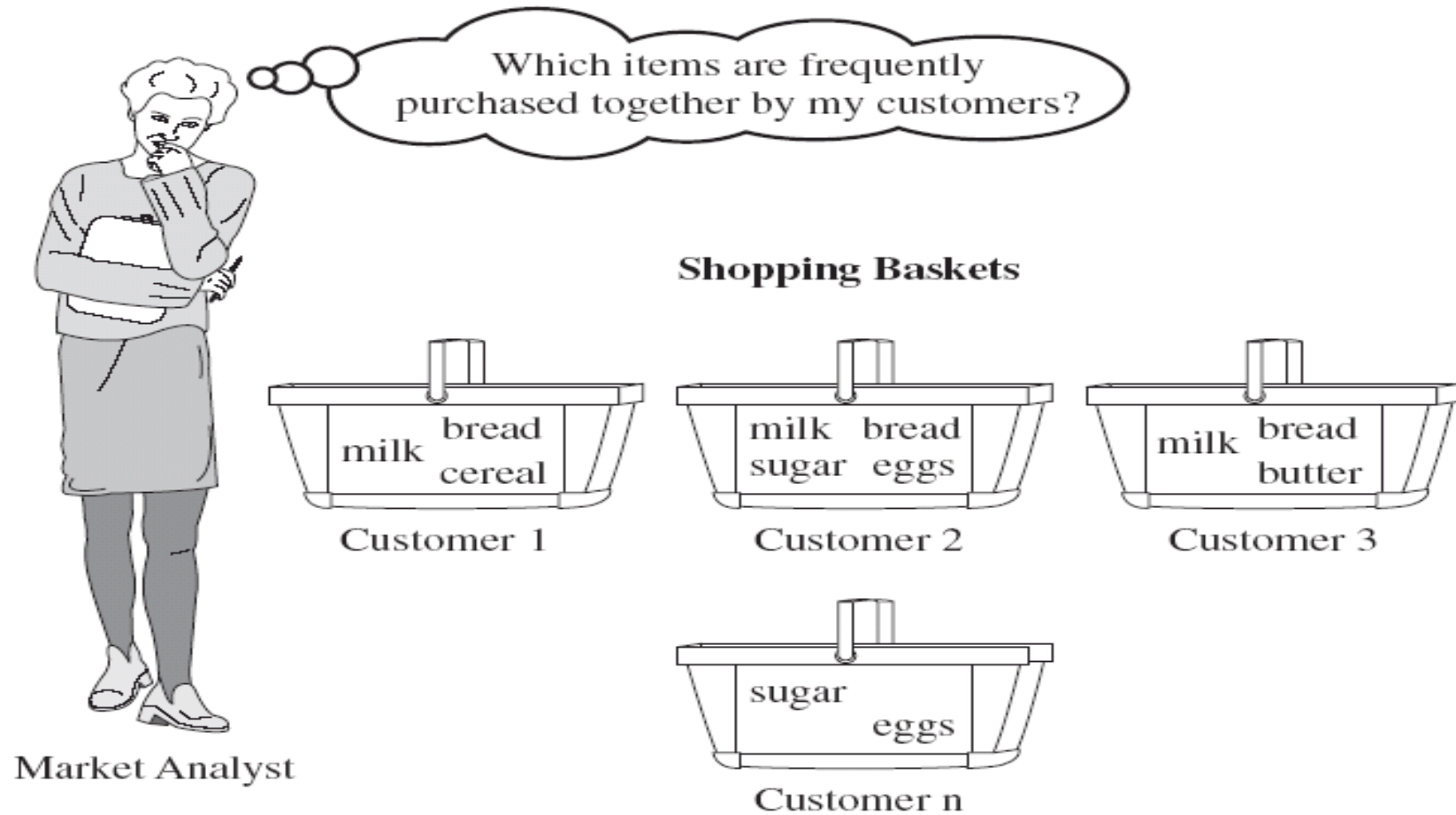
Some retail stores are **making more money out of your data**

(using it, selling it, etc.)

than they are through their normal day-to-day business

Be careful what you agree to when signing up to “Special Offers”

Quick Exercise – What are the most commonly bought combination of products



Association Rules: The Problem of Lots of Data

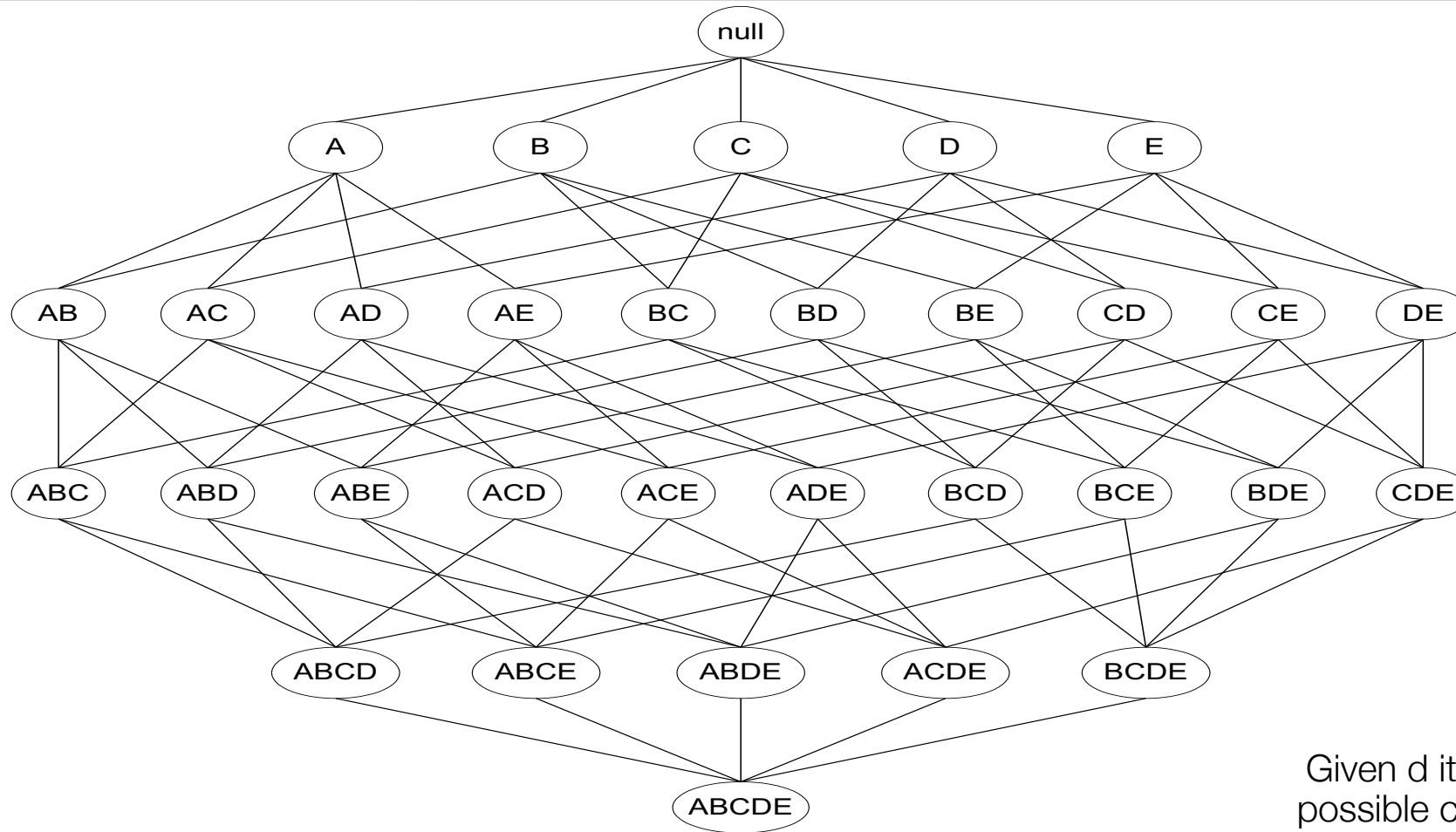
- Fast Food Restaurant...could have 100 items on its menu
 - How many combinations are there with 3 different menu items?
 - 161,700 !
- Supermarket...10,000 or more unique items
 - 50 million 2-item combinations
 - 100 billion 3-item combinations
- Use of product hierarchies (groupings) helps address this common issue
- Finally, know that the number of transactions in a given time-period could also be huge (hence expensive to analyse)

Try writing some SQL queries to find frequently items sets for these situations!

Itemsets & Frequent Itemsets

- An **itemset** is a set of items
- A **k -itemset** is an itemset that contains k items
- The **occurrence frequency of an itemset** is the number of transactions that contain the itemset
 - This is also known more simply as the **frequency, support count or count**
- An itemset is said to be **frequent** if the support count satisfies a **minimum support count threshold**
- The set of frequent itemsets is denoted L_k

Frequent Itemset Generation



Given d items, there are 2^d possible candidate itemsets

Association Rule Mining

- In general association rule mining can be reduced to the following two steps:
 1. Find all frequent itemsets
 - Each itemset will occur at least as frequently as as a minimum support count
 2. Generate strong association rules from the frequent itemsets
 - These rules will satisfy minimum support and confidence measures

- Apriori pruning principle: If there is any itemset which is infrequent, its superset should not be generated/tested!



The Apriori Algorithm

- Join Step: C_k is generated by joining L_{k-1} with itself
- Prune Step: Any $(k-1)$ -itemset that is not frequent cannot be a subset of a frequent k -itemset

- Pseudo-code:

C_k : Candidate itemset of size k

L_k : frequent itemset of size k

$L_1 = \{\text{frequent items}\};$

for $(k = 1; L_k \neq \emptyset; k++)$ **do begin**

C_{k+1} = candidates generated from L_k ;

for each transaction t in database **do**

increment the count of all candidates in C_{k+1} that are contained in t

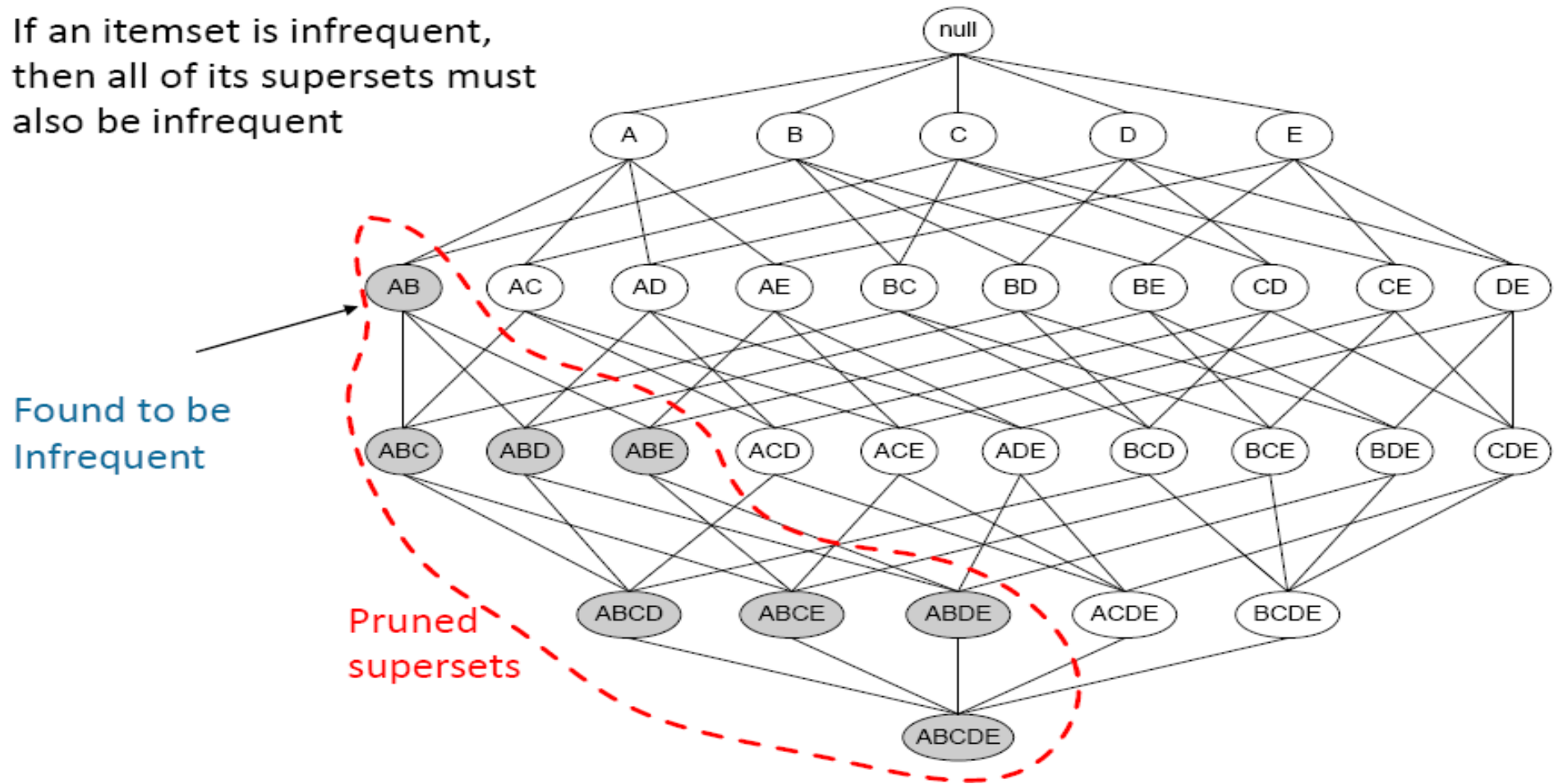
L_{k+1} = candidates in C_{k+1} with min_support

end

return $\cup_k L_k$;

Apriori principle for pruning candidates

If an itemset is infrequent, then all of its supersets must also be infrequent



Generating Association Rules

- Once all **frequent itemsets** have been **found**, the **association rules** can be generated
- Strong association rules from a frequent itemset are generated by calculating the confidence in each possible rule arising from that itemset and testing it against a minimum confidence threshold

The Apriori Algorithm — Example

Database D

TID	Items
100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5

Scan D

C_1

itemset	sup.
{1}	2
{2}	3
{3}	3
{4}	1
{5}	3

L_1

itemset	sup.
{1}	2
{2}	3
{3}	3
{5}	3

L_2

itemset	sup
{1 3}	2
{2 3}	2
{2 5}	3
{3 5}	2

C_2

itemset	sup
{1 2}	1
{1 3}	2
{1 5}	1
{2 3}	2
{2 5}	3
{3 5}	2

Scan D

C_2

itemset
{1 2}
{1 3}
{1 5}
{2 3}
{2 5}
{3 5}

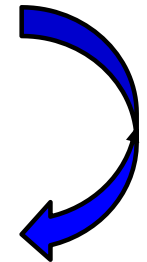
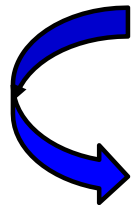
C_3

itemset
{2 3 5}

Scan D

L_3

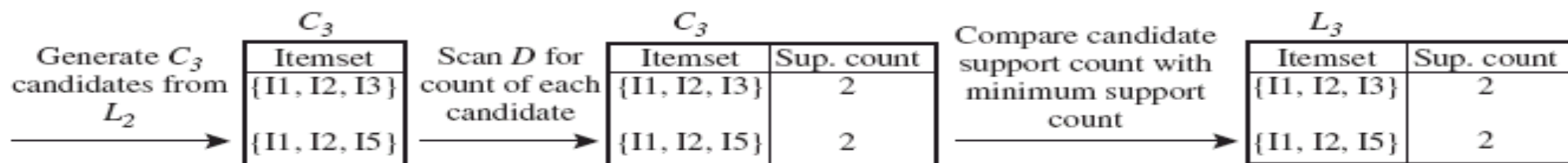
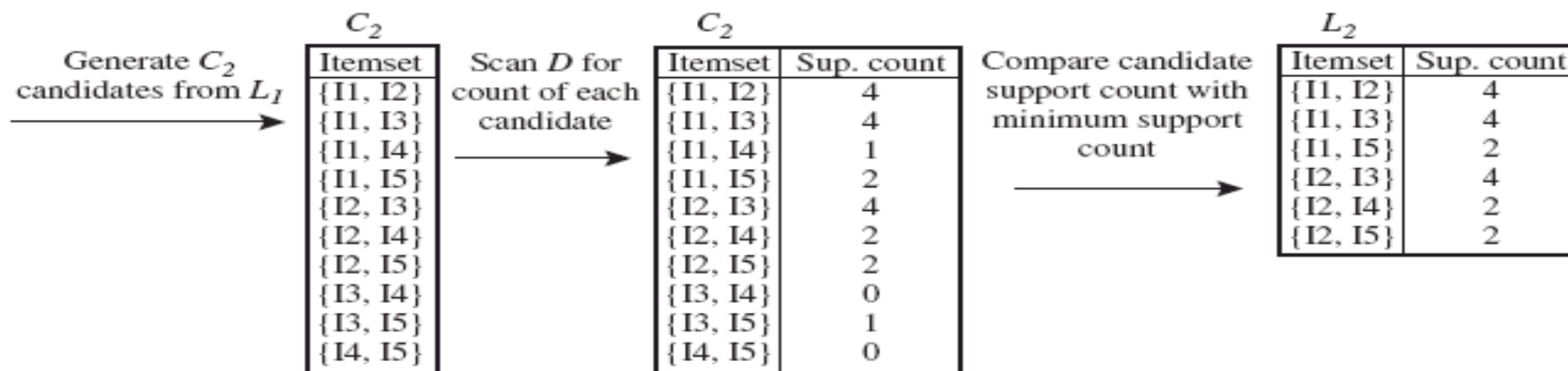
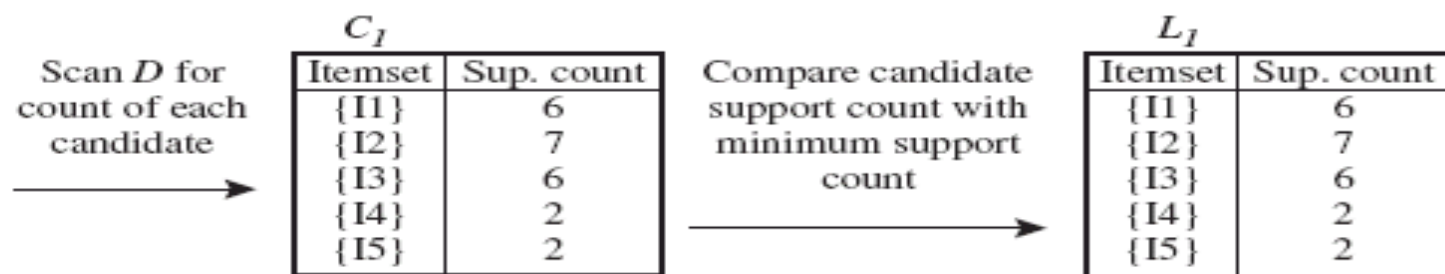
itemset	sup
{2 3 5}	2



Another Example

TID	List of item_IDs
T100	Beer, Crisps, Milk
T200	Crisps, Bread
T300	Crisps, Nappies
T400	Beer, Crisps, Bread
T500	Beer, Nappies
T600	Crisps, Nappies
T700	Beer, Nappies
T800	Beer, Crisps, Nappies, Milk
T900	Beer, Crisps, Nappies

ID	Item
I1	Beer
I2	Crisps
I3	Nappies
I4	Bread
I5	Milk



Association Rule Support & Confidence

- We say that an association rule $A \Rightarrow B$ holds in the transaction set D with **support**, S , and **confidence**, C
- The **support** of the association rule is given as the percentage of transactions in D that contain both A and B (or $A \cup B$)
 - So, the support can be considered the probability $P(A \cup B)$
- The **confidence** of the association rule is given as the percentage of transactions in D containing A that also contain B
- So, the confidence can be considered the conditional probability $P(B | A)$
- Association rules that satisfy minimum support and confidence values are said to be **strong**

Support & Confidence Again

- Support and confidence values can be calculated as follows:

$$\begin{aligned} \text{support}(A \Rightarrow B) &= P(A \cup B) \\ &= \frac{\text{support_count}(A \cup B)}{\text{count}(\quad)} \end{aligned}$$

$$\begin{aligned} \text{confidence}(A \Rightarrow B) &= P(B \mid A) \\ &= \frac{\text{support}(A \cup B)}{\text{support}(A)} \\ &= \frac{\text{support_count}(A \cup B)}{\text{support_count}(A)} \end{aligned}$$

Rule Measures: Support and Confidence

- Find all the rules $X \& Y \Rightarrow Z$ with minimum confidence and support
 - **support, s** , probability that a transaction contains $\{X \& Y \& Z\}$
 - **confidence, c** , conditional probability that a transaction having $\{X \& Y\}$ also contains Z

Transaction ID	Items Bought
2000	A, B, C
1000	A, C
4000	A, D
5000	B, E, F

Let minimum support 50%, and minimum confidence 50%, we have

$A \Rightarrow C$ (50%, 66.6%)

$C \Rightarrow A$ (50%, 100%)

Rule Measures: Support and Confidence

- Find all the rules $X \& Y \Rightarrow Z$ with minimum confidence and support
 - support, s , probability that a transaction contains $\{X \& Y \& Z\}$
 - confidence, c , conditional probability that a transaction having $\{X \& Y\}$ also contains Z

Transaction ID	Items Bought
2000	A, B, C
1000	A, C
4000	A, D
5000	B, E, F

Let minimum support 50%, and minimum confidence 50%, we have

$A \Rightarrow C$ (50%, 66.6%)

$C \Rightarrow A$ (50%, 100%)

Rule Measures: Support and Confidence

- Find all the rules $X \& Y \Rightarrow Z$ with minimum confidence and support
 - support, s , probability that a transaction contains $\{X \& Y \& Z\}$
 - confidence, c , conditional probability that a transaction having $\{X \& Y\}$ also contains Z

Transaction ID	Items Bought
2000	A, B, C
1000	A, C
4000	A, D
5000	B, E, F

Let minimum support 50%, and minimum confidence 50%, we have

$A \Rightarrow C$ (50%, 66.6%)

$C \Rightarrow A$ (50%, 100%)

Mining Association Rules: An Example

Transaction-id	Items bought
10	A, B, C
20	A, C
30	A, D
40	B, E, F

Frequent pattern	Support
{A}	75%
{B}	50%
{C}	50%
{A, C}	50%


$$\begin{aligned} \text{support}(A \Rightarrow C) &= \frac{\text{support_count}(\{A\} \cup \{C\})}{\text{count}()} \\ &= 50\% \end{aligned}$$

$$\begin{aligned} \text{confidence}(A \Rightarrow C) &= \frac{\text{support_count}(\{A\} \cup \{C\})}{\text{support_count}(\{A\})} \\ &= 66.7\% \end{aligned}$$

The **Apriori** principle:
Any subset of a frequent itemset must be frequent

Mining Association Rules: An Example

Transaction-id	Items bought
10	A, B, C
20	A, C
30	A, D
40	B, E, F



Frequent pattern	Support
{A}	75%
{B}	50%
{C}	50%
{A, C}	50%

$$\begin{aligned} \text{support}(C \Rightarrow A) &= \frac{\text{support_count}(\{C\} \cup \{A\})}{\text{count}()} \\ &= 50\% \end{aligned}$$

$$\begin{aligned} \text{confidence}(C \Rightarrow A) &= \frac{\text{support_count}(\{C\} \cup \{A\})}{\text{support_count}(\{C\})} \\ &= 100\% \end{aligned}$$

Support

- The rule $X \Rightarrow Y$ holds with supports if $s\%$ of transactions in D contain XUY .

TID	Items	Support = Occurrence / Total Support
1	ABC	Total Support = 5 Support {AB} = $2 / 5 = 40\%$ Support {BC} = $3 / 5 = 60\%$ Support {ABC} = $1 / 5 = 20\%$
2	ABD	
3	BC	
4	AC	
5	BCD	

- Rules that have a S greater than a user-specified support is said to have minimum support.
- **Support:** Support of a rule is a measure of how frequently the items involved in it occur together. Using probability notation: support (A implies B) = $P(A, B)$.

Confidence

- The rule $X \Rightarrow Y$ holds with confidence c if $c\%$ of the transactions in D that contain X also contain Y .

TID	Items	Given $X \Rightarrow Y$	
1	ABC	Confidence = Occurrence {Y} / Occurrence {X}	
2	ABD		
3	BC		Confidence $\{A \Rightarrow B\} = 2 / 3 = 66\%$
4	AC		Confidence $\{B \Rightarrow C\} = 3 / 4 = 75\%$
5	BCD		Confidence $\{AB \Rightarrow C\} = 1 / 2 = 50\%$

- Rules that have a C greater than a user-specified confidence is said to have minimum confidence.
- **Confidence:** Confidence of a rule is the conditional probability of B given A . Using probability notation: confidence (A implies B) = $P(B \text{ given } A)$.

$$\text{lift}(X \Rightarrow Y) = \frac{\text{supp}(X \cup Y)}{\text{supp}(X) \times \text{supp}(Y)}$$

Lift

- Lift
 - Lift indicates the strength of a rule over the random co-occurrence of the antecedent and the consequent, given their individual support. It provides information about the improvement, the increase in probability of the consequent given the antecedent. Lift is defined as follows.
 - $(\text{Rule Support}) / (\text{Support}(\text{Antecedent}) * \text{Support}(\text{Consequent}))$
 - This can also be defined as the confidence of the combination of items divided by the support of the consequent.
- Example
 - Convenience store customers who buy orange juice also buy milk with a 75% confidence.
 - The combination of milk and orange juice has a support of 30%.

 - This at first sounds like an excellent rule, and in most cases, it would be. It has high confidence and high support. However, what if convenience store customers in general buy milk 90% of the time? In that case, orange juice customers are actually less likely to buy milk than customers in general.

 - in our milk example, assuming that 40% of the customers buy orange juice, the improvement would be:
 - $30\% / (40\% * 90\%) = 0.83$ – an improvement of less than 1.
 - Any rule with an improvement of less than 1 does not indicate a real cross-selling opportunity, no matter how high its support and confidence, because it actually offers less ability to predict a purchase than does random chance.

 - If lift > 1, then items are positively correlated
 - lift < 1, then negatively correlated
 - lift = 1, then are independent

Sequence Databases and Sequential Pattern Analysis

- Frequent patterns vs. (frequent) sequential patterns
- Applications of sequential pattern mining
 - Customer shopping sequences:
 - First buy computer, then CD-ROM, and then digital camera, within 3 months.
 - Medical treatment
 - Natural disasters (e.g., earthquakes)
 - Science & engineering processes,
 - Stocks Markets
 - Telephone calling patterns,
 - Weblog click streams
 - DNA sequences and gene structures
 - ...



But

What about **privacy issues**?

Does this matter?

Forbes

TECH | 2/16/2012 @ 11:02AM | 2,337,246 views

How Target Figured Out A Teen Girl Was Pregnant Before Her Father Did

334 comments, 173 called-out

+ Comment Now + Follow Comments

Every time you go shopping, you share intimate details about your consumption patterns with retailers. And many of those retailers are studying those details to figure out what you like, what you need, and which coupons are most likely to make you happy. Target, for example, has figured out how to data-mine its way into your womb, to figure out whether you have a baby on the way long before you need to start buying diapers.

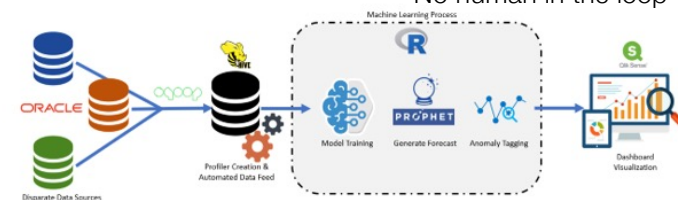
Charles Duhigg outlines in the New York Times how Target tries to hook parents-to-be at that crucial moment before they turn into rampant – and loyal – buyers of all things pastel, plastic, and miniature. He talked to Target statistician Andrew Pole – before Target freaked out and cut off all communications – about the clues to a customer’s impending bundle of joy. Target assigns every customer a Guest ID number, tied to their credit card, name, or email address that becomes a bucket that stores a history of everything they’ve bought and any demographic information Target has collected from them or bought from other sources. Using that, Pole looked at historical buying data for all the ladies who had signed up for Target baby registries in the past. From the NYT:



Target has got you in its aim



Automated Pattern Discovery
No human in the loop



<https://www.forbes.com/sites/kashmirhill/2012/02/16/how-target-figured-out-a-teen-girl-was-pregnant-before-her-father-did/?sh=7385e4a16668>

Maciek Wasiak commented on this



Anthony O'Neill

Director, Analytics Centre of Excellence at eir
10 hrs

When predictive analytics goes wrong...

Jac Rayner (@GirlFromBlupo) tweeted at 8:22 AM on Fri, Apr 06, 2018:

Dear Amazon, I bought a toilet seat because I needed one. Necessity, not desire. I do not collect them. I am not a toilet seat addict. No matter how temptingly you email me, I'm not going to think, oh go on then, just one more toilet seat, I'll treat myself.

13 Likes • 1 Comment



Like



Comment



Share



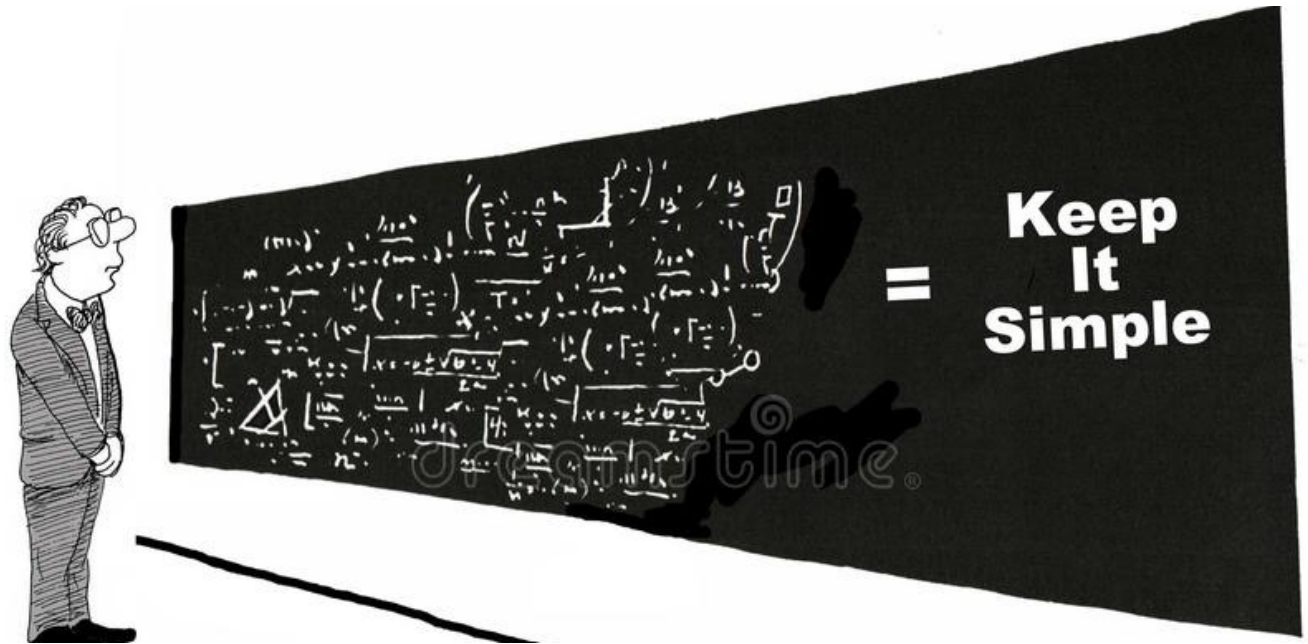
Maciek Wasiak Lol, recommender systems can be tricky to build but by now everyone would expect from Amazon to notice the difference between one-off and repeat purchases. How many garden sheds do I need, like...



Like



Reply





Time for an
Example

Any Questions ?

What Now/Next ?