

# Association Rule Mining with R \*

Yanchang Zhao

<http://www.RDataMining.com>

R and Data Mining Course

Beijing University of Posts and Telecommunications,  
Beijing, China

July 2019

---

\*Chapter 9 - Association Rules, in *R and Data Mining: Examples and Case Studies*. <http://www.rdatamining.com/docs/RDataMining-book.pdf>

# Contents

## Association Rules: Concept and Algorithms

- Basics of Association Rules

- Algorithms: Apriori, ECLAT and FP-growth

- Interestingness Measures

- Applications

## Association Rule Mining with R

- Mining Association Rules

- Removing Redundancy

- Interpreting Rules

- Visualizing Association Rules

- Wrap Up

## Further Readings and Online Resources

## Exercise

# Association Rules

- ▶ To discover association rules showing itemsets that occur together frequently [Agrawal et al., 1993].
- ▶ Widely used to analyze retail basket or transaction data.
- ▶ An association rule is of the form  $A \Rightarrow B$ , where  $A$  and  $B$  are itemsets or attribute-value pair sets and  $A \cap B = \emptyset$ .
- ▶ A: antecedent, left-hand-side or LHS
- ▶ B: consequent, right-hand-side or RHS
- ▶ The rule means that those database tuples having the items in the left hand of the rule are also likely to having those items in the right hand.
- ▶ Examples of association rules:
  - ▶ *bread*  $\Rightarrow$  *butter*
  - ▶ *computer*  $\Rightarrow$  *software*
  - ▶ *age in [25,35] & income in [80K,120K]*  $\Rightarrow$  *buying up-to-date mobile handsets*

# Association Rules

Association rules are rules presenting association or correlation between itemsets.

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

where  $P(A)$  is the percentage (or probability) of cases containing  $A$ .

# An Example

- ▶ Assume there are 100 students.
- ▶ 10 out of them know data mining techniques, 8 know R language and 6 know both of them.
- ▶  $R \Rightarrow DM$ : If a student knows R, then he or she knows data mining.

# An Example

- ▶ Assume there are 100 students.
- ▶ 10 out of them know data mining techniques, 8 know R language and 6 know both of them.
- ▶  $R \Rightarrow DM$ : If a student knows R, then he or she knows data mining.
- ▶ support =

## An Example

- ▶ Assume there are 100 students.
- ▶ 10 out of them know data mining techniques, 8 know R language and 6 know both of them.
- ▶  $R \Rightarrow DM$ : If a student knows R, then he or she knows data mining.
- ▶  $\text{support} = P(R \wedge DM) = 6/100 = 0.06$

## An Example

- ▶ Assume there are 100 students.
- ▶ 10 out of them know data mining techniques, 8 know R language and 6 know both of them.
- ▶  $R \Rightarrow DM$ : If a student knows R, then he or she knows data mining.
- ▶  $\text{support} = P(R \wedge DM) = 6/100 = 0.06$
- ▶  $\text{confidence} =$



## An Example

- ▶ Assume there are 100 students.
- ▶ 10 out of them know data mining techniques, 8 know R language and 6 know both of them.
- ▶  $R \Rightarrow DM$ : If a student knows R, then he or she knows data mining.
- ▶  $\text{support} = P(R \wedge DM) = 6/100 = 0.06$
- ▶  $\text{confidence} = \text{support} / P(R) = 0.06/0.08 = 0.75$

## An Example

- ▶ Assume there are 100 students.
- ▶ 10 out of them know data mining techniques, 8 know R language and 6 know both of them.
- ▶  $R \Rightarrow DM$ : If a student knows R, then he or she knows data mining.
- ▶  $\text{support} = P(R \wedge DM) = 6/100 = 0.06$
- ▶  $\text{confidence} = \text{support} / P(R) = 0.06/0.08 = 0.75$
- ▶  $\text{lift} =$

## An Example

- ▶ Assume there are 100 students.
- ▶ 10 out of them know data mining techniques, 8 know R language and 6 know both of them.
- ▶  $R \Rightarrow DM$ : If a student knows R, then he or she knows data mining.
- ▶  $\text{support} = P(R \wedge DM) = 6/100 = 0.06$
- ▶  $\text{confidence} = \text{support} / P(R) = 0.06/0.08 = 0.75$
- ▶  $\text{lift} = \text{confidence} / P(DM) = 0.75/0.1 = 7.5$

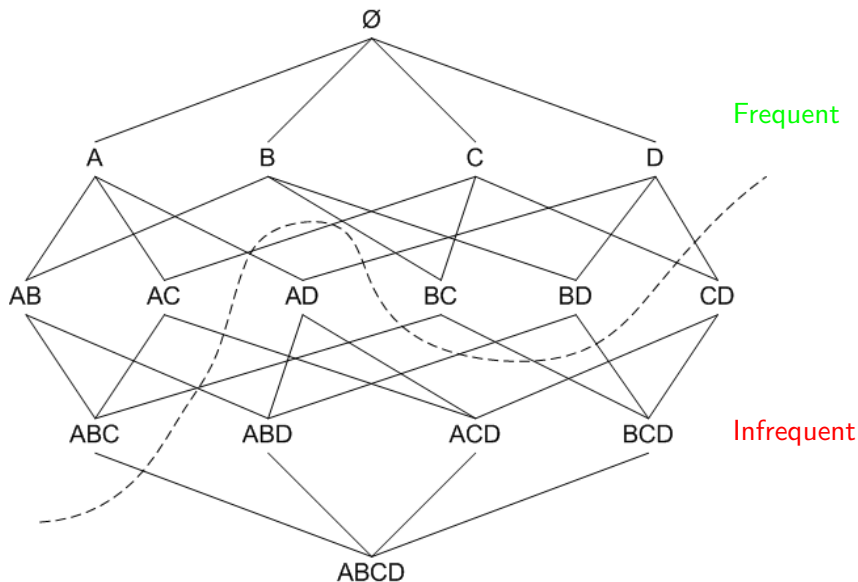
# Association Rule Mining

- ▶ Association Rule Mining is normally composed of two steps:
  - ▶ Finding all frequent itemsets whose supports are no less than a minimum support threshold;
  - ▶ From above frequent itemsets, generating association rules with confidence above a minimum confidence threshold.
- ▶ The second step is straightforward, but the first one, frequent itemset generation, is computing intensive.
- ▶ The number of possible itemsets is  $2^n - 1$ , where  $n$  is the number of unique items.
- ▶ Algorithms: Apriori, ECLAT, FP-Growth

# Downward-Closure Property

- ▶ Downward-closure property of support, a.k.a. anti-monotonicity
- ▶ For a frequent itemset, all its subsets are also frequent.  
if  $\{A,B\}$  is frequent, then both  $\{A\}$  and  $\{B\}$  are frequent.
- ▶ For an infrequent itemset, all its super-sets are infrequent.  
if  $\{A\}$  is infrequent, then  $\{A,B\}$ ,  $\{A,C\}$  and  $\{A,B,C\}$  are infrequent.
- ▶ Useful to prune candidate itemsets

# Itemset Lattice



# Apriori

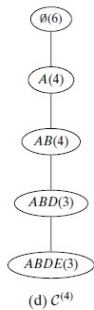
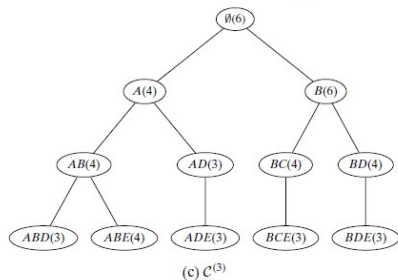
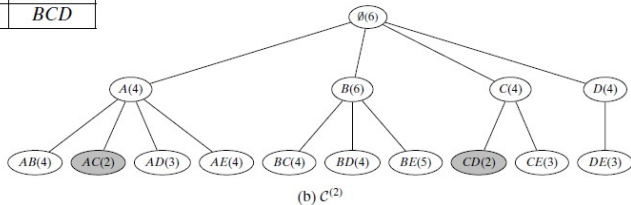
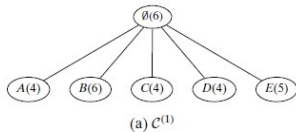
- ▶ Apriori [Agrawal and Srikant, 1994]: a classic algorithm for association rule mining
- ▶ A level-wise, breadth-first algorithm
- ▶ Counts transactions to find frequent itemsets
- ▶ Generates candidate itemsets by exploiting downward closure property of support

# Apriori Process

1. Find all frequent 1-itemsets  $L_1$
2. Join step: generate candidate  $k$ -itemsets by joining  $L_{k-1}$  with itself
3. Prune step: prune candidate  $k$ -itemsets using downward-closure property
4. Scan the dataset to count frequency of candidate  $k$ -itemsets and select frequent  $k$ -itemsets  $L_k$
5. Repeat above process, until no more frequent itemsets can be found.



$t$	$\mathbf{i}(t)$
1	<i>ABDE</i>
2	<i>BCE</i>
3	<i>ABDE</i>
4	<i>ABCE</i>
5	<i>ABCDE</i>
6	<i>BCD</i>



# FP-growth

- ▶ FP-growth: frequent-pattern growth, which mines frequent itemsets without candidate generation [Han et al., 2004]
- ▶ Compresses the input database creating an FP-tree instance to represent frequent items.
- ▶ Divides the compressed database into a set of conditional databases, each one associated with one frequent pattern.
- ▶ Each such database is mined separately.
- ▶ It reduces search costs by looking for short patterns recursively and then concatenating them in long frequent patterns.<sup>†</sup>

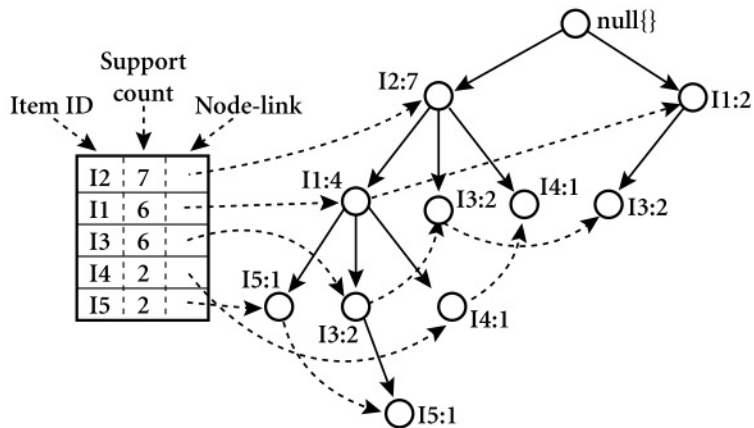
---

<sup>†</sup>[https://en.wikibooks.org/wiki/Data\\_Mining\\_Algorithms\\_In\\_R/Frequent\\_Pattern\\_Mining/The\\_FP-Growth\\_Algorithm](https://en.wikibooks.org/wiki/Data_Mining_Algorithms_In_R/Frequent_Pattern_Mining/The_FP-Growth_Algorithm)

# FP-tree

- ▶ The frequent-pattern tree (FP-tree) is a compact structure that stores quantitative information about frequent patterns in a dataset. It has two components:
  - ▶ A root labeled as “null” with a set of item-prefix subtrees as children
  - ▶ A frequent-item header table
- ▶ Each node has three attributes:
  - ▶ Item name
  - ▶ Count: number of transactions represented by the path from root to the node
  - ▶ Node link: links to the next node having the same item name
- ▶ Each entry in the frequent-item header table also has three attributes:
  - ▶ Item name
  - ▶ Head of node link: point to the first node in the FP-tree having the same item name
  - ▶ Count: frequency of the item

# FP-tree



From [Han, 2005]

# The FP-growth Algorithm

- ▶ In the first pass, the algorithm counts occurrence of items (attribute-value pairs) in the dataset, and stores them to a header table.
- ▶ In the second pass, it builds the FP-tree structure by inserting instances.
- ▶ Items in each instance have to be sorted by descending order of their frequency in the dataset, so that the tree can be processed quickly.
- ▶ Items in each instance that do not meet minimum coverage threshold are discarded.
- ▶ If many instances share most frequent items, FP-tree provides high compression close to tree root.

# The FP-growth Algorithm

- ▶ Recursive processing of this compressed version of main dataset grows large item sets directly, instead of generating candidate items and testing them against the entire database.
- ▶ Growth starts from the bottom of the header table (having longest branches), by finding all instances matching given condition.
- ▶ New tree is created, with counts projected from the original tree corresponding to the set of instances that are conditional on the attribute, with each node getting sum of its children counts.
- ▶ Recursive growth ends when no individual items conditional on the attribute meet minimum support threshold, and processing continues on the remaining header items of the original FP-tree.
- ▶ Once the recursive process has completed, all large item sets with minimum coverage have been found, and association rule creation begins.

# ECLAT

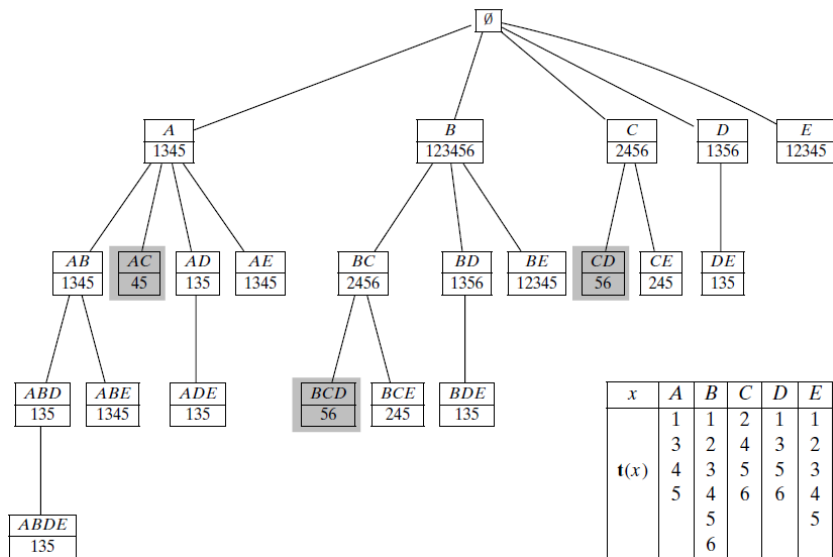
- ▶ ECLAT: equivalence class transformation [Zaki et al., 1997]
- ▶ A depth-first search algorithm using set intersection
- ▶ Idea: use tid (transaction ID) set intersection to compute the support of a candidate itemset, avoiding the generation of subsets that does not exist in the prefix tree.
- ▶  $t(AB) = t(A) \cap t(B)$ , where  $t(A)$  is the set of IDs of transactions containing A.
- ▶  $support(AB) = |t(AB)|$
- ▶ Eclat intersects the tidsets only if the frequent itemsets share a common prefix.
- ▶ It traverses the prefix search tree in a way of depth-first searching, processing a group of itemsets that have the same prefix, also called a prefix equivalence class.

# ECLAT

- ▶ It works recursively.
- ▶ The initial call uses all single items with their tid-sets.
- ▶ In each recursive call, it verifies each itemset tid-set pair  $(X, t(X))$  with all the other pairs to generate new candidates. If the new candidate is frequent, it is added to the set  $P_X$ .
- ▶ Recursively, it finds all frequent itemsets in the  $X$  branch.



# ECLAT



From [?]

# Interestingness Measures

- ▶ Which rules or patterns are interesting (and useful)?
- ▶ Two types of rule interestingness measures: subjective and objective [Freitas, 1998, Silberschatz and Tuzhilin, 1996].
- ▶ Objective measures, such as *lift*, *odds ratio* and *conviction*, are often data-driven and give the interestingness in terms of statistics or information theory.
- ▶ Subjective (user-driven) measures, such as *unexpectedness* and *actionability*, focus on finding interesting patterns by matching against a given set of user beliefs.

# Objective Interestingness Measures

- ▶ Support, confidence and lift are the most widely used objective measures to select interesting rules.
- ▶ Many other objective measures introduced by Tan et al. [Tan et al., 2002], such as  *$\phi$ -coefficient*, *odds ratio*, *kappa*, *mutual information*, *J-measure*, *Gini index*, *laplace*, *conviction*, *interest* and *cosine*.
- ▶ Different measures have different intrinsic properties and there is no measure that is better than others in all application domains.
- ▶ In addition, any-confidence, all-confidence and bond, are designed by Omiecinski [Omiecinski, 2003].
- ▶ *Utility* is used by Chan et al. [Chan et al., 2003] to find top-*k* objective-directed rules.
- ▶ *Unexpected Confidence Interestingness* and *Isolated Interestingness* are designed by Dong and Li [Dong and Li, 1998] by considering its unexpectedness in terms of other association rules in its neighbourhood.

# Subjective Interestingness Measures

- ▶ A pattern is unexpected if it is new to a user or contradicts the user's experience or domain knowledge.
- ▶ A pattern is actionable if the user can do something with it to his/her advantage [Silberschatz and Tuzhilin, 1995].
- ▶ Liu and Hsu [Liu and Hsu, 1996] proposed to rank learned rules by matching against expected patterns provided by the user.
- ▶ Ras and Wieczorkowska [Ras and Wieczorkowska, 2000] designed action-rules which show “what actions should be taken to improve the profitability of customers”. The attributes are grouped into “hard attributes” which cannot be changed and “soft attributes” which are possible to change with reasonable costs. The status of customers can be moved from one to another by changing the values of soft ones.

# Interestingness Measures - I

#	Measure	Formula
1	$\phi$ -coefficient	$\frac{P(A,B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
2	Goodman-Kruskal's ( $\lambda$ )	$\frac{\sum_j \max_k P(A_j, B_k) + \sum_k \max_j P(A_j, B_k) - \max_j P(A_j) - \max_k P(B_k)}{2 - \max_j P(A_j) - \max_k P(B_k)}$
3	Odds ratio ( $\alpha$ )	$\frac{P(A,B)P(\bar{A},\bar{B})}{P(A,\bar{B})P(\bar{A},B)}$
4	Yule's $Q$	$\frac{P(A,B)P(\bar{A}\bar{B}) - P(A,\bar{B})P(\bar{A},B)}{P(A,B)P(\bar{A}\bar{B}) + P(A,\bar{B})P(\bar{A},B)} = \frac{\alpha - 1}{\alpha + 1}$
5	Yule's $Y$	$\frac{\sqrt{P(A,B)P(\bar{A}\bar{B})} - \sqrt{P(A,\bar{B})P(\bar{A},B)}}{\sqrt{P(A,B)P(\bar{A}\bar{B})} + \sqrt{P(A,\bar{B})P(\bar{A},B)}} = \frac{\sqrt{\alpha} - 1}{\sqrt{\alpha} + 1}$
6	Kappa ( $\kappa$ )	$\frac{P(A,B) + P(\bar{A},\bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}$
7	Mutual Information ( $M$ )	$\frac{\sum_i \sum_j P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}}{\min(-\sum_i P(A_i) \log P(A_i), -\sum_j P(B_j) \log P(B_j))}$
8	J-Measure ( $J$ )	$\max \left( P(A, B) \log \left( \frac{P(B A)}{P(B)} \right) + P(\bar{A}\bar{B}) \log \left( \frac{P(\bar{B} \bar{A})}{P(\bar{B})} \right), \right.$
9	Gini index ( $G$ )	$\left. P(A, B) \log \left( \frac{P(A B)}{P(A)} \right) + P(\bar{A}B) \log \left( \frac{P(\bar{A} B)}{P(\bar{A})} \right) \right)$ $\max \left( P(A)[P(B A)^2 + P(\bar{B} A)^2] + P(\bar{A})[P(B \bar{A})^2 + P(\bar{B} \bar{A})^2] \right.$ $\left. - P(B)^2 - P(\bar{B})^2, \right.$ $\left. P(B)[P(A B)^2 + P(\bar{A} B)^2] + P(\bar{B})[P(A \bar{B})^2 + P(\bar{A} \bar{B})^2] \right.$ $\left. - P(A)^2 - P(\bar{A})^2 \right)$

From [Tan et al., 2002]

# Interestingness Measures - II

10	Support ( $s$ )	$P(A, B)$
11	Confidence ( $c$ )	$\max(P(B A), P(A B))$
12	Laplace ( $L$ )	$\max\left(\frac{NP(A,B)+1}{NP(A)+2}, \frac{NP(A,B)+1}{NP(B)+2}\right)$
13	Conviction ( $V$ )	$\max\left(\frac{P(A)P(\bar{B})}{P(AB)}, \frac{P(B)P(\bar{A})}{P(BA)}\right)$
14	Interest ( $I$ )	$\frac{P(A,B)}{P(A)P(B)}$
15	cosine ( $IS$ )	$\frac{P(A,B)}{\sqrt{P(A)P(B)}}$
16	Piatetsky-Shapiro's ( $PS$ )	$P(A, B) - P(A)P(B)$
17	Certainty factor ( $F$ )	$\max\left(\frac{P(B A)-P(B)}{1-P(B)}, \frac{P(A B)-P(A)}{1-P(A)}\right)$
18	Added Value ( $AV$ )	$\max(P(B A) - P(B), P(A B) - P(A))$
19	Collective strength ( $S$ )	$\frac{P(A,B)+P(\bar{A}\bar{B})}{P(A)P(B)+P(\bar{A})P(\bar{B})} \times \frac{1-P(A)P(B)-P(\bar{A})P(\bar{B})}{1-P(A,B)-P(\bar{A}\bar{B})}$
20	Jaccard ( $\zeta$ )	$\frac{P(A,B)}{P(A)+P(B)-P(A,B)}$
21	Klosgen ( $K$ )	$\sqrt{P(A, B) \max(P(B A) - P(B), P(A B) - P(A))}$

From [Tan et al., 2002]

# Applications

- ▶ Market basket analysis
  - ▶ Identifying associations between items in shopping baskets, i.e., which items are frequently purchased together
  - ▶ Can be used by retailers to understand customer shopping habits, do selective marketing and plan shelf space
- ▶ Churn analysis and selective marketing
  - ▶ Discovering demographic characteristics and behaviours of customers who are likely/unlikely to switch to other telcos
  - ▶ Identifying customer groups who are likely to purchase a new service or product
- ▶ Credit card risk analysis
  - ▶ Finding characteristics of customers who are likely to default on credit card or mortgage
  - ▶ Can be used by banks to reduce risks when assessing new credit card or mortgage applications

# Applications (cont.)

- ▶ Stock market analysis
  - ▶ Finding relationships between individual stocks, or between stocks and economic factors
  - ▶ Can help stock traders select interesting stocks and improve trading strategies
- ▶ Medical diagnosis
  - ▶ Identifying relationships between symptoms, test results and illness
  - ▶ Can be used for assisting doctors on illness diagnosis or even on treatment



# Contents

## Association Rules: Concept and Algorithms

- Basics of Association Rules

- Algorithms: Apriori, ECLAT and FP-growth

- Interestingness Measures

- Applications

## Association Rule Mining with R

- Mining Association Rules

- Removing Redundancy

- Interpreting Rules

- Visualizing Association Rules

- Wrap Up

## Further Readings and Online Resources

## Exercise

# Association Rule Mining Algorithms in R

- ▶ Apriori [Agrawal and Srikant, 1994]
  - ▶ A level-wise, breadth-first algorithm which counts transactions to find frequent itemsets and then derive association rules from them
  - ▶ `apriori()` in package *arules*
- ▶ ECLAT [Zaki et al., 1997]
  - ▶ Finds frequent itemsets with equivalence classes, depth-first search and set intersection instead of counting
  - ▶ `ec1at()` in package *arules*

# The Titanic Dataset

- ▶ The Titanic dataset in the *datasets* package is a 4-dimensional table with summarized information on the fate of passengers on the Titanic according to social class, sex, age and survival.
- ▶ To make it suitable for association rule mining, we reconstruct the raw data as `titanic.raw`, where each row represents a person.
- ▶ The reconstructed raw data can also be downloaded at <http://www.rdatamining.com/data/titanic.raw.rdata>.

# Pipe Operations in R

- ▶ Load library magrittr for pipe operations
- ▶ Avoid nested function calls
- ▶ Make code easy to understand
- ▶ Supported by dplyr and ggplot2

```
library(magrittr)  ## for pipe operations
## traditional way
b <- fun3(fun2(fun1(a), p2))
## the above can be rewritten to
b <- a %>% fun1() %>% fun2(p2) %>% fun3()
```

```
## download data
download.file(url="http://www.rdatamining.com/data/titanic.raw.rdata",
              destfile="./data/titanic.raw.rdata")
```

```
library(magrittr)  ## for pipe operations
## load data, and the name of the R object is titanic.raw
load("../data/titanic.raw.rdata")
## dimensionality
titanic.raw %>% dim()
## [1] 2201      4

## structure of data
titanic.raw %>% str()
## 'data.frame': 2201 obs. of  4 variables:
## $ Class      : Factor w/ 4 levels "1st","2nd","3rd",...: 3 3 3...
## $ Sex        : Factor w/ 2 levels "Female","Male": 2 2 2 2 2 ...
## $ Age        : Factor w/ 2 levels "Adult","Child": 2 2 2 2 2 ...
## $ Survived   : Factor w/ 2 levels "No","Yes": 1 1 1 1 1 1 1...
```

```
## draw a random sample of 5 records
idx <- 1:nrow(titanic.raw) %>% sample(5)
titanic.raw[idx, ]
```

```
##      Class      Sex   Age Survived
## 2080    2nd Female Adult      Yes
## 1162   Crew    Male Adult      No
## 954    Crew    Male Adult      No
## 2172   3rd Female Adult      Yes
## 456    3rd    Male Adult      No
```

```
## a summary of the dataset
titanic.raw %>% summary()
```

```
##      Class      Sex      Age      Survived
## 1st :325   Female: 470   Adult:2092   No :1490
## 2nd :285   Male  :1731   Child: 109   Yes: 711
## 3rd :706
## Crew:885
```

# Function apriori()

- ▶ Mine frequent itemsets, association rules or association hyperedges using the Apriori algorithm.
- ▶ The Apriori algorithm employs level-wise search for frequent itemsets.
- ▶ Default settings:
  - ▶ minimum support: `supp=0.1`
  - ▶ minimum confidence: `conf=0.8`
  - ▶ maximum length of rules: `maxlen=10`

```

## mine association rules
library(arules) ## load required library
rules.all <- titanic.raw %>% apriori() ## run the APRIORI algorithm
## Apriori
##
## Parameter specification:
## confidence minval smax arem aval originalSupport maxtime
##           0.8     0.1     1 none FALSE             TRUE       5
## support minlen maxlen target   ext
##           0.1       1      10 rules FALSE
##
## Algorithmic control:
## filter tree heap memopt load sort verbose
##       0.1 TRUE TRUE  FALSE TRUE     2     TRUE
##
## Absolute minimum support count: 220
##
## set item appearances ...[0 item(s)] done [0.00s].
## set transactions ...[10 item(s), 2201 transaction(s)] done ...
## sorting and recoding items ... [9 item(s)] done [0.00s].
## creating transaction tree ... done [0.00s].
## checking subsets of size 1 2 3 4 done [0.00s].
## writing ... [27 rule(s)] done [0.00s].
## creating S4 object ... done [0.00s].

```



```
rules.all %>% length() ## number of rules discovered
## [1] 27
```

```
rules.all %>% inspect() ## print all rules
```

##	lhs	rhs	support	confidence	...
## [1]	{}	=> {Age=Adult}	0.9504771	0.9504771	1...
## [2]	{Class=2nd}	=> {Age=Adult}	0.1185825	0.9157895	0...
## [3]	{Class=1st}	=> {Age=Adult}	0.1449341	0.9815385	1...
## [4]	{Sex=Female}	=> {Age=Adult}	0.1930940	0.9042553	0...
## [5]	{Class=3rd}	=> {Age=Adult}	0.2848705	0.8881020	0...
## [6]	{Survived=Yes}	=> {Age=Adult}	0.2971377	0.9198312	0...
## [7]	{Class=Crew}	=> {Sex=Male}	0.3916402	0.9740113	1...
## [8]	{Class=Crew}	=> {Age=Adult}	0.4020900	1.0000000	1...
## [9]	{Survived=No}	=> {Sex=Male}	0.6197183	0.9154362	1...
## [10]	{Survived=No}	=> {Age=Adult}	0.6533394	0.9651007	1...
## [11]	{Sex=Male}	=> {Age=Adult}	0.7573830	0.9630272	1...
## [12]	{Sex=Female,				...
##	Survived=Yes}	=> {Age=Adult}	0.1435711	0.9186047	0...
## [13]	{Class=3rd,				...
##	Sex=Male}	=> {Survived=No}	0.1917310	0.8274510	1...
## [14]	{Class=3rd,				...
##	Survived=No}	=> {Age=Adult}	0.2162653	0.9015152	0...
## [15]	{Class=3rd,				...
##	Sex=Male}	=> {Age=Adult}	0.2000046	0.9058824	0...

- ▶ Suppose we want to find patterns of survival and non-survival
- ▶ `verbose=F`: suppress progress report
- ▶ `minlen=2`: find rules that contain at least two items
- ▶ Use lower thresholds for support and confidence
- ▶ `rhs=c(...)`: find rules whose right-hand sides are in the list
- ▶ `default="lhs"`: use default setting for left-hand side
- ▶ `quality(...)`: interestingness measures

```
## run APRIORI again to find rules with rhs containing "Survived" only
rules.surv <- titanic.raw %>% apriori(
  control = list(verbose=F),
  parameter = list(minlen=2, supp=0.005, conf=0.8),
  appearance = list(rhs=c("Survived=No",
                           "Survived=Yes"),
                    default="lhs"))

## keep three decimal places
quality(rules.surv) <- rules.surv %>% quality() %>% round(digits=3)
## sort rules by lift
rules.surv.sorted <- rules.surv %>% sort(by="lift")
```

```

rules.surv.sorted %>% inspect() ## print rules
##      lhs                      rhs          support confidence lif...
## [1]  {Class=2nd,                      ...
##      Age=Child} => {Survived=Yes}    0.011      1.000 3.09...
## [2]  {Class=2nd,                      ...
##      Sex=Female,                      ...
##      Age=Child} => {Survived=Yes}    0.006      1.000 3.09...
## [3]  {Class=1st,                      ...
##      Sex=Female} => {Survived=Yes}    0.064      0.972 3.01...
## [4]  {Class=1st,                      ...
##      Sex=Female,                      ...
##      Age=Adult} => {Survived=Yes}    0.064      0.972 3.01...
## [5]  {Class=2nd,                      ...
##      Sex=Female} => {Survived=Yes}    0.042      0.877 2.71...
## [6]  {Class=Crew,                      ...
##      Sex=Female} => {Survived=Yes}    0.009      0.870 2.69...
## [7]  {Class=Crew,                      ...
##      Sex=Female,                      ...
##      Age=Adult} => {Survived=Yes}    0.009      0.870 2.69...
## [8]  {Class=2nd,                      ...
##      Sex=Female,                      ...
##      Age=Adult} => {Survived=Yes}    0.036      0.860 2.66...
## [9]  {Class=2nd,                      ...
##      Sex=Male,                      ...

```

# Redundant Rules

- ▶ There are often too many association rules discovered from a dataset.
- ▶ It is necessary to remove redundant rules before a user is able to study the rules and identify interesting ones from them.

# Redundant Rules

```
## redundant rules
rules.surv.sorted[1:2] %>% inspect()
##      lhs                                rhs      support confidence lift...
## [1] {Class=2nd,                               ...
##      Age=Child} => {Survived=Yes}    0.011      1 3.096...
## [2] {Class=2nd,                               ...
##      Sex=Female,                               ...
##      Age=Child} => {Survived=Yes}    0.006      1 3.096...
```

- ▶ Rule #2 provides no extra knowledge in addition to rule #1, since rule #1 tells us that all 2nd-class children survived.
- ▶ When a rule (such as #2) is a super rule of another rule (#1) and the former has the same or a lower lift, the former rule (#2) is considered to be redundant.
- ▶ Other redundant rules in the above result are rules #4, #7 and #8, compared respectively with #3, #6 and #5.

# Remove Redundant Rules

```
## find redundant rules
subset.matrix <- is.subset(rules.surv.sorted, rules.surv.sorted)
subset.matrix[lower.tri(subset.matrix, diag = T)] <- F
redundant <- colSums(subset.matrix) >= 1
```

```
## which rules are redundant
redundant %>% which()

## {Class=2nd,Sex=Female,Age=Child,Survived=Yes}
## 2
## {Class=1st,Sex=Female,Age=Adult,Survived=Yes}
## 4
## {Class=Crew,Sex=Female,Age=Adult,Survived=Yes}
## 7
## {Class=2nd,Sex=Female,Age=Adult,Survived=Yes}
## 8
```

```
## remove redundant rules
rules.surv.pruned <- rules.surv.sorted[!redundant]
```

## Remaining Rules

```
rules.surv.pruned %>% inspect() ## print rules
```

##	lhs	rhs	support	confidence	lift...
## [1]	{Class=2nd,				...
##	Age=Child}	=> {Survived=Yes}	0.011	1.000	3.096...
## [2]	{Class=1st,				...
##	Sex=Female}	=> {Survived=Yes}	0.064	0.972	3.010...
## [3]	{Class=2nd,				...
##	Sex=Female}	=> {Survived=Yes}	0.042	0.877	2.716...
## [4]	{Class=Crew,				...
##	Sex=Female}	=> {Survived=Yes}	0.009	0.870	2.692...
## [5]	{Class=2nd,				...
##	Sex=Male,				...
##	Age=Adult}	=> {Survived=No}	0.070	0.917	1.354...
## [6]	{Class=2nd,				...
##	Sex=Male}	=> {Survived=No}	0.070	0.860	1.271...
## [7]	{Class=3rd,				...
##	Sex=Male,				...
##	Age=Adult}	=> {Survived=No}	0.176	0.838	1.237...
## [8]	{Class=3rd,				...
##	Sex=Male}	=> {Survived=No}	0.192	0.827	1.222...

```
rules.surv.pruned[1] %>% inspect() ## print rules
##      lhs                                rhs      support confidence
## [1] {Class=2nd, Age=Child} => {Survived=Yes} 0.011      1
##      lift  count
## [1] 3.096 24
```

- ▶ Did children have a higher survival rate than adults?
- ▶ Did children of the 2nd class have a higher survival rate than other children?



```
rules.surv.pruned[1] %>% inspect() ## print rules
##      lhs                                rhs      support confidence
## [1] {Class=2nd, Age=Child} => {Survived=Yes} 0.011      1
##      lift  count
## [1] 3.096 24
```

- ▶ Did children have a higher survival rate than adults?
- ▶ Did children of the 2nd class have a higher survival rate than other children?
- ▶ The rule states only that all children of class 2 survived, but provides no information at all about the survival rates of other classes.

# Find Rules about Age Groups

- ▶ Use lower thresholds to find all rules for children of different classes
- ▶ `verbose=F`: suppress progress report
- ▶ `minlen=3`: find rules that contain at least three items
- ▶ Use lower thresholds for support and confidence
- ▶ `rhs=c(...)`, `rhs=c(...)`: find rules whose left/right-hand sides are in the list
- ▶ `quality(...)`: interestingness measures

```
## mine rules about class and age group
rules.age <- titanic.raw %>% apriori(control = list(verbose=F),
  parameter = list(minlen=3, supp=0.002, conf=0.2),
  appearance = list(default="none", rhs=c("Survived=Yes"),
    lhs=c("Class=1st", "Class=2nd", "Class=3rd",
      "Age=Child", "Age=Adult")))
rules.age <- sort(rules.age, by="confidence")
```

# Rules about Age Groups

```
rules.age %>% inspect()  ## print rules
```

##	lhs	rhs	support
## [1]	{Class=2nd, Age=Child}	=> {Survived=Yes}	0.010904134
## [2]	{Class=1st, Age=Child}	=> {Survived=Yes}	0.002726034
## [3]	{Class=1st, Age=Adult}	=> {Survived=Yes}	0.089504771
## [4]	{Class=2nd, Age=Adult}	=> {Survived=Yes}	0.042707860
## [5]	{Class=3rd, Age=Child}	=> {Survived=Yes}	0.012267151
## [6]	{Class=3rd, Age=Adult}	=> {Survived=Yes}	0.068605179

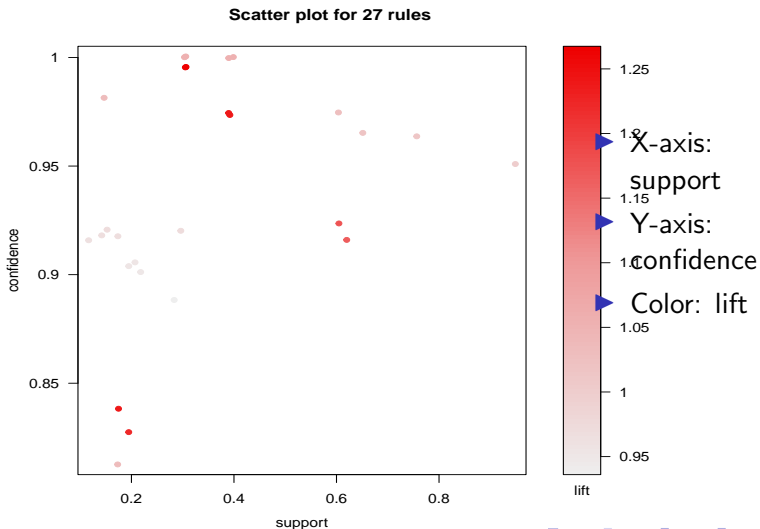
##	confidence	lift	count
## [1]	1.0000000	3.0956399	24
## [2]	1.0000000	3.0956399	6
## [3]	0.6175549	1.9117275	197
## [4]	0.3601533	1.1149048	94
## [5]	0.3417722	1.0580035	27
## [6]	0.2408293	0.7455209	151

```
## average survival rate
```

titanic.raw\$Survived	%>% table()	%>% prop.table()
## .		
##	No	Yes
##	0.676965	0.323035

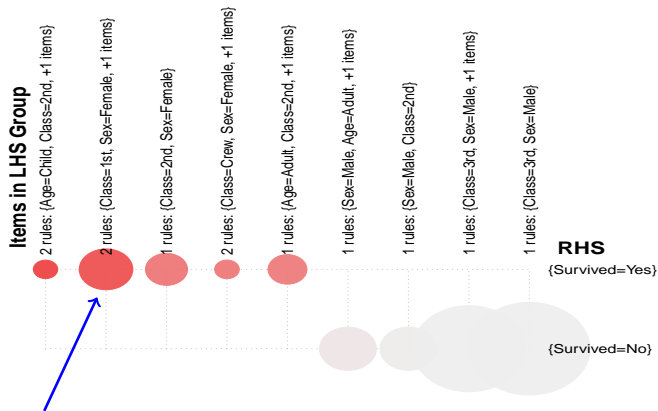
```
## rule visualisation  
library(arulesViz)  
rules.all %>% plot()
```



```
rules.surv %>% plot(method = "grouped")
```

### Grouped Matrix for 12 Rules

Size: support  
Color: lift



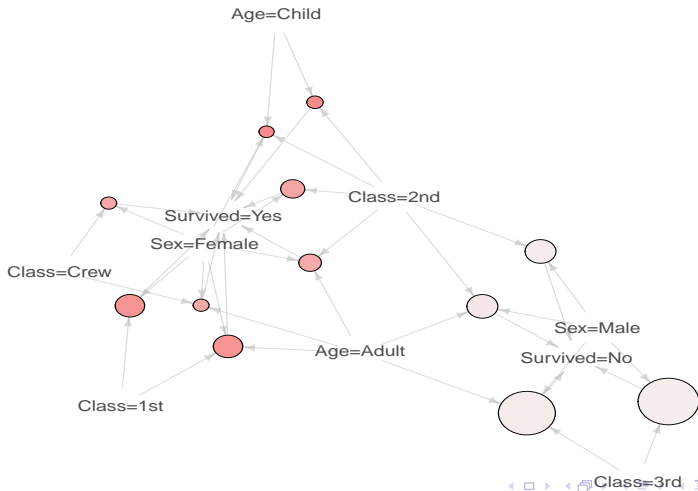
{Class=1st, Sex=Female, +1 items}  $\Rightarrow$  {Survived=Yes}

```
rules.surv %>% plot(method="graph",  
  control=list(layout=igraph::with_fr()))
```

### Graph for 12 rules

size: support (0.006 – 0.192)

color: lift (1.222 – 3.096)

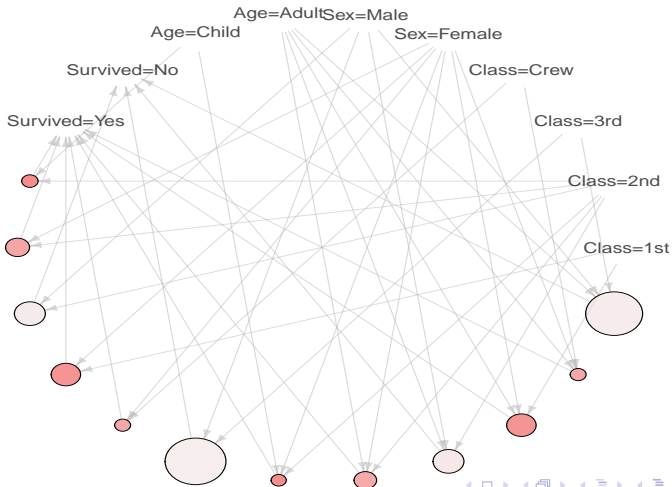


```
rules.surv %>% plot(method="graph",  
  control=list(layout=igraph::in_circle()))
```

### Graph for 12 rules

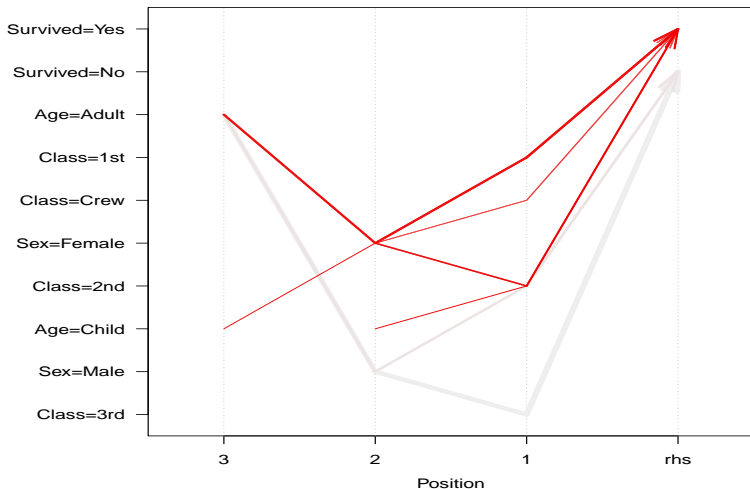
size: support (0.006 – 0.192)

color: lift (1.222 – 3.096)



```
rules.surv %>% plot(method="paracoord",  
                    control=list(reorder=T))
```

**Parallel coordinates plot for 12 rules**





# Interactive Plots and Reorder rules

```
rules.all %>% plot(interactive = T)
```

```
interactive = TRUE
```

- ▶ Selecting and inspecting one or multiple rules
- ▶ Zooming
- ▶ Filtering rules with an interesting measure

```
rules.surv %>% plot(method = "paracoord", control = list(reorder = T))
```

```
reorder = TRUE
```

- ▶ To improve visualisation by reordering rules and minimizing crossovers
- ▶ The visualisation is likely to change from run to run.

# Wrap Up

- ▶ Starting with a high support, to get a small set of rules quickly
- ▶ Setting constraints to left and/or right hand side of rules, to focus on rules that you are interested in
- ▶ Digging down data to find more associations with lower thresholds of support and confidence
- ▶ Rules of low confidence / lift can be interesting and useful.
- ▶ Be cautious when interpreting rules

# Contents

## Association Rules: Concept and Algorithms

- Basics of Association Rules

- Algorithms: Apriori, ECLAT and FP-growth

- Interestingness Measures

- Applications

## Association Rule Mining with R

- Mining Association Rules

- Removing Redundancy

- Interpreting Rules

- Visualizing Association Rules

- Wrap Up

## Further Readings and Online Resources

## Exercise

## Further Readings

- ▶ Association Rule Learning

[https://en.wikipedia.org/wiki/Association\\_rule\\_learning](https://en.wikipedia.org/wiki/Association_rule_learning)

- ▶ Data Mining Algorithms In R: Apriori

[https://en.wikibooks.org/wiki/Data\\_Mining\\_Algorithms\\_In\\_R/Frequent\\_Pattern\\_Mining/The\\_Apriori\\_Algorithm](https://en.wikibooks.org/wiki/Data_Mining_Algorithms_In_R/Frequent_Pattern_Mining/The_Apriori_Algorithm)

- ▶ Data Mining Algorithms In R: ECLAT

[https://en.wikibooks.org/wiki/Data\\_Mining\\_Algorithms\\_In\\_R/Frequent\\_Pattern\\_Mining/The\\_Eclat\\_Algorithm](https://en.wikibooks.org/wiki/Data_Mining_Algorithms_In_R/Frequent_Pattern_Mining/The_Eclat_Algorithm)

- ▶ Data Mining Algorithms In R: FP-Growth

[https://en.wikibooks.org/wiki/Data\\_Mining\\_Algorithms\\_In\\_R/Frequent\\_Pattern\\_Mining/The\\_FP-Growth\\_Algorithm](https://en.wikibooks.org/wiki/Data_Mining_Algorithms_In_R/Frequent_Pattern_Mining/The_FP-Growth_Algorithm)

- ▶ FP-Growth Implementation by Christian Borgelt

<http://www.borgelt.net/fpgrowth.html>

- ▶ Frequent Itemset Mining Implementations Repository

<http://fimi.ua.ac.be/data/>

## Further Readings

- ▶ More than 20 interestingness measures, such as chi-square, conviction, gini and leverage  
Tan, P.-N., Kumar, V., and Srivastava, J. (2002). Selecting the right interestingness measure for association patterns. In Proc. of KDD '02, pages 32-41, New York, NY, USA. ACM Press.
- ▶ More reviews on interestingness measures:  
[Silberschatz and Tuzhilin, 1996], [Tan et al., 2002] and [Omiecinski, 2003]
- ▶ Post mining of association rules, such as selecting interesting association rules, visualization of association rules and using association rules for classification [Zhao et al., 2009]  
Yanchang Zhao, et al. (Eds.). "Post-Mining of Association Rules: Techniques for Effective Knowledge Extraction", ISBN 978-1-60566-404-0, May 2009. Information Science Reference.
- ▶ Package *arulesSequences*: mining sequential patterns  
<http://cran.r-project.org/web/packages/arulesSequences/>

# Contents

## Association Rules: Concept and Algorithms

- Basics of Association Rules

- Algorithms: Apriori, ECLAT and FP-growth

- Interestingness Measures

- Applications

## Association Rule Mining with R

- Mining Association Rules

- Removing Redundancy

- Interpreting Rules

- Visualizing Association Rules

- Wrap Up

## Further Readings and Online Resources

## Exercise

# The Mushroom Dataset I

- ▶ The mushroom dataset includes descriptions of hypothetical samples corresponding to 23 species of gilled mushrooms ‡.
- ▶ A csv file with 8,124 observations on 23 categorical variables:
  1. class: edible=e, poisonous=p
  2. cap-shape: bell=b,conical=c,convex=x,flat=f, knobbed=k,sunken=s
  3. cap-surface: fibrous=f,grooves=g,scaly=y,smooth=s
  4. cap-color: brown=n,buff=b,cinnamon=c,gray=g,green=r, pink=p,purple=u,red=e,white=w,yellow=y
  5. bruises?: bruises=t,no=f
  6. odor: almond=a,anise=l,creosote=c,fishy=y,foul=f, musty=m,none=n,pungent=p,spicy=s
  7. gill-attachment: attached=a,descending=d,free=f,notched=n
  8. gill-spacing: close=c,crowded=w,distant=d
  9. gill-size: broad=b,narrow=n
  10. gill-color: black=k,brown=n,buff=b,chocolate=h,gray=g, green=r,orange=o,pink=p,purple=u,red=e, white=w,yellow=y

# The Mushroom Dataset II

- 11. stalk-shape: enlarging=e,tapering=t
- 12. stalk-root: bulbous=b,club=c,cup=u,equal=e,  
rhizomorphs=z,rooted=r,missing=?
- 13. stalk-surface-above-ring: fibrous=f,scaly=y,silky=k,smooth=s
- 14. stalk-surface-below-ring: fibrous=f,scaly=y,silky=k,smooth=s
- 15. stalk-color-above-ring:  
brown=n,buff=b,cinnamon=c,gray=g,orange=o,  
pink=p,red=e,white=w,yellow=y
- 16. stalk-color-below-ring:  
brown=n,buff=b,cinnamon=c,gray=g,orange=o,  
pink=p,red=e,white=w,yellow=y
- 17. veil-type: partial=p,universal=u
- 18. veil-color: brown=n,orange=o,white=w,yellow=y
- 19. ring-number: none=n,one=o,two=t
- 20. ring-type: cobwebby=c,evanescent=e,flaring=f,large=l,  
none=n,pendant=p,sheathing=s,zone=z



# The Mushroom Dataset III

- 21. spore-print-color:  
black=k,brown=n,buff=b,chocolate=h,green=r,  
orange=o,purple=u,white=w,yellow=y
- 22. population: abundant=a,clustered=c,numerous=n,  
scattered=s,several=v,solitary=y
- 23. habitat: grasses=g,leaves=l,meadows=m,paths=p,  
urban=u,waste=w,woods=d

---

<sup>‡</sup><https://archive.ics.uci.edu/ml/datasets/Mushroom>

# Load Mushroom Dataset

```
## load mushroom data from UCI the Machine Learning Repository
url <- paste0("http://archive.ics.uci.edu/ml/",
  "machine-learning-databases/mushroom/agaricus-lepiota.data")
```

```
mushrooms <- read.csv(file = url, header = FALSE)
names(mushrooms) <- c("class", "cap-shape", "cap-surface",
  "cap-color", "bruises", "odor", "gill-attachment", "gill-spacing",
  "gill-size", "gill-color", "stalk-shape", "stalk-root",
  "stalk-surface-above-ring", "stalk-surface-below-ring",
  "stalk-color-above-ring", "stalk-color-below-ring",
  "veil-type", "veil-color", "ring-number", "ring-type",
  "spore-print-color", "population", "habitat")
table(mushrooms$class, useNA="ifany")
##
##      e      p
## 4208 3916
```

# The Mushroom Dataset

```
str(mushrooms)

## 'data.frame': 8124 obs. of  23 variables:
## $ class                : Factor w/ 2 levels "e","p": 2 ...
## $ cap-shape             : Factor w/ 6 levels "b","c","f"...
## $ cap-surface           : Factor w/ 4 levels "f","g","s"...
## $ cap-color             : Factor w/ 10 levels "b","c","e...
## $ bruises               : Factor w/ 2 levels "f","t": 2 ...
## $ odor                  : Factor w/ 9 levels "a","c","f"...
## $ gill-attachment       : Factor w/ 2 levels "a","f": 2 ...
## $ gill-spacing           : Factor w/ 2 levels "c","w": 1 ...
## $ gill-size             : Factor w/ 2 levels "b","n": 2 ...
## $ gill-color            : Factor w/ 12 levels "b","e","g...
## $ stalk-shape           : Factor w/ 2 levels "e","t": 1 ...
## $ stalk-root            : Factor w/ 5 levels "?","b","c"...
## $ stalk-surface-above-ring: Factor w/ 4 levels "f","k","s"...
## $ stalk-surface-below-ring: Factor w/ 4 levels "f","k","s"...
## $ stalk-color-above-ring : Factor w/ 9 levels "b","c","e"...
## $ stalk-color-below-ring : Factor w/ 9 levels "b","c","e"...
## $ veil-type             : Factor w/ 1 level "p": 1 1 1 1...
## $ veil-color            : Factor w/ 4 levels "n","o","w"...
## $ ring-number           : Factor w/ 3 levels "n","o","t"...
## $ ring-type             : Factor w/ 5 levels "e" "f" "l"
```

# Exercise

- ▶ From the mushroom data, find association rules that can be used to identify the edibility of a mushroom
- ▶ Think about parameters: length of rules, minimum support, minimum confidence
- ▶ How to find only rules relevant to edibility?
- ▶ Which interestingness measures to use?
- ▶ Any redundant rules? How to remove them?
- ▶ What are characteristics of edible mushrooms? And characteristics of poisonous ones?

# Mining Association Rules from Mushroom Dataset

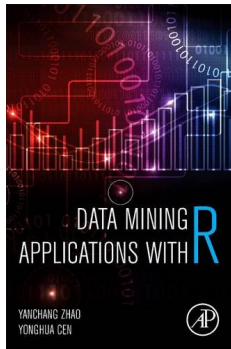
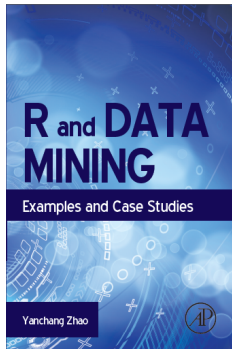
```
## find associatin rules from the mushroom dataset
rules <- apriori(mushrooms, control = list(verbose=F),
                 parameter = list(minlen=2, maxlen=5),
                 appearance = list(rhs=c("class=p", "class=e"),
                                   default="lhs"))
quality(rules) <- round(quality(rules), digits=3)
rules.sorted <- sort(rules, by="confidence")
inspect(head(rules.sorted))
```

##	lhs	rhs	support	confidence
## [1]	{ring-type=l}	=> {class=p}	0.160	1
## [2]	{gill-color=b}	=> {class=p}	0.213	1
## [3]	{odor=f}	=> {class=p}	0.266	1
## [4]	{gill-size=b,gill-color=n}	=> {class=e}	0.108	1
## [5]	{odor=n,stalk-root=e}	=> {class=e}	0.106	1
## [6]	{bruises=f,stalk-root=e}	=> {class=e}	0.106	1
##	lift	count		
## [1]	2.075	1296		
## [2]	2.075	1728		
## [3]	2.075	2160		
## [4]	1.931	880		
## [5]	1.931	864		
## [6]	1.931	864		

# Online Resources

- ▶ Book titled *R and Data Mining: Examples and Case Studies* [Zhao, 2012]  
<http://www.rdatamining.com/docs/RDataMining-book.pdf>
- ▶ R Reference Card for Data Mining  
<http://www.rdatamining.com/docs/RDataMining-reference-card.pdf>
- ▶ Free online courses and documents  
<http://www.rdatamining.com/resources/>
- ▶ RDataMining Group on LinkedIn (27,000+ members)  
<http://group.rdatamining.com>
- ▶ Twitter (3,300+ followers)  
@RDataMining

# The End



Thanks!

Email: [yanchang\(at\)RDataMining.com](mailto:yanchang(at)RDataMining.com)

Twitter: @RDataMining

# How to Cite This Work

## ► Citation

Yanchang Zhao. R and Data Mining: Examples and Case Studies. ISBN 978-0-12-396963-7, December 2012. Academic Press, Elsevier. 256 pages. URL: <http://www.rdatamining.com/docs/RDataMining-book.pdf>.

## ► BibTex

```
@BOOK{Zhao2012R,  
  title = {R and Data Mining: Examples and Case Studies},  
  publisher = {Academic Press, Elsevier},  
  year = {2012},  
  author = {Yanchang Zhao},  
  pages = {256},  
  month = {December},  
  isbn = {978-0-123-96963-7},  
  keywords = {R, data mining},  
  url = {http://www.rdatamining.com/docs/RDataMining-book.pdf}  
}
```



# References I



Agrawal, R., Imielinski, T., and Swami, A. (1993).

Mining association rules between sets of items in large databases.

In *Proc. of the ACM SIGMOD International Conference on Management of Data*, pages 207–216, Washington D.C. USA.



Agrawal, R. and Srikant, R. (1994).

Fast algorithms for mining association rules in large databases.

In *Proc. of the 20th International Conference on Very Large Data Bases*, pages 487–499, Santiago, Chile.



Chan, R., Yang, Q., and Shen, Y.-D. (2003).

Mining high utility itemsets.

In *Data Mining, 2003. ICDM 2003. Third IEEE International Conference on*, pages 19–26.



Dong, G. and Li, J. (1998).

Interestingness of discovered association rules in terms of neighborhood-based unexpectedness.

In *PAKDD '98: Proceedings of the Second Pacific-Asia Conference on Research and Development in Knowledge Discovery and Data Mining*, pages 72–86, London, UK. Springer-Verlag.



Freitas, A. A. (1998).

On objective measures of rule surprisingness.

In *PKDD '98: Proceedings of the Second European Symposium on Principles of Data Mining and Knowledge Discovery*, pages 1–9, London, UK. Springer-Verlag.



Han, J. (2005).

*Data Mining: Concepts and Techniques*.

Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.



Han, J., Pei, J., Yin, Y., and Mao, R. (2004).

Mining frequent patterns without candidate generation.

*Data Mining and Knowledge Discovery*, 8:53–87.

# References II



Liu, B. and Hsu, W. (1996).

Post-analysis of learned rules.

In *Proceedings of the 13th National Conference on Artificial Intelligence (AAAI-96)*, pages 828–834, Portland, Oregon, USA.



Omiecinski, E. R. (2003).

Alternative interest measures for mining associations in databases.

*IEEE Transactions on Knowledge and Data Engineering*, 15(1):57–69.



Ras, Z. W. and Wieczorkowska, A. (2000).

Action-rules: How to increase profit of a company.

In *PKDD '00: Proceedings of the 4th European Conference on Principles of Data Mining and Knowledge Discovery*, pages 587–592, London, UK. Springer-Verlag.



Silberschatz, A. and Tuzhilin, A. (1995).

On subjective measures of interestingness in knowledge discovery.

In *Knowledge Discovery and Data Mining*, pages 275–281.



Silberschatz, A. and Tuzhilin, A. (1996).

What makes patterns interesting in knowledge discovery systems.

*IEEE Transactions on Knowledge and Data Engineering*, 8(6):970–974.



Tan, P.-N., Kumar, V., and Srivastava, J. (2002).

Selecting the right interestingness measure for association patterns.

In *KDD '02: Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 32–41, New York, NY, USA. ACM Press.



Zaki, M. J., Parthasarathy, S., Ogihara, M., and Li, W. (1997).

New algorithms for fast discovery of association rules.

Technical Report 651, Computer Science Department, University of Rochester, Rochester, NY 14627.

# References III



Zhao, Y. (2012).

*R and Data Mining: Examples and Case Studies*, ISBN 978-0-12-396963-7.  
Academic Press, Elsevier.



Zhao, Y., Zhang, C., and Cao, L., editors (2009).

*Post-Mining of Association Rules: Techniques for Effective Knowledge Extraction*, ISBN 978-1-60566-404-0.

Information Science Reference, Hershey, PA.