

Association Rules

This lab presents examples of association rule mining with R. It starts with basic concepts of association rules, and then demonstrates association rules mining with R. After that, it presents examples of pruning redundant rules and interpreting and visualizing association rules.

Basics of Association Rules

Association rules are rules presenting association or correlation between itemsets. An association rule is in the form of $A \Rightarrow B$, where A and B are two disjoint itemsets, referred to respectively as the lhs (left-hand side) and rhs (right-hand side) of the rule. The three most widely-used measures for selecting interesting rules are support, confidence and lift. Support is the percentage of cases in the data that contains both A and B, confidence is the percentage of cases containing A that also contain B, and lift is the ratio of confidence to the percentage of cases containing B. The formulae to calculate them are:

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

where $P(A)$ is the percentage (or probability) of cases containing A. In addition to support, confidence and lift, there are many other interesting measures, such as chi-square, conviction, gini and leverage. An introduction to over 20 measures can be found in Tan et al.'s work [Tan et al., 2002].

The Titanic Dataset

The Titanic dataset 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, where each row represents a person. The reconstructed raw data can also be downloaded as file “titanic.raw.rdata” from the lab folder. Load this file into R as `titanic.raw` then run the `summary` command.

```
> # have a look at the 1st 5 lines
> readLines("./titanic.raw.rdata ", n=5)
[1] "1st  adult male   yes" "1st  adult male
[4] "1st  adult male   yes" "1st  adult male
yes" "1st  adult male   yes"
yes"
> # read it into R
> titanic.raw <- read.table("./titanic.raw.rdata ",
header=F)
> names(titanic) <- c("Class", "Sex", "Age", "Survived")

> summary(titanic.raw)
  Class      Sex      Age      Survived
1st :325   Female: 470   Adult:2092   No :1490
2nd :285   Male  :1731   Child: 109   Yes: 711
3rd :706
Crew:885
```

Now we have a dataset where each row stands for a person, and it can be used for association rule mining.

Association Rule Mining

A classic algorithm for association rule mining is APRIORI [Agrawal and Srikant, 1994]. It is a level-wise, breadth-first algorithm which counts transactions to find frequent itemsets and then derive association rules from them. An implementation of it is function `apriori()` in package `arules` [Hahsler et al., 2011]. Another algorithm for association rule mining is the ECLAT algorithm [Zaki, 2000], which finds frequent itemsets with equivalence classes, depth-first search and set intersection instead of counting. It is implemented as function `eclat()` in the same package.

Below we demonstrate association rule mining with `apriori()`. With the function, the default settings are:

- 1) `supp=0.1`, which is the minimum support of rules;
- 2) `conf=0.8`, which is the minimum confidence of rules;

3) maxlen=10, which is the maximum length of rules.

```
> library(arules)
> # find association rules with default settings
> rules.all <- apriori(titanic.raw)
parameter specification:
  confidence minval smax arem  aval originalSupport support minlen
maxlen target
      0.8      0.1      1 none FALSE                TRUE      0.1      1
algorithmic control:
  filter tree heap memopt load sort verbose
    0.1 TRUE TRUE  FALSE TRUE     2     TRUE
apriori - find association rules with the apriori algorithm
version 4.21 (2004.05.09)      (c) 1996-2004  Christian Borgelt
set item appearances ...[0 item(s)] done [0.00s].
set transactions ...[10 item(s), 2201 transaction(s)] done [0.00s].
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
set of 27 rules

> inspect(rules.all)
  lhs                rhs                support confidence    lift
1  {}                => {Age=Adult}    0.9504771  0.9504771  1.0000000
2  {Class=2nd}       => {Age=Adult}    0.1185825  0.9157895  0.9635051
3  {Class=1st}       => {Age=Adult}    0.1449341  0.9815385  1.0326798
4  {Sex=Female}      => {Age=Adult}    0.1930940  0.9042553  0.9513700
5  {Class=3rd}       => {Age=Adult}    0.2848705  0.8881020  0.9343750
6  {Survived=Yes}    => {Age=Adult}    0.2971377  0.9198312  0.9677574
7  {Class=Crew}      => {Sex=Male}    0.3916402  0.9740113  1.2384742
8  {Class=Crew}      => {Age=Adult}    0.4020900  1.0000000  1.0521033
9  {Survived=No}     => {Sex=Male}    0.6197183  0.9154362  1.1639949
10 {Survived=No}     => {Age=Adult}    0.6533394  0.9651007  1.0153856
11 {Sex=Male}        => {Age=Adult}    0.7573830  0.9630272  1.0132040
12 {Sex=Female,
    Survived=Yes}    => {Age=Adult}    0.1435711  0.9186047  0.9664669
13 {Class=3rd,
    Sex=Male}        => {Survived=No}  0.1917310  0.8274510  1.2222950
14 {Class=3rd,
    Survived=No}     => {Age=Adult}    0.2162653  0.9015152  0.9484870
15 {Class=3rd,
    Sex=Male}        => {Age=Adult}
16 {Sex=Male,
    Survived=Yes}    => {Age=Adult}
17 {Class=Crew,
    Survived=No}     => {Sex=Male}
18 {Class=Crew,
    0.2099046  0.9058824  0.9530818
    0.1535666  0.9209809  0.9689670
    0.3044071  0.9955423  1.2658514
    Survived=No}     => {Age=Adult}    0.3057701  1.0000000  1.0521033
19 {Class=Crew,
    Sex=Male}        => {Age=Adult}    0.3916402  1.0000000  1.0521033
```

20	{Class=Crew, Age=Adult}	=> {Sex=Male}	0.3916402	0.9740113	1.2384742
21	{Sex=Male, Survived=No}	=> {Age=Adult}	0.6038164	0.9743402	1.0251065
22	{Age=Adult, Survived=No}	=> {Sex=Male}	0.6038164	0.9242003	1.1751385
23	{Class=3rd, Sex=Male, Survived=No}	=> {Age=Adult}	0.1758292	0.9170616	0.9648435
24	{Class=3rd, Age=Adult, Survived=No}	=> {Sex=Male}	0.1758292	0.8130252	1.0337773
25	{Class=3rd, Sex=Male, Age=Adult}	=> {Survived=No}	0.1758292	0.8376623	1.2373791
26	{Class=Crew, Sex=Male, Survived=No}	=> {Age=Adult}	0.3044071	1.0000000	1.0521033
27	{Class=Crew, Age=Adult, Survived=No}	=> {Sex=Male}	0.3044071	0.9955423	1.2658514

As a common phenomenon for association rule mining, many rules generated above are un-interesting. Suppose that we are interested in only rules with rhs indicating survival, so we set

`rhs=c("Survived=No", "Survived=Yes")` in appearance to make sure that only "Survived=No" and "Survived=Yes" will appear in the rhs of rules. All other items can appear in the lhs, as set with `default="lhs"`. In the above result `rules.all`, we can also see that the left-hand side (lhs) of the first rule is empty. To exclude such rules, we set `minlen` to 2 in the code below. Moreover, the details of progress are suppressed with `verbose=F`. After association rule mining, rules are sorted by lift to make high-lift rules appear first.

```
> # rules with rhs containing "Survived" only
> rules <- apriori(titanic.raw, control = list(verbose=F),
+ parameter = list(minlen=2, supp=0.005, conf=0.8),
+ appearance = list(rhs=c("Survived=No", "Survived=Yes"),
+ default="lhs"))
> quality(rules) <- round(quality(rules), digits=3)
> rules.sorted <- sort(rules, by="lift")
> inspect(rules.sorted)
```

	lhs	rhs	support	confidence	lift
1	{Class=2nd, Age=Child}	=> {Survived=Yes}	0.006	1.000	3.096
2	{Class=2nd, Sex=Female, Age=Child}	=> {Survived=Yes}	0.064	0.972	3.010
3	{Class=1st, Sex=Female}	=> {Survived=Yes}	0.064	0.972	3.010
	...etc				

When other settings are unchanged, with a lower minimum support, more rules will be produced, and the associations between itemsets shown in the rules will be more likely to be by chance. In the above code, the minimum support is set to 0.005, so each rule is supported at least by 12

(=ceiling(0.005 * 2201)) cases, which is acceptable for a population of 2201.

Support, confidence and lift are three common measures for selecting interesting association rules. Besides them, there are many other interestingness measures, such as chi-square, conviction, gini and leverage [Tan et al., 2002]. More than twenty measures can be calculated with function `interestMeasure()` in the `arules` package.

Removing Redundancy

Some rules generated in the previous section provide little or no extra information when some other rules are in the result. For example, the above rule 2 provides no extra knowledge in addition to rule 1, since rule 1 tells us that all 2nd-class children survived. Generally speaking, when a rule (such as rule 2) is a super rule of another rule (such as rule 1) and the former has the same or a lower lift, the former rule (rule 2) is considered to be redundant. Other redundant rules in the above result are rules 4, 7 and 8, compared respectively with rules 3, 6 and 5.

Below we prune redundant rules. Note that the rules have already been sorted descending by lift.

```
> # find redundant rules
> subset.matrix <- is.subset(rules.sorted, rules.sorted)
> subset.matrix[lower.tri(subset.matrix, diag=T)] <- NA
> redundant <- colSums(subset.matrix, na.rm=T) >= 1
> which(redundant)
[1] 2 4 7 8
> # remove redundant rules
> rules.pruned <- rules.sorted[!redundant]
> inspect(rules.pruned)
  lhs                rhs                support confidence lift
1 {Class=2nd, Age=Child} => {Survived=Yes}
2 {Class=1st, Sex=Female} => {Survived=Yes}
3 {Class=2nd, Sex=Female} => {Survived=Yes}
4 {Class=Crew, Sex=Female} => {Survived=Yes}
5 {Class=2nd, Sex=Male, Age=Adult} => {Survived=No}
6 {Class=2nd, Sex=Male} => {Survived=No}
7 {Class=3rd, Sex=Male, Age=Adult} => {Survived=No}
8 {Class=3rd, Sex=Male} => {Survived=No}
```

In the code above, function `is.subset(r1, r2)` `r2` is a superset of `r1`. Function `lower.tri()` returns a logical matrix with TRUE in lower triangle. From the above results, we can see that rules 2, 4, 7 and 8 (before redundancy removal) are successfully pruned.

Interpreting Rules

While it is easy to find high-lift rules from data, it is not an easy job to understand the identified rules. It is not uncommon that the association

rules are misinterpreted to find their business meanings. For instance, in the above rule list `rules.pruned`, the first rule "{Class=2nd, Age=Child} => {Survived=Yes}" has a confidence of one and a lift of three and there are no rules on children of the 1st or 3rd classes. Therefore, it might be interpreted by users as children of the 2nd class had a higher survival rate than other children. This is wrong! The rule states only that all children of class 2 survived, but provides no information at all to compare the survival rates of different classes. To investigate the above issue, we run the code below to find rules whose rhs is "Survived=Yes" and lhs contains "Class=1st", "Class=2nd", "Class=3rd", "Age=Child" and "Age=Adult" only, and which contains no other items (default="none"). We use lower thresholds for both support and confidence than before to find all rules for children of different classes.

```
> rules <- apriori(titanic.raw,
+   parameter = list (minlen=3, sup=0.002, conf=0.2),
+   appearance = list(rhs=c("Survived=Yes"),
+                       lhs=c("Class=1st", "Class=2nd", "Class=3rd", "Age=Child",
+                             "Age=Adult"),
+                       default="none")
+   control = list(verbose=F))
> rules.sorted <- sort(rules, by="confidence")
> inspect(rules.sorted)
```

	lhs	rhs	support	confidence	lift
1	{Class=2nd, Age=Child}	=> {Survived=Yes}	0.010904134	1.0000000	3.0956399
2	{Class=1st, Age=Child}	=> {Survived=Yes}	0.002726034	1.0000000	3.0956399
3	{Class=1st, Age=Adult}	=> {Survived=Yes}	0.089504771	0.6175549	1.9117275
4	{Class=2nd, Age=Adult}	=> {Survived=Yes}	0.042707860	0.3601533	1.1149048
5	{Class=3rd, Age=Child}	=> {Survived=Yes}	0.012267151	0.3417722	1.0580035
6	{Class=3rd, Age=Adult}	=> {Survived=Yes}	0.068605179	0.2408293	0.7455209

In the above result, the first two rules show that children of the 1st class are of the same survival rate as children of the 2nd class and that all of them survived. The rule of 1st-class children didn't appear before, simply because of its support was below the threshold specified in Section 9.3. Rule 5 presents a sad fact that children of class 3 had a low survival rate of 34%, which is comparable with that of 2nd-class adults (see rule 4) and much lower than 1st-class adults (see rule 3).