

# 数据挖掘和大数据分析



# Outline

① Review Association Rules

② Apriori Algorithm



③ Lift & Association Rule Example



## Which knowledge is used in this example?

- A Association Rules of DM
- B Classification of DM
- C Regression of DM

### Frequently Bought Together



Price For All Three: **\$166.83**



Add all three to Cart

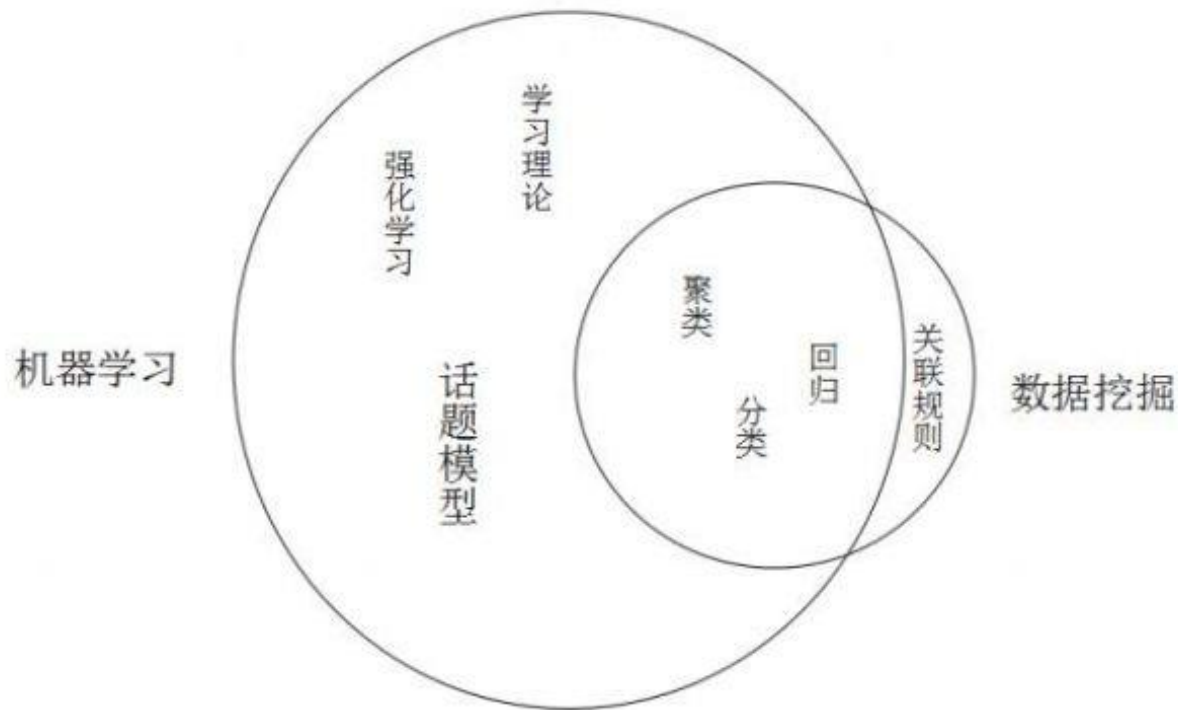
Add all three to Wish List

[Show availability and shipping details](#)

- ☒ **This item:** Data Mining: Practical Machine Learning Tools and Techniques, Second Edition (Morgan Kaufmann Series in Data Management Systems) by Eibe Frank
- ☒ [The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition \(Springer Series in Statistics\)](#) by Robert Tibshirani

提交

# Review Association Rules



# Review Association Rules

## Frequent Itemsets



## Association Rules



Support( Bread  $\rightarrow$  Milk ) = ?

Confidence( Bread  $\rightarrow$  Milk ) = ?

- ☐ A 1/8 & 1/3
- ☐ B 2/8 & 2/3
- ☒ C 2/8 & 1/3
- ☐ D 2/8 & 2/4

| Transactions | Items                        |
|--------------|------------------------------|
| 1            | Bread, Jelly, Peanut, Butter |
| 2            | Bread, Butter                |
| 3            | Bread, Jelly                 |
| 4            | Bread, Milk, Butter          |
| 5            | Chips, Milk                  |
| 6            | Bread, Chips                 |
| 7            | Bread, Milk                  |
| 8            | Chips, Jelly                 |

提交

Support( Milk  $\rightarrow$  Bread ) = ?

Confidence( Milk  $\rightarrow$  Bread ) = ?

- ☐ A 1/8 & 1/3
- ☒ B 2/8 & 2/3
- ☐ C 2/8 & 1/3
- ☐ D 3/8 & 1/3

| Transactions | Items                        |
|--------------|------------------------------|
| 1            | Bread, Jelly, Peanut, Butter |
| 2            | Bread, Butter                |
| 3            | Bread, Jelly                 |
| 4            | Bread, Milk, Butter          |
| 5            | Chips, Milk                  |
| 6            | Bread, Chips                 |
| 7            | Bread, Milk                  |
| 8            | Chips, Jelly                 |

提交

# Review Association Rules

| Transactions | Items                        |
|--------------|------------------------------|
| 1            | Bread, Jelly, Peanut, Butter |
| 2            | Bread, Butter                |
| 3            | Bread, Jelly                 |
| 4            | Bread, Milk, Butter          |
| 5            | Chips, Milk                  |
| 6            | Bread, Chips                 |
| 7            | Bread, Milk                  |
| 8            | Chips, Jelly                 |

| Itemset | Support |
|---------|---------|
| Bread   | 6/8     |
| Butter  | 3/8     |
| Chips   | 2/8     |
| Jelly   | 3/8     |
| Milk    | 3/8     |
| Peanut  | 1/8     |

Bread → Milk

**Support:** 2/8

**Confidence:** 1/3

Milk → Bread

**Support:** 2/8

**Confidence:** 2/3

Searching for rules in the form of: **Bread** → **Butter**



# Review Association Rules

## ❖ Support and Confidence are bounded by thresholds:

- Minimum support  $\sigma$
- Minimum confidence  $\Phi$

❖ A frequent (large) itemset is an itemset with support larger than  $\sigma$ .


❖ **A strong rule is a rule that is frequent and its confidence is higher than  $\Phi$ .**

## ❖ Association Rule Problem

- Given **I, D,  $\sigma$  and  $\Phi$** , to find all strong rules in the form of  **$X \rightarrow Y$** .

❖ **The number of all possible association rules is huge.**

- Brute force strategy is infeasible.
- A smart way is to find frequent itemsets first.



# **DATA ANALYTICS:** **DATA MINING AND BIG DATA**



—— Mining Rules 2

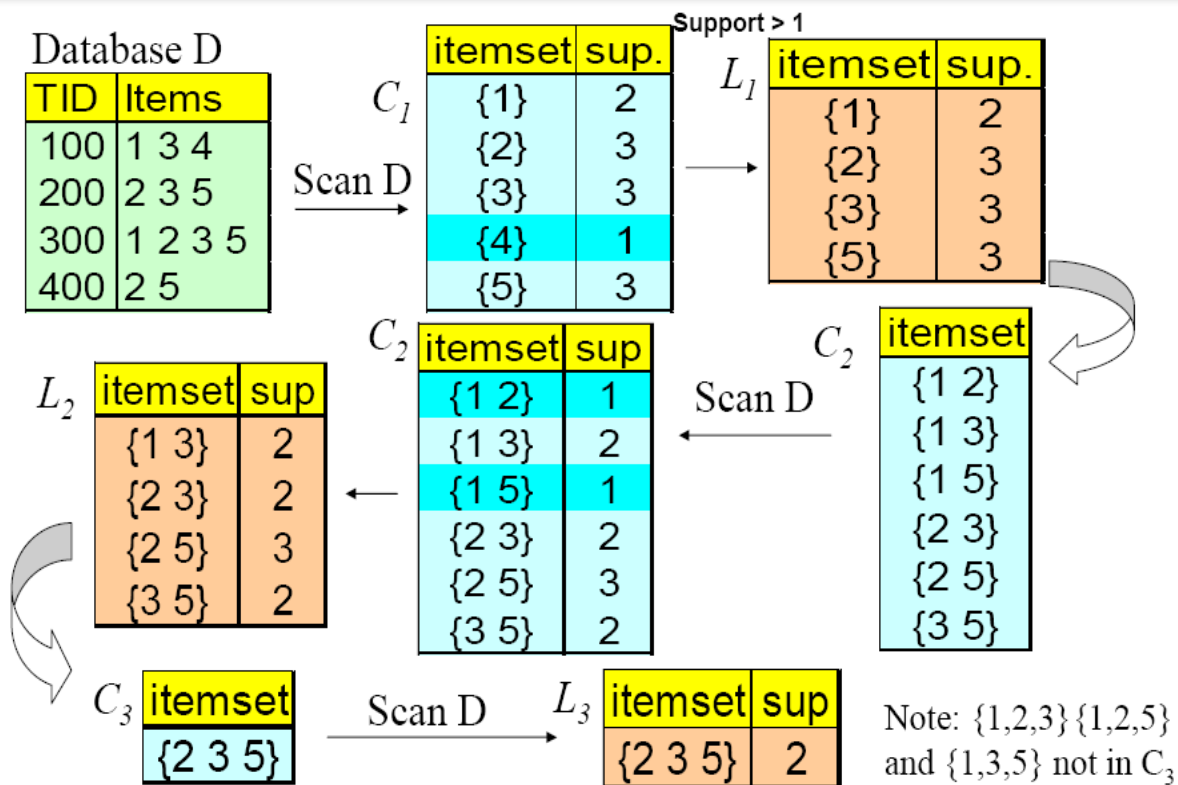


# Apriori Algorithm

❖ For each  $k$ , we construct two sets of  **$k$ -sets** (sets of size  $k$ ):

- $C_k = \text{candidate } k\text{-sets}$  = those that might be frequent sets (support  $\geq s$ ) based on information from the pass for  $k-1$ .
- $L_k$  = the set of truly frequent  $k$ -sets.

# Apriori Algorithm



| A_Rule   | Con   |
|----------|-------|
| 2->{3,5} | 66.7% |
| 3->{3,5} | 66.7% |
| 5->{2,3} | 66.7% |
| {3,5}->2 | 1     |
| {2,5}->3 | 66.7% |
| {2,3}->5 | 1     |

>50%

# Apriori Algorithm



$(T1, T2, T3, T4, T5, T6, T7, T8, T9)$

$\{\{I1, I2, I5\}, \{I2, I4\}, \{I2, I3\}, \{I1, I2, I4\}, \{I1, I3\}, \{I2, I3\}, \{I1, I3\}, \{I1, I2, I3, I5\}, \{I1, I2, I3\}\}$

# Apriori Algorithm

C1

| Items | Support |
|-------|---------|
| {I1}  | 6       |
| {I2}  | 7       |
| {I3}  | 6       |
| {I4}  | 2       |
| {I5}  | 2       |

Min\_Support = 2

# Apriori Algorithm

C2

| Items    | Support |
|----------|---------|
| {I1, I2} | 4       |
| {I1, I3} | 4       |
| {I1, I4} | 1       |
| {I1, I5} | 2       |
| {I2, I3} | 4       |
| {I2, I4} | 2       |
| {I2, I5} | 2       |
| {I3, I4} | 0       |
| {I3, I5} | 1       |
| {I4, I5} | 0       |

Min\_Support = 2

Do you finish the task?

A

Yes

B

No

提交



# Apriori Algorithm



C3

| Items        | Support |
|--------------|---------|
| {I1, I2, I3} | 2       |
| {I1, I2, I5} | 2       |
| {I1, I3, I5} | 1       |
| {I1, I2, I4} | 1       |

Min\_Support = 2

C4

| Items            | Support |
|------------------|---------|
| {I1, I2, I3, I5} | 1       |

# Apriori Algorithm

$C_k$ : Candidate itemset of size  $k$

$L_k$ : Frequent itemset of size  $k$

$L_1 \leftarrow \{frequent\ items\}$

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

$C_{k+1} \leftarrow candidate(L_k)$

candidates

    for each transaction  $t$

$Q \leftarrow \{c \mid c \in C_{k+1} \wedge c \subseteq t\}$

$count[c] \leftarrow count[c] + 1, \quad \forall c \in Q$

counting

    end for

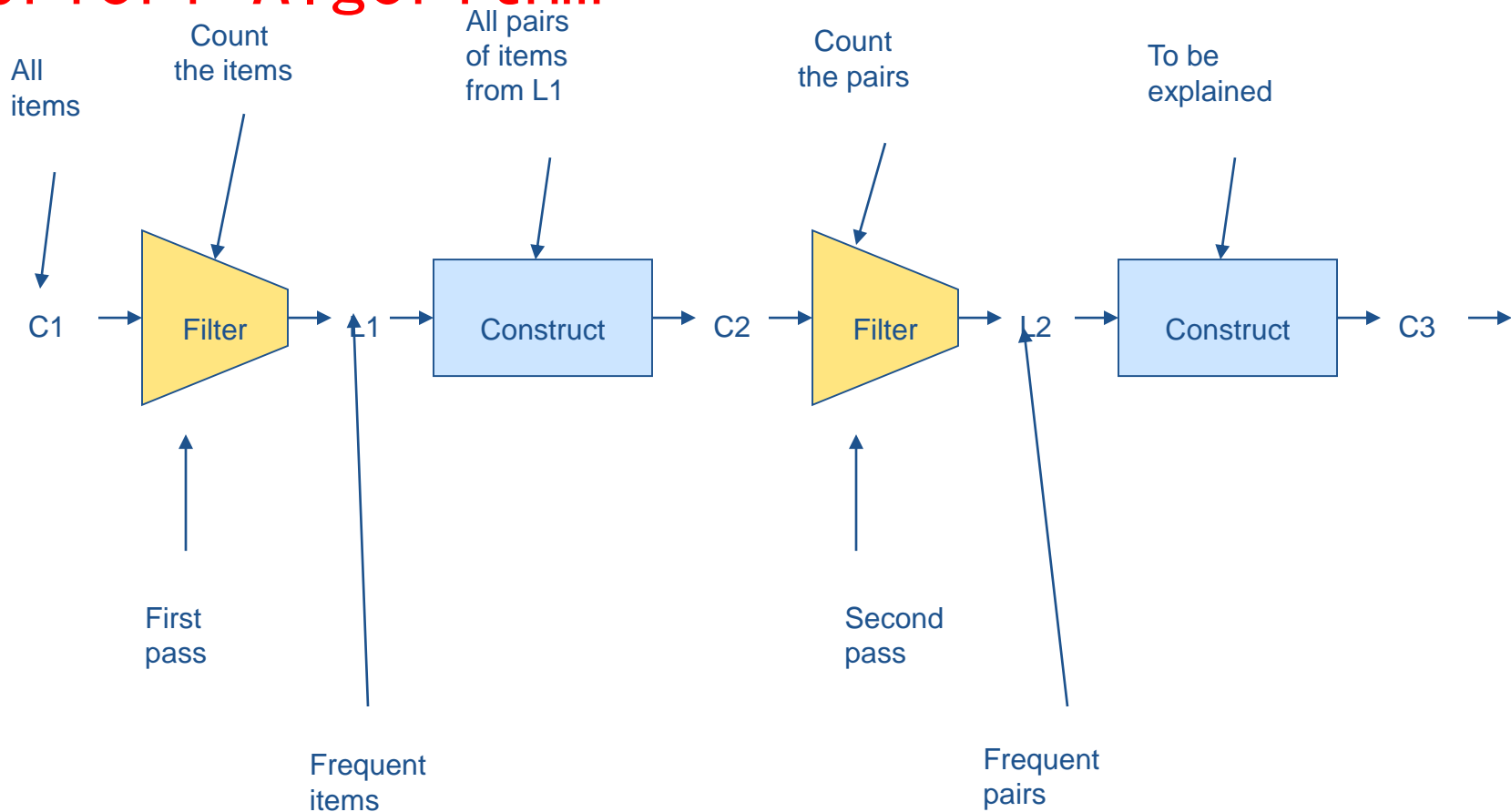
$L_{k+1} \leftarrow \{c \mid c \in C_{k+1} \wedge count[c] / N \geq \sigma\}$

filtering

end for

return  $\bigcup_k L_k$

# Apriori Algorithm



# Apriori Algorithm

```
from akapriori import Apriori
dataset =
[("ariocot","apple","cherry","plum","banana"),("strawberry","plum","cherry"),("persimmon",
"peach","banana","apple"),("kiwi
fruit","apple","pear"),("cherry","pear","banana"),("watermelon","apple"),("plum","banana",
"apple"),("pear","peach","cherry","banana","apricot"),("pineapple","apple","plum",("banana",
"plum","peach"),("grape","cherry"),("mandarin","plum"),("melon","apple","persimmon",
"plum"),("peach","cherry","apple"),("apple","mandarin","plum","persimmon")]
rules = apriori(dataset,support = 0.05,confidence = 0.3,lift = 2)
rules_sorted = sorted(rules,key=lambda x:[x[4],x[3],x[2]],reverse = True)
```

Do you finish the task?

A

Yes

B

No

提交

# Lift

$$\text{Lift}(A \Rightarrow B) > 1$$



$$\text{Lift}(A \Rightarrow B) = \frac{\text{Support}(A \Rightarrow B)}{\text{Support}(A) * \text{Support}(B)}$$

$$\text{Lift}(A \Rightarrow B) \leq 1$$



$$\text{Lift}(A \Rightarrow B) = 1$$



# Lift

| 篮子 | 商品1  | 商品2  | 商品3  |
|----|------|------|------|
| 1  | 香草威化 | 香蕉   | 狗粮   |
| 2  | 香蕉   | 面包   | 酸奶   |
| 3  | 香蕉   | 苹果   | 酸奶   |
| 4  | 香草威化 | 香蕉   | 生奶油  |
| 5  | 面包   | 香草威化 | 酸奶   |
| 6  | 牛奶   | 面包   | 香蕉   |
| 7  | 香草威化 | 苹果   | 香蕉   |
| 8  | 酸奶   | 苹果   | 香草威化 |
| 9  | 香草威化 | 香蕉   | 牛奶   |
| 10 | 香蕉   | 面包   | 花生酱  |

$$\text{Lift (Vanilla Wafer} \Rightarrow \text{Banana)} =$$

$$\text{Lift (Banana} \Rightarrow \text{Vanilla Wafer)} =$$

$$\text{Confidence (Vanilla Wafer} \Rightarrow \text{Banana)} = 4/6 = 67\%$$

$$\text{Confidence (Banana} \Rightarrow \text{Vanilla Wafer)} = 4/8 = 50\%$$

# Association Rule Example

## 最佳组合



Cloud Computing Bible



气象灾害防护指引：暴雨

## 最佳组合



YINGFA英发 OK3800AF  
近视泳镜 大镜框 舒适款

奇海平光防紫外线防雾游  
泳镜2500M黑色（镜片防

✓ ¥49.00

奥浪均码男士泳裤8320均  
码

✓ ¥59.00

侨丰电动气泵

✓ ¥29.00

奇海平光防紫外线防雾游  
泳镜2500M蓝色（镜片防

✓ ¥49.00

## Customers Who Bought This Item Also Bought



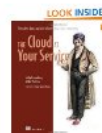
Cloud Computing  
Explained:  
Implementation  
Handbook... by John  
Rhoton

★★★★★ (17)



Cloud Computing  
Architected: Solution  
Design Handbook by John  
Rhoton

★★★★★ (3)  
\$26.37



The Cloud at Your  
Service by Jothy Rosenberg

★★★★★ (5)  
\$19.79



# Association Rule Example



## 关于抽查作业，你的想法？

A

需要不停补作业，老师很烦人

B

督促学生按时完成作业，不至于积累很多

C

没有什么意思，该抄还是抄

D

作业又不计分，没有必要检查

提交



贵在坚持！