

# Homework 3

## Pattern Mining and Social Network Analysis

BOUYSSOU Gatien , de POURTALES Caroline, LAMBA Ankit

17 novembre, 2020

### Contents

<b>Parameters in association rules</b>	<b>3</b>
Support . . . . .	3
Confidence . . . . .	3
Lift . . . . .	3
Leverage . . . . .	3
Conviction . . . . .	3
Coverage . . . . .	4
Jaccard Index . . . . .	4
Loevinger . . . . .	4
Laplace . . . . .	4
Rule interest . . . . .	4
Least Contradiction . . . . .	4
Tile . . . . .	5
<b>APRIORI algorithm</b>	<b>6</b>
Definition . . . . .	6
Example on Groceries data on R . . . . .	6
<b>Using Frequent itemsets to find rules</b>	<b>9</b>
Concept . . . . .	9
The different types of frequent itemsets mining . . . . .	9
Example on Adult data on R . . . . .	9
Example on mushroom data on Python with scikit-learn . . . . .	13
<b>Clustering with APRIORI algorithm</b>	<b>17</b>
Concept . . . . .	17
Affinity dissimilarity . . . . .	17
Example on tennis data on R . . . . .	17
Cluster the items . . . . .	20
Cluster the rules . . . . .	21
<b>Association Rule Classification</b>	<b>22</b>
Classification Based on Associations : CBA Algorithm . . . . .	22
Concept . . . . .	22
Example on tennis data on R . . . . .	22
Recall from Homework 1 . . . . .	22
Classification using homemade transactions . . . . .	22
Classification using rules . . . . .	23
Conclusion . . . . .	23

Regularized Class Association Rules for Multi-class Problems : RCAR Algorithm . . . . .	24
Concept . . . . .	24
Example on tennis data on R . . . . .	24
First Order Inductive Learner : FOIL Algorithm . . . . .	25
Concept . . . . .	25
Laplace accuracy . . . . .	25
Example on tennis data on R . . . . .	25
Classification Based on Multiple Class-association Rules : CMAR Algorithm . . . . .	26
Classification based on Predictive Association Rules : CPAR Algorithm . . . . .	26

## Parameters in association rules

There are parameters controlling the number of rules to be generated.

### Support

Support is an indication of how frequently the itemset appears in the dataset.

$$Support(A \rightarrow B) = \frac{\text{Number of transaction with both A and B}}{\text{Total Number of transaction}} = P(A \cap B)$$

### Confidence

Confidence is an indication of how often the rule has been found to be true.

This says how likely B is induced by A.

$$Confidence(A \rightarrow B) = \frac{\text{Number of transaction with both A and B}}{\text{Total Number of transaction with A}} = \frac{P(A \cap B)}{P(A)}$$

### Lift

Lift is the factor by which, the co-occurrence of A and B exceeds the expected probability of A and B co-occurring, had they been independent. So, higher the lift, higher the chance of A and B occurring together.

$$Lift(A \rightarrow B) = \frac{P(A \cap B)}{P(A) \times P(B)}$$

### Leverage

The leverage compares the frequency of A and B appearing together and the frequency that would be expected if A and B were independent.

$$Leverage(A \rightarrow B) = P(A \cap B) - P(A) \times P(B)$$

Therefore, if A and B independent :

$$Leverage(A \rightarrow B) = 0$$

### Conviction

Conviction compares the probability that A appears without B if they were dependent with the actual frequency of the appearance of A without B. If A and B are independent, then, conviction(A, B) = 1. On the other hand, when  $P(A \cap B)$  tends toward  $P(A)$ , conviction(A,B) tends toward infinity.

$$Conviction(A \rightarrow B) = \frac{P(A) \times P(\bar{B})}{P(A \cap \bar{B})}$$

or

$$Conviction(A \rightarrow B) = \frac{1 - P(B)}{1 - \frac{P(A \cap B)}{P(A)}}$$

## Coverage

The coverage of an association rule is the number of instances for which it predicts correctly.

$$Coverage(A, B) = \frac{P(A \cap B) - P(A \cap \bar{B})}{P(A \cap B)}$$

## Jaccard Index

The Jaccard coefficient assesses the distance between A and B as the fraction of cases covered by both with respect to the fraction of cases covered by A.

$$Jaccard(A, B) = \frac{P(A \cap B)}{P(A) + P(B) - P(A \cap B)}$$

or

$$Jaccard(A, B) = \frac{P(A \cap B)}{P(A \cup B)}$$

We can notice that :  $0 \leq Jaccard(A, B) \leq 1$

With the Jaccard measure, if A and B are not similar at all the Jaccard(A, B) will be equal to 0; Indeed, when  $P(A \cap B)$  increases then  $P(A \cup B)$  decrease assuming that P(A) and P(B) are constant. So, if  $P(A \cap B)$  increases Jaccard(A, B) should increase as well. The closer Jaccard(A, B) comes to 1 the more similar A to B.

## Loevinger

$$Loevinger(A \rightarrow B) = 1 - \frac{1}{Conviction(A \rightarrow B)} = 1 - \frac{P(A \cap \bar{B})}{P(A)P(\bar{B})}$$

When A and B are independent  $P(A \cap \neg B) = P(A) * P(\neg B)$

So,  $Loevinger(A, B) = 0$

When A and B are dependent  $P(A \cap \neg B)$  should tend towards 0 and  $Loevinger(A, B)$  towards 1.

## Laplace

$$Laplace(A \rightarrow B) = \frac{Support(A \rightarrow B) + 1}{P(A) + 2}$$

With Laplace when the A and B are independent  $Support(A \rightarrow B) = 0$  and :

$$Laplace(A, B) = \frac{1}{P(A)+2} \approx 0 \text{ when } P(A) \text{ is big.}$$

On the other end, when A and B are dependent the Laplace formula should tend towards 1.

## Rule interest

$$RI(A, B) = P(A)(P(A|B) - P(B))$$

or

$$RI(A, B) = P(A \cap B) - P(A)P(B)$$

If A and B are independent then  $P(A \cap B) = P(A) * P(B)$  and  $RI(A, B) = 0$

## Least Contradiction

$$LC(A, B) = \frac{P(A \cap B) - P(A \cap \bar{B})}{P(B)}$$

## Tile

It uses itemsets AND rows.

$$tile(X) = N_{X=1} * |X|$$

For example :  $tile(\{A\}) = N_{\{A=1\}} * 1$  or  $tile(AB) = N_{A=1, B=1} * 2$

# APRIORI algorithm

## Definition

APRIORI searches for frequent itemset browsing the lattice of itemsets in breadth.

The database is scanned at each level of lattice. Additionally, APRIORI uses a pruning technique based on the properties of the itemsets, which are: If an itemset is frequent, all its sub-sets are frequent and does not need to be considered.

## Example on Groceries data on R

The Groceries data set contains 30 days of real-world point-of-sale transaction data from a typical local grocery outlet. The data set contains 9835 transactions and the items are aggregated to 169 categories.

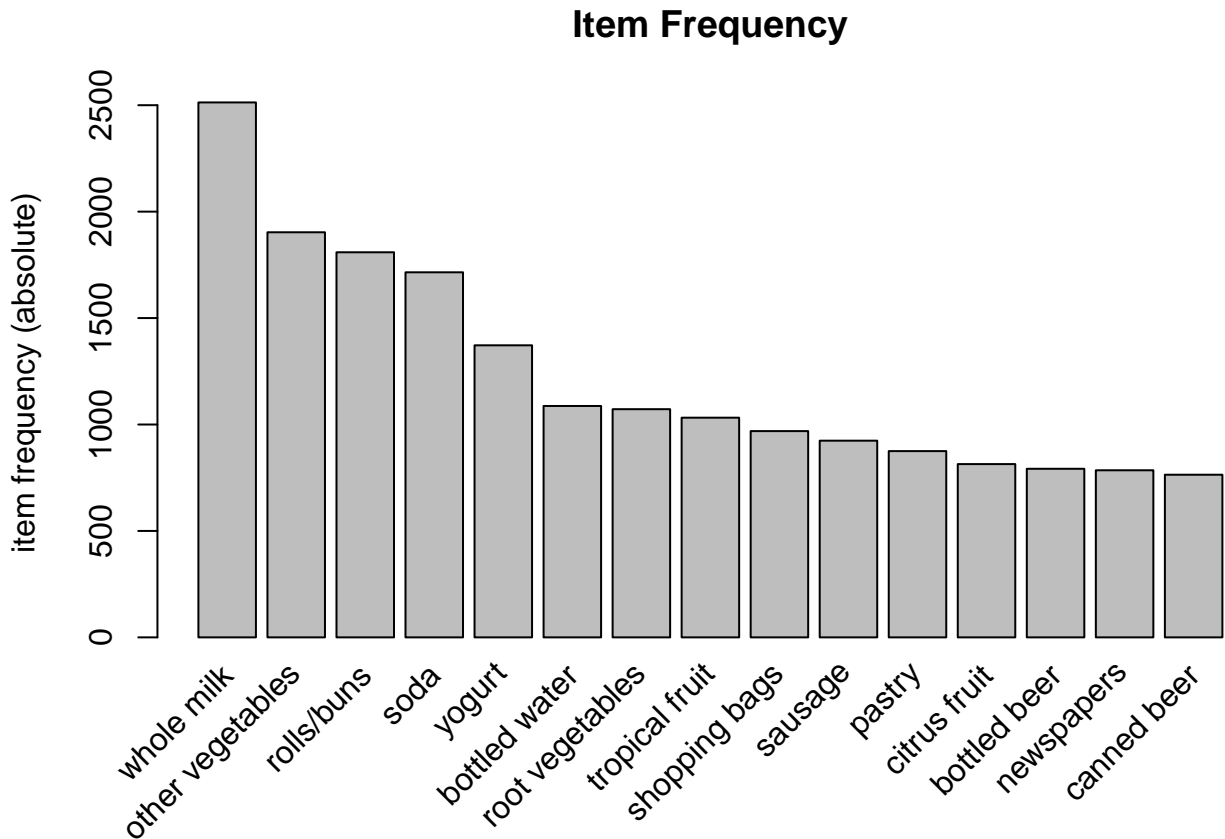
We can see the class of the dataset is :

```
## [1] "transactions"
## attr("package")
## [1] "arules"
```

Looking at some examples of transaction :

```
##      items
## [1] {citrus fruit,
##      semi-finished bread,
##      margarine,
##      ready soups}
## [2] {tropical fruit,
##      yogurt,
##      coffee}
## [3] {whole milk}
```

We can find the 15 most common variables.



Let's apply APRIORI algorithm on the dataset :

```
## Apriori
##
## Parameter specification:
## confidence minval smax arem aval originalSupport maxtime support minlen
##      0.2      0.1    1 none FALSE                TRUE      5   0.005      1
## maxlen target  ext
##      10   rules TRUE
##
## Algorithmic control:
## filter tree heap memopt load sort verbose
##      0.1 TRUE TRUE  FALSE TRUE    2    TRUE
##
## Absolute minimum support count: 49
##
## set item appearances ...[0 item(s)] done [0.00s].
## set transactions ...[169 item(s), 9835 transaction(s)] done [0.01s].
## sorting and recoding items ... [120 item(s)] done [0.00s].
## creating transaction tree ... done [0.01s].
## checking subsets of size 1 2 3 4 done [0.01s].
## writing ... [873 rule(s)] done [0.00s].
## creating S4 object ... done [0.00s].
```

We have a set of associations rules.

## set of 873 rules

If we look at the 3 rules with highest confidence, we have these rules :

##	lhs	rhs	support	confidence	coverage	lift	count
## [1]	{tropical fruit, root vegetables, yogurt}	=> {whole milk}	0.005693950	0.700	0.008134215	2.739554	56
## [2]	{pip fruit, root vegetables, other vegetables}	=> {whole milk}	0.005490595	0.675	0.008134215	2.641713	54
## [3]	{butter, whipped/sour cream}	=> {whole milk}	0.006710727	0.660	0.010167768	2.583008	66

However “whole milk” is the most frequent item in the data set and this frequency plays a role in confidence.

So we can look at the 3 rules with highest lift (A and B occurring together), and we have these rules :

##	lhs	rhs	support	confidence	coverage	lift	count
## [1]	{citrus fruit, other vegetables, whole milk}	=> {root vegetables}	0.005795628	0.4453125	0.01301474	4.085493	57
## [2]	{other vegetables, butter}	=> {whipped/sour cream}	0.005795628	0.2893401	0.02003050	4.036397	57
## [3]	{herbs}	=> {root vegetables}	0.007015760	0.4312500	0.01626843	3.956477	69

We can also look at the items which induce “soda”. Then we can sort them by confidence and look at the first 3 (so the 3 rules with highest confidence).

##	lhs	rhs	support	confidence
## [1]	{bottled water, fruit/vegetable juice}	=> {soda}	0.005185562	0.3642857
## [2]	{sausage, shopping bags}	=> {soda}	0.005693950	0.3636364
## [3]	{yogurt, bottled water}	=> {soda}	0.007422471	0.3230088

##	coverage	lift	count
## [1]	0.01423488	2.089067	51
## [2]	0.01565836	2.085343	56
## [3]	0.02297916	1.852357	73

Looking at the confidence, we see that for a third of the people buying : - bottled water and fruit/vegetable juice or - sausage and shopping bags or - yogurt and bottled water

it also induces buying soda.



# Using Frequent itemsets to find rules

## Concept

### Eclat algorithm :

It mines frequent itemsets

This algorithm uses simple intersection operations for equivalence class clustering along with bottom-up lattice traversal. Then looking at the most frequent itemsets, we can find rules between the items inside these itemsets.

## The different types of frequent itemsets mining

- **Max itemsets** : An itemset X is a maximal frequent itemset (or max-itemset) in set S if X is frequent, and there exists no super-itemset Y such that XY and Y is frequent in S.
- **Closed itemsets** : An itemset X is a closed frequent itemset in set S if X is both closed and frequent in S.
- **Free or Generator itemsets** : Free itemsets are itemsets that are not included in any closure of their proper sub-set. It means that it has no subset with the same support.
- **Largest tiles** : It is the itemset with the biggest area in a dataset. The area is equal to the frequency multiplied by the size of the dataset.

## Example on Adult data on R

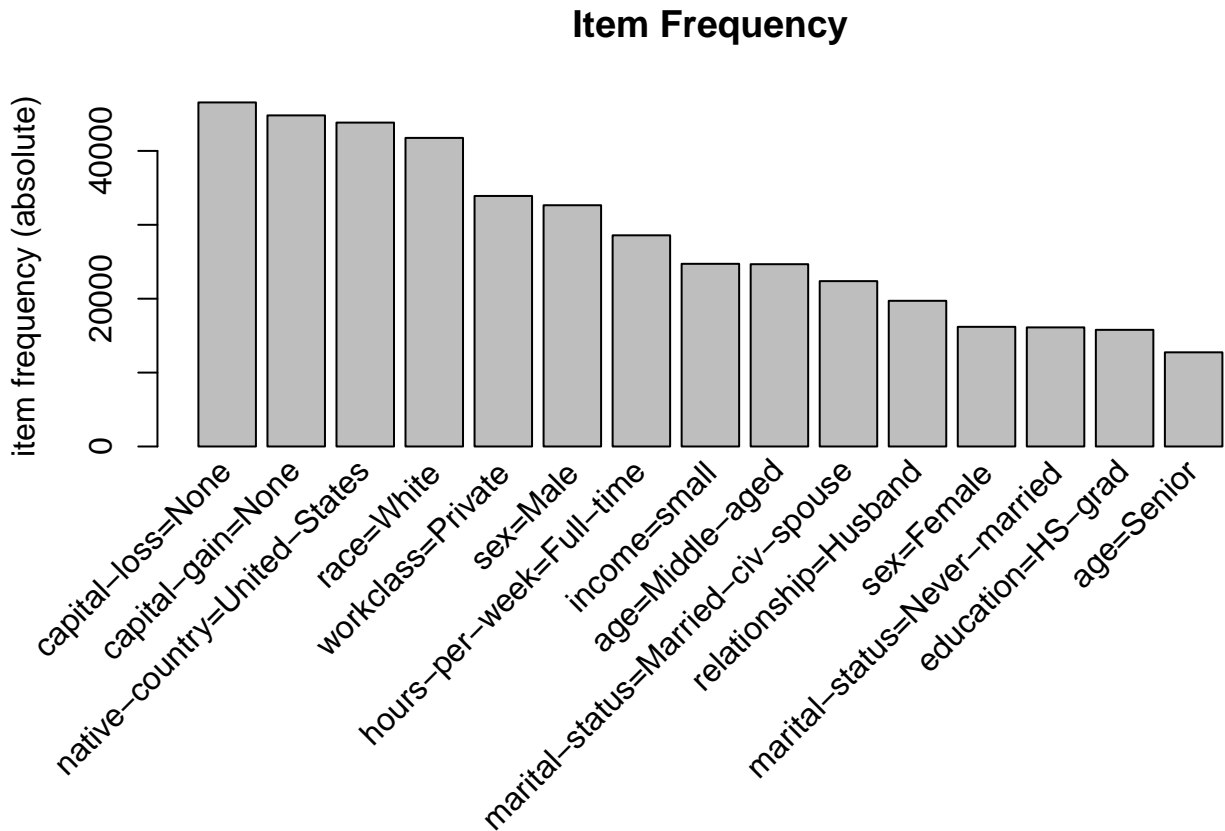
The Adult data set from R contains 48842 observations on the 15 variables (age, workclass, ...).

```
## [1] "transactions"
## attr(,"package")
## [1] "arules"
```

We can look at the first transaction to see what are the items in a transaction.

```
##      items                                transactionID
## [1] {age=Middle-aged,
##      workclass=State-gov,
##      education=Bachelors,
##      marital-status=Never-married,
##      occupation=Adm-clerical,
##      relationship=Not-in-family,
##      race=White,
##      sex=Male,
##      capital-gain=Low,
##      capital-loss=None,
##      hours-per-week=Full-time,
##      native-country=United-States,
##      income=small}                                1
```

We can find the 15 most common variables.



We apply ECLAT algorithm on the data set.  
This returns the most frequent itemsets along with their support.  
Let's look at 3 of these itemsets :

```
## Eclat
##
## parameter specification:
## tidLists support minlen maxlen          target  ext
##   FALSE    0.01      1    100 frequent itemsets TRUE
##
## algorithmic control:
## sparse sort verbose
##      7    -2    TRUE
##
## Absolute minimum support count: 488
##
## create itemset ...
## set transactions ...[115 item(s), 48842 transaction(s)] done [0.09s].
## sorting and recoding items ... [67 item(s)] done [0.01s].
## creating bit matrix ... [67 row(s), 48842 column(s)] done [0.01s].
## writing ... [80228 set(s)] done [0.37s].
## Creating S4 object ... done [0.03s].

##      items                                support  transIdenticalToItemsets
## [1] {education=5th-6th,capital-loss=None} 0.01009377 493
## [2] {education=Doctorate,capital-loss=None} 0.01076942 526
## [3] {education=Doctorate,race=White}       0.01076942 526
##      count
```

```
## [1] 493
## [2] 526
## [3] 526
```

We can also find the rules from the most frequent itemsets.

We use the ruleInduction function from R. We can set the method with the argument “method”.

If in control method = “APRIORI” is used, a very simple rule induction method is used. All rules are mined from the transactions data set using APRIORI with the minimal support found in itemsets. Then, all rules which do not stem from one of the itemsets are removed. The drawback of this procedure is that it is very slow in many cases.

	lhs	rhs	support	confidence	lift
## [1]	{marital-status=Married-civ-spouse, sex=Female, capital-gain=None, native-country=United-States, income=large}	=> {relationship=Wife}	0.01095369	0.9870849	20.68263
## [2]	{marital-status=Married-civ-spouse, race=White, sex=Female, capital-gain=None, income=large}	=> {relationship=Wife}	0.01076942	0.9868668	20.67806
## [3]	{marital-status=Married-civ-spouse, race=White, sex=Female, native-country=United-States, income=large}	=> {relationship=Wife}	0.01238688	0.9837398	20.61254

If in control method = “ptree” is used, the transactions are counted into a prefix tree and then the rules are selectively generated using the counts in the tree. This is usually faster than the above approach.

We can also find the rules with a specific given result.

For example, let’s answer the question :

How to be rich ?

```
## Eclat
##
## parameter specification:
## tidLists support minlen maxlen          target ext
## FALSE      0.01      1      200 frequent itemsets TRUE
##
## algorithmic control:
## sparse sort verbose
##      7      -2      TRUE
##
## Absolute minimum support count: 488
##
## create itemset ...
## set transactions ... [115 item(s), 48842 transaction(s)] done [0.09s].
## sorting and recoding items ... [67 item(s)] done [0.01s].
## creating bit matrix ... [67 row(s), 48842 column(s)] done [0.01s].
## writing ... [80228 set(s)] done [0.40s].
## Creating S4 object ... done [0.04s].
```

We take the 3 best rules according to lift.

```
## set of 14 rules
```

	lhs	rhs	support	confidence	lift
## [1]	{capital-loss=None, hours-per-week=Over-time, income=large}	=> {capital-gain=High}	0.01148602	0.1817887	5.253802
## [2]	{race=White, capital-loss=None, hours-per-week=Over-time, income=large}	=> {capital-gain=High}	0.01052373	0.1779778	5.143665
## [3]	{capital-loss=None, hours-per-week=Over-time, native-country=United-States, income=large}	=> {capital-gain=High}	0.01046231	0.1779248	5.142132

We see a pattern for people with a high capital gain : they have often a large income, work over-time and have no capital loss.

Differences between ECLAT and APRIORI:

- Apriori algorithm is a classical algorithm used to mining the frequent item sets in a given dataset.
- Coming to Eclat algorithm also mining the frequent itemsets but in vertical manner and it follows the depth first search of a graph.
- As per the speed,Eclat is faster than the Apriori algorithm.
- Apriori works on larger datasets where as Eclat algorithm works on smaller datasets.

## Example on mushroom data on Python with scikit-learn

This database contains a lot of mushrooms with a set of characteristics. Each mushroom is classified either as edible or poisonous. The database has been found in kaggle and is available here : <https://www.kaggle.com/uciml/mushroom-classification>.

First, we want to have an overview of the data.

```
##   class cap-shape cap-surface ... spore-print-color population habitat
## 0      p         x           s ...                k           s       u
## 1      e         x           s ...                n           n       g
##
## [2 rows x 23 columns]
```

As we can see, each column contains values that are single characters. Their meaning is given by the file values\_name.txt.

```
## 8124
```

Now, we want to know the data repartition for each columns.

```
## e      4208
## p      3916
## Name: class, dtype: int64

## <bound method IndexOpsMixin.value_counts of 0      e
## 1         c
## 2         c
## 3         e
## 4         e
##      ..
## 8119      ?
## 8120      ?
## 8121      ?
## 8122      ?
## 8123      ?
## Name: stalk-root, Length: 8124, dtype: object>
```

We can't print the distribution for each column because it would take too much place. We have just displayed two features. As you can see, there is almost as much poisonous as edible mushrooms. Moreover, the dataset contains some unknown values in the column stalk-root. We are going to discard those rows to keep lines that are complete.

```
## 5644
```

```
## e      3488
## p      2156
## Name: class, dtype: int64
```

Even without the discarded lines the dataset still have plenty of data and the class label is almost balanced. Before feeding the APRIORI algorithm with our data, we need to use the TransactionEncoder provided by mlxtend. This class transforms our data into a matrix where :

- each possible value for each feature will become a column
- for each mushroom and each column we assign a boolean that correspond to whether or not the feature is contained by the mushroom.

For example, such a dataset :

Will be changed into this matrix :

```
Columns : odor
0         pungent
1         almond
2         anise
3         pungent
4         none
```

```
Columns : pungent  almond  anise  none
0         True     False   False False
1         False    True    False False
2         False    False   True  False
3         True     False   False False
4         False    False   False True
```

The mushroom dataset contains character values. In order to have columns that are a bit more intelligible, we will replace the character values by their full name.

This function below split for one feature the character values from the full name values. It returns two arrays with each type of values.

This function goes through all the features and maps the feature's names with the character and full name values. Therefore it changes this line :

```
cap-shape: bell=b,conical=c,convex=x,flat=f, knobbed=k,sunken=s
```

into that dictionary :

```
{'cap-shape': [['b', 'c', 'x', 'f', 'k', 's'], ['bell', 'conical', 'convex', 'flat', 'knobbed', 'sunken']]}
```

Then, we will replace the values in the first array by the values of the second one for a given feature.

```
##      class cap-shape cap-surface ... spore-print-color population  habitat
## 0  poisonous    convex    smooth ...           black  scattered    urban
## 1    edible    convex    smooth ...           brown   numerous  grasses
## 2    edible     bell     smooth ...           brown   numerous  meadows
##
## [3 rows x 23 columns]

##      abundant  almond  anise  attached  bell  ...    two  urban  white  woods  yellow
## 0      False   False  False    False  False ...   False  True   True  False  False
## 1      False    True  False    False  False ...   False  False  True  False   True
##
## [2 rows x 64 columns]

##      support                                itemsets  length
## 0      0.875266                                (broad)      1
## 1      0.642452                                (brown)       1
## 2      0.669029                                (bulbous)     1
## 3      0.818568                                (close)      1
## 4      0.618001                                (edible)     1
## ..      ...                                ...          ...
## 186  0.616584      (partial, pendant, free, white, smooth)  5
## 187  0.608079  (partial, bulbous, broad, close, free, white)  6
## 188  0.711552      (partial, broad, close, free, white, one)  6
## 189  0.600992      (partial, broad, free, white, one, smooth) 6
## 190  0.603827  (partial, bulbous, close, free, white, one)  6
##
## [191 rows x 3 columns]
```

Thanks to the APRIORI algorithm it is possible to associate some feature together.

```
##      support      itemsets  length
## 4    0.618001      (edible)      1
## 31   0.618001    (free, edible)    2
## 32   0.609497    (one, edible)    2
## 33   0.618001    (partial, edible) 2
## 34   0.618001    (edible, white)  2
## 89   0.609497    (one, edible, free) 3
## 90   0.618001    (partial, edible, free) 3
## 91   0.618001    (free, edible, white) 3
## 92   0.609497    (one, edible, partial) 3
## 93   0.609497    (one, edible, white) 3
## 94   0.618001    (partial, edible, white) 3
## 151  0.609497    (one, edible, partial, free) 4
## 152  0.609497    (one, edible, white, free) 4
## 153  0.618001    (partial, edible, white, free) 4
## 154  0.609497    (one, edible, partial, white) 4
## 184  0.609497    (partial, edible, free, white, one) 5

## Empty DataFrame
## Columns: [support, itemsets, length]
## Index: []
```

With the APRIORI algorithm, we can see some associations containing the *edible* feature with a support around 0.6. Also, it seems that the APRIORI haven't found any associations with the *poisonous* feature with a support above 0.6.

The result given by APRIORI will be used by the `association_rules` function given by `mixtend`.

```
##      antecedents      consequents ... leverage conviction
## 19      (edible)      (free) ... 0.001971      inf
## 20      (edible)      (one) ... 0.008577    2.008505
## 21      (edible)    (partial) ... 0.000000      inf
## 22      (edible)      (white) ... 0.000000      inf
## 152    (one, edible)      (free) ... 0.001944      inf
## ..      ...      ... ...      ...      ...
## 818 (partial, edible)    (free, white, one) ... 0.008577    2.008505
## 819    (free, edible)    (partial, white, one) ... 0.008577    2.008505
## 820    (edible, white)    (partial, one, free) ... 0.008577    2.008505
## 821    (one, edible)    (partial, white, free) ... 0.001944      inf
## 822      (edible) (partial, white, one, free) ... 0.008577    2.008505
##
## [65 rows x 9 columns]
```

Here we are listing all the rules that are implied by *edible*. We need to know the rules where *edible* is implied (ie the rules where *edible* is contained by the *consequents* column). But before searching for those rules, we are going to try out another algorithm, named `fpgrowth`, to see if we can obtain different results.

The results obtained by `fpgrowth` look similar to the results obtained by the APRIORI algorithm. Therefore, now we can look for the rules that implies *edible*.

```
## Empty DataFrame
## Columns: [antecedents, consequents, antecedent support, consequent support, support, confidence, lift]
## Index: []

##      antecedents      consequents ... leverage conviction
## 1284      (free)      (edible) ... 0.001971    1.005203
```

```

## 1286      (partial)                (edible) ... 0.000000 1.000000
## 1289      (white)                (edible) ... 0.000000 1.000000
## 1290      (one)                (edible) ... 0.008577 1.023637
## 1293 (partial, free)                (edible) ... 0.001971 1.005203
## ...      ...                ... ... ...
## 1408      (one, white)      (partial, edible, free) ... 0.008577 1.023637
## 1409      (partial)      (free, edible, white, one) ... 0.000000 1.000000
## 1411      (free)      (partial, edible, white, one) ... 0.001944 1.005019
## 1412      (white)      (partial, edible, one, free) ... 0.000000 1.000000
## 1413      (one)      (partial, edible, white, free) ... 0.008577 1.023637
##
## [65 rows x 9 columns]

```

As we can see above if the threshold is above 0.6 the `association_rules` function does not find any rules where edible is implied. The confidence and the consequent support of the rules lies around 60%. Therefore, they cannot be considered as reliable.



# Clustering with APRIORI algorithm

## Concept

We can find a dissimilarity between transactions so we can compare the data. Then this dissimilarity is used as distance measure in clustering.

So a direct approach to cluster itemsets is to define a distance metric between two itemsets  $X_i$  and  $X_j$ .

## Affinity dissimilarity

A good choice is the Affinity defined as :

$$A(X_i, X_j) = \frac{\text{Support}(X_i, X_j)}{P(X_i) + P(X_j) - \text{Support}(X_i, X_j)} = \frac{P(X_i \cap X_j)}{P(X_i \cup X_j)}$$

Here this means that affinity is the Jaccard similarity between items.

The Jaccard distance defined as :

$$J(X_i, X_j) = \frac{|X_i \cap X_j|}{|X_i \cup X_j|}$$

The distance simply is the number of items that  $X_i$  and  $X_j$  have in common divided by the number of unique items in both sets.

## Example on tennis data on R

We use a dataset from the Wimbledon tennis tournament for Women in 2013. We will predict the result for player 1 (win=1 or loose=0) based on : the number of aces won by each player, and, the number of unforced errors committed by both players. The data set is a subset of a data set from <https://archive.ics.uci.edu/ml/datasets/Tennis+Major+Tournament+Match+Statistics>.

```
##      Result ACE.1 UFE.1 ACE.2 UFE.2
## 1         0     2    18     3     14
## 2         0     0    10     4     14
## 3         1     1    13     2     29
## 4         1     4    30     0     45
## 5         0     2    28     6     19
## 6         0     6    42    11     40
```

We can transform the tennis data set into a transaction data set.

We can look at the 3 first transactions.

```
##      items                                     transactionID
## [1] {Result=0,ACE.1=Low,UFE.1=Low,ACE.2=Low,UFE.2=Low}      1
## [2] {Result=0,ACE.1=None,UFE.1=Low,ACE.2=High,UFE.2=Low}    2
## [3] {Result=1,ACE.1=Low,UFE.1=Low,ACE.2=Low,UFE.2=High}    3
```

We can restrict the rules to the result rhs="Result=1" which means Player-1 winner.

The associations rules for Player-1 winning are :

```
## Apriori
##
## Parameter specification:
## confidence minval smax arem aval originalSupport maxtime support minlen
```

```

##      0.3    0.1    1 none FALSE          TRUE      5    0.15    1
## maxlen target ext
##      10 rules TRUE
##
## Algorithmic control:
## filter tree heap memopt load sort verbose
##      0.1 TRUE TRUE  FALSE TRUE    2    TRUE
##
## Absolute minimum support count: 17
##
## set item appearances ...[1 item(s)] done [0.00s].
## set transactions ...[12 item(s), 118 transaction(s)] done [0.00s].
## sorting and recoding items ... [11 item(s)] done [0.00s].
## creating transaction tree ... done [0.00s].
## checking subsets of size 1 2 3 4 done [0.00s].
## writing ... [13 rule(s)] done [0.00s].
## creating S4 object ... done [0.00s].

```

These are the 5 rules with highest lift.

```

## set of 13 rules

##      lhs                                rhs      support  confidence coverage
## [1] {ACE.1=High,UFE.1=Low} => {Result=1} 0.1525424 0.8181818 0.1864407
## [2] {ACE.1=High}           => {Result=1} 0.2881356 0.6938776 0.4152542
## [3] {ACE.1=High,UFE.2=Low} => {Result=1} 0.1694915 0.6451613 0.2627119
## [4] {ACE.2=Low}            => {Result=1} 0.2457627 0.6170213 0.3983051
## [5] {UFE.1=Low}           => {Result=1} 0.3220339 0.6129032 0.5254237
##      lift      count
## [1] 1.532468 18
## [2] 1.299644 34
## [3] 1.208397 20
## [4] 1.155691 29
## [5] 1.147977 38

```

These rules look correct : either a player-1 winning make a lot of aces and few unforced errors or the player-2 make few aces.

We can also restrict the rules to the result rhs="Result=0" which means Player-1 loosing.

The associations rules for Player-1 loosing :

```

## Apriori
##
## Parameter specification:
## confidence minval smax arem aval originalSupport maxtime support minlen
##      0.3    0.1    1 none FALSE          TRUE      5    0.15    1
## maxlen target ext
##      10 rules TRUE
##
## Algorithmic control:
## filter tree heap memopt load sort verbose
##      0.1 TRUE TRUE  FALSE TRUE    2    TRUE
##
## Absolute minimum support count: 17
##
## set item appearances ...[1 item(s)] done [0.00s].

```

```
## set transactions ...[12 item(s), 118 transaction(s)] done [0.00s].
## sorting and recoding items ... [11 item(s)] done [0.00s].
## creating transaction tree ... done [0.00s].
## checking subsets of size 1 2 3 4 done [0.00s].
## writing ... [10 rule(s)] done [0.00s].
## creating S4 object ... done [0.00s].
```

These are the 3 rules with highest lift.

```
## set of 10 rules
```

```
##      lhs                rhs      support  confidence coverage lift      count
## [1] {ACE.2=High} => {Result=0} 0.2372881 0.6086957 0.3898305 1.305929 28
## [2] {UFE.1=High} => {Result=0} 0.2627119 0.5535714 0.4745763 1.187662 31
## [3] {ACE.1=Low}  => {Result=0} 0.2372881 0.5283019 0.4491525 1.133448 28
```

These rules look correct : either player-1 is loosing because player-2 makes a lot of aces or because he does a lot of unforced errors or player-2 makes a lot of aces.

Now let's look at all the associations rules leading to "Result".

All the rules with Result as association :

```
## Apriori
##
## Parameter specification:
## confidence minval smax arem aval originalSupport maxtime support minlen
##      0.4      0.1      1 none FALSE              TRUE        5      0.1      1
## maxlen target ext
##      10 rules TRUE
##
## Algorithmic control:
## filter tree heap memopt load sort verbose
##      0.1 TRUE TRUE  FALSE TRUE      2      TRUE
##
## Absolute minimum support count: 11
##
## set item appearances ...[2 item(s)] done [0.00s].
## set transactions ...[12 item(s), 118 transaction(s)] done [0.00s].
## sorting and recoding items ... [12 item(s)] done [0.00s].
## creating transaction tree ... done [0.00s].
## checking subsets of size 1 2 3 4 done [0.00s].
## writing ... [46 rule(s)] done [0.00s].
## creating S4 object ... done [0.00s].
```

These are the 5 rules with highest lift.

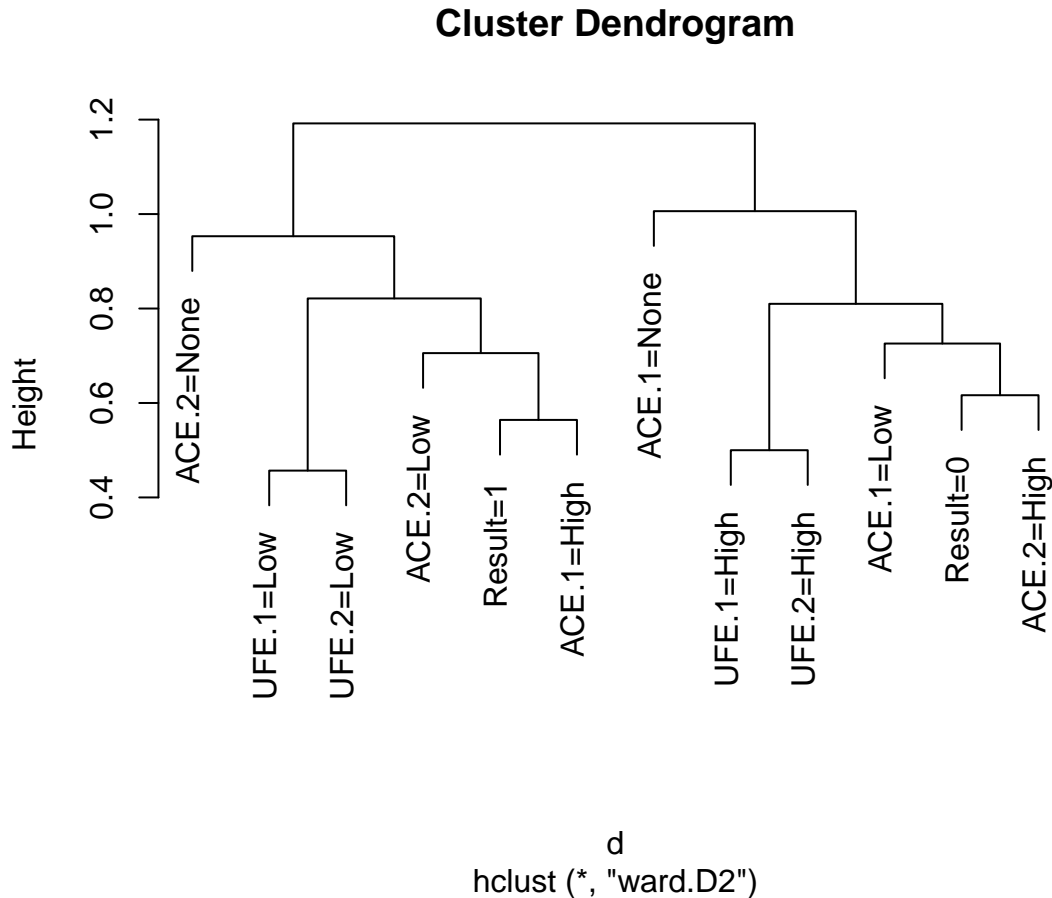
```
## set of 46 rules
```

```
##      lhs                rhs      support  confidence
## [1] {ACE.1=None}          => {Result=0} 0.1016949 0.7500000
## [2] {ACE.1=High,UFE.1=Low} => {Result=1} 0.1525424 0.8181818
## [3] {UFE.1=High,ACE.2=High} => {Result=0} 0.1186441 0.7000000
## [4] {ACE.2=High,UFE.2=Low}  => {Result=0} 0.1355932 0.6956522
## [5] {ACE.1=High,UFE.1=Low,UFE.2=Low} => {Result=1} 0.1271186 0.7894737
##      coverage lift      count
## [1] 0.1355932 1.609091 12
## [2] 0.1864407 1.532468 18
## [3] 0.1694915 1.501818 14
```

```
## [4] 0.1949153 1.492490 16
## [5] 0.1610169 1.478697 15
```

Firstly let's look at the clustering of items

Cluster the items



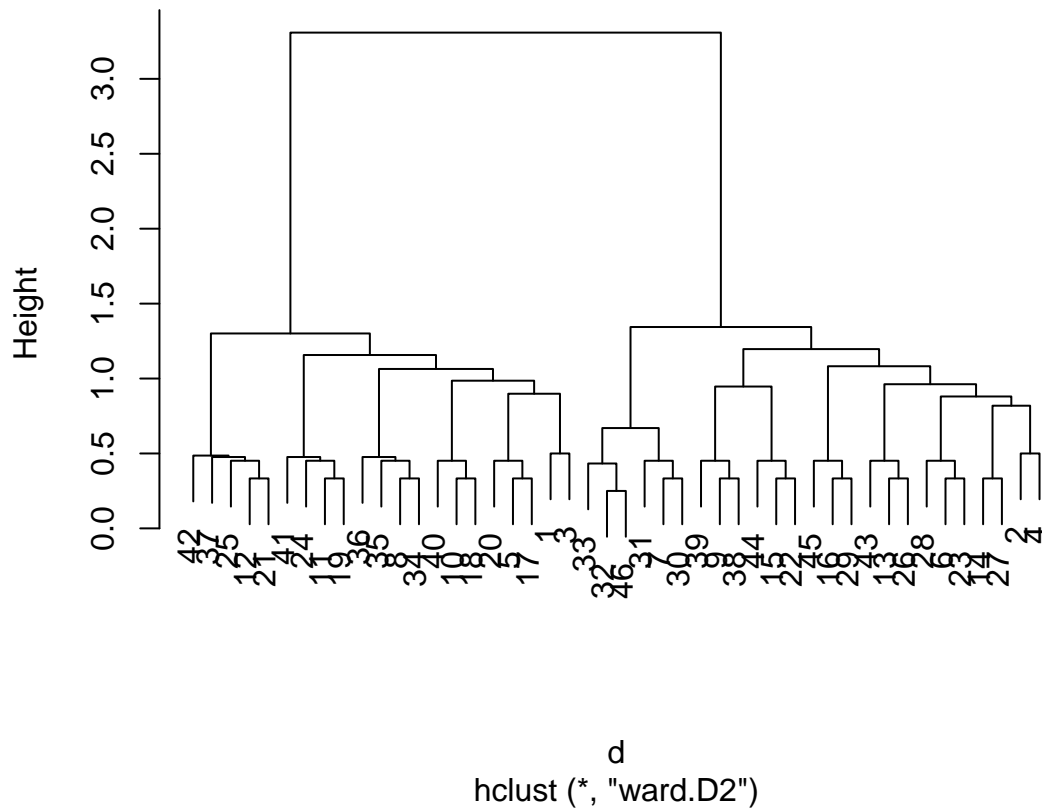
We can see two clusters on the dendrogram. One contains Result=1 and the other contains Result=0. Moreover, in the cluster with the branch Result=1, we can also see that the number of aces made by player-1 is high and the number of unforced errors is low. Also, in this cluster player-2 made few or none aces. In the other cluster, it is the opposite. Player-1 makes few or none aces whereas player-2 makes a lot.

So, it seems that this clustering manage to cluster data linked to the result together.

Now let's try to cluster the rules.

Cluster the rules

## Cluster Dendrogram



If we cut the dendrogram in two clusters. We can look at the first cluster. (We only print 5 items from the cluster, look at the code for the whole cluster)

```
##      lhs      rhs      support  confidence coverage lift      count
## [1] {}      => {Result=0} 0.4661017 0.4661017 1.0000000 1.0000000 55
## [2] {ACE.1=None} => {Result=0} 0.1016949 0.7500000 0.1355932 1.6090909 12
## [3] {ACE.2=High} => {Result=0} 0.2372881 0.6086957 0.3898305 1.3059289 28
## [4] {ACE.1=Low}  => {Result=0} 0.2372881 0.5283019 0.4491525 1.1334477 28
## [5] {UFE.2=High} => {Result=0} 0.1864407 0.4000000 0.4661017 0.8581818 22
```

And at the second cluster. (We only print 5 items from the cluster, look at the code for the whole cluster)

```
##      lhs      rhs      support  confidence coverage lift      count
## [1] {}      => {Result=1} 0.5338983 0.5338983 1.0000000 1.0000000 63
## [2] {ACE.2=None} => {Result=1} 0.1355932 0.6400000 0.2118644 1.1987302 16
## [3] {ACE.2=Low}  => {Result=1} 0.2457627 0.6170213 0.3983051 1.1556906 29
## [4] {ACE.1=High} => {Result=1} 0.2881356 0.6938776 0.4152542 1.2996437 34
## [5] {ACE.1=Low}  => {Result=1} 0.2118644 0.4716981 0.4491525 0.8834981 25
```

This clustering regroups Player-1 winner together and Player-2 winner together.

# Association Rule Classification

## Classification Based on Associations : CBA Algorithm

### Concept

CBA (Classification Based on Associations) Algorithm build a classifier based on association rules mined for an input dataset.

Candidate classification association rules (CARs) are mined with the standard APRIORI algorithm. Rules are ranked by confidence, support and size.

It uses either M1 or M2 pruning strategy.

Explication pruning M1 and M2 techniques so, either the M1 or M2 algorithm are used to perform database coverage pruning and to determine the number of rules to use and the default class.

M1 is the naive version of the algorithm. M2 is the improved version of the algorithm.

M1 is not always the right choice for pruning the mined rules. It traverses the dataset multiple times and keeps rules in memory so this may not be the optimal solution for very large datasets.

### Example on tennis data on R

**Recall from Homework 1** With Random Forest, the accuracy rate was 0.6931818.

With Logistic regression it was 0.7667.

CBA can take as input a classic non-transaction dataset as tennis. We just have to choose the discretization method in parameter.

However, discretization of integers is difficult so we did the transactions set ourselves.

**Classification using homemade transactions** We can also use the transactions we created before to train the classifier.

ACE.1 and ACE.2 take for value either None, or Low or High.

UFE.1 and UFE.2 take for value either Low or High.

```
## CBA Classifier Object
## Class:
## Default Class: NA
## Number of rules: 15
## Classification method: first
## Description: CBA algorithm (Liu et al., 1998)
```

The 4 rules with highest confidence are :

##	lhs	rhs	support	confidence		
## [1]	{ACE.1=Low,UFE.1=Low,UFE.2=High}	=> {Result=1}	0.10638298	1		
## [2]	{UFE.1=Low,ACE.2=Low,UFE.2=High}	=> {Result=1}	0.06382979	1		
## [3]	{ACE.1=High,UFE.1=Low,ACE.2=None}	=> {Result=1}	0.05319149	1		
## [4]	{ACE.1=Low,ACE.2=High,UFE.2=Low}	=> {Result=0}	0.05319149	1		
##	coverage	lift	count	size	coveredTransactions	totalErrors
## [1]	0.10638298	1.740741	10	4	10	40
## [2]	0.06382979	1.740741	6	4	2	40
## [3]	0.05319149	1.740741	5	4	5	37
## [4]	0.05319149	2.350000	5	4	5	35

We have the following confusion matrix :

```
##
## classifierTransactions.prediction 0 1
```

```
##                0 9 1
##                1 6 8
```

The accuracy rate is :

```
## [1] 0.7083333
```

So the accuracy rate is good.

This classification is also almost as good as logistic regression.

**Classification using rules** CBA\_ruleset creates a new object of class CBA using the provides rules as the rule base.

With the method “first” for “first found rule”:

```
## CBA Classifier Object
## Class: Result=0, Result=1
## Default Class: Result=1
## Number of rules: 46
## Classification method: first
## Description: Custom rule set
```

We have the following confusion matrix :

```
##
## classifierRules.prediction  0  1
##                0 10  1
##                1  5  8
```

The accuracy rate is :

```
## [1] 0.75
```

With the method “majority” :

Majority selection of the class label requires selecting a group of good quality rules matching the case to be classified, and assigning the appropriate class with **simple majority voting** among selected rules.

Let’s look at the method behind CBA\_ruleset with “majority” as classification method :

```
## CBA Classifier Object
## Class: Result=0, Result=1
## Default Class: Result=1
## Number of rules: 46
## Classification method: majority
## Description: Custom rule set
```

We have the following confusion matrix :

```
##
## classifierRules.prediction  0  1
##                0 11  1
##                1  4  8
```

The accuracy rate is :

```
## [1] 0.7916667
```

## Conclusion

The best classification using CBA is the classification using rules with majority method.

## Regularized Class Association Rules for Multi-class Problems : RCAR Algorithm

### Concept

Regularized Class Association Rules (RCAR) is an algorithm which produces rules based classifier in a categorical data space. The main goal of RCAR algorithm is to build classifiers which are as accurate as the state of the art algorithms, while improving the interpretability and allowing end-users to maintain and understand its outcome easily and without statistical modeling background.

### Example on tennis data on R

The elastic net mixing parameter for  $\alpha = 1$  is the lasso penalty (default RCAR), and for  $\alpha = 0$  it is the ridge penalty.

RCAR uses logistic regression.

Let's look at the method behind RCAR :

```
## CBA Classifier Object
## Class:
## Default Class: NA
## Number of rules: 73
## Classification method: logit
## Description: RCAR+ based on RCAR (Azmi et al., 2019)
```

The 4 rules with highest confidence are :

```
##      lhs                                rhs      support  confidence
## [1] {ACE.1=High,UFE.1=Low,ACE.2=None} => {Result=1} 0.05319149 1
## [2] {ACE.1=Low,ACE.2=High,UFE.2=Low}   => {Result=0} 0.05319149 1
## [3] {UFE.1=Low,ACE.2=Low,UFE.2=High}   => {Result=1} 0.06382979 1
## [4] {ACE.1=Low,UFE.1=Low,UFE.2=High}   => {Result=1} 0.10638298 1
##      coverage  lift    count weight  oddsratio
## [1] 0.05319149 1.740741  5      0.01537972 1.015499
## [2] 0.05319149 2.350000  5      0.02122279 1.021450
## [3] 0.06382979 1.740741  6      0.01548662 1.015607
## [4] 0.10638298 1.740741 10      0.01632068 1.016455
```

We have the following confusion matrix :

The accuracy rate is :

```
## [1] 0.6923077
```

The accuracy rate is not very good.



## First Order Inductive Learner : FOIL Algorithm

### Concept

FOIL learns rules and then use them as a classifier.

For each class, we find the positive and negative examples and learn the rules using FOIL. Then, the rules for all classes are combined and sorted by Laplace accuracy on the training data.

### Laplace accuracy

Laplace accuracy is used to measure the accuracy of the rules. Given a rule  $r$  it is defined as follows:

$$LaplaceAccuracy(r) = \frac{N_c + 1}{N_{total} + m}$$

where  $m$  is the number of classes,  $N_{total}$  is the total number of examples that satisfies the rule's body and  $N_c$  is the number of examples belonging to the predicted class  $c$  of the rule.

We classify new examples by 1. select all the rules whose bodies are satisfied by the example; 2. from the rules select the best  $k$  rules per class (highest expected Laplace accuracy); 3. average the expected Laplace accuracy per class and choose the class with the highest average.

### Example on tennis data on R

We can use our homemade transactions as input data in FOIL algorithms.

Let's look at the method behind FOIL :

```
##      lhs      rhs      support  confidence lift
## [1] {ACE.2=Low,UFE.2=High} => {Species=1} 0.14285714 0.8750000 1.559091
## [2] {ACE.1=High,ACE.2=Low}  => {Species=1} 0.15306122 0.8333333 1.484848
## [3] {ACE.1=High,UFE.2=High} => {Species=1} 0.11224490 0.7333333 1.306667
## [4] {ACE.1=None,UFE.2=Low}  => {Species=0} 0.06122449 0.7500000 1.709302
## [5] {ACE.1=None,ACE.2=High} => {Species=0} 0.03061224 0.7500000 1.709302
##      laplace
## [1] 0.8333333
## [2] 0.8000000
## [3] 0.7058824
## [4] 0.7000000
## [5] 0.6666667
```

The accuracy rate is :

```
## [1] 0.75
```

The accuracy rate is better than the one from RCAR and is enough good.

## **Classification Based on Multiple Class-association Rules : CMAR Algorithm**

CMAR selects a small set of high confidence, highly related rules and analyzes the correlation among those rules. To avoid bias, we develop a new technique, called weighted  $\chi^2$ , which derives a good measure on how strong the rule is under both conditional support and class distribution. An extensive performance study shows that CMAR in general has higher prediction accuracy than CBA.

## **Classification based on Predictive Association Rules : CPAR Algorithm**

CPAR inherits the basic idea of FOIL in rule generation and integrates the features of associative classification in predictive rule analysis. In comparison with associative classification, CPAR has the following advantages:

- CPAR generates a much smaller set of highquality predictive rules directly from the dataset;
- To avoid generating redundant rules, CPAR generates each rule by considering the set of “already generated” rules; and
- When predicting the class label of an example, CPAR uses the best  $k$  rules that the example satisfies.