

# Data Management

## Lab 02 - Analyse de données

Robin Chappatte

Frédéric Montet

Brian Nydegger

Rendu le 28 novembre 2016

à Lausanne

### **Professeurs :**

Dr. Laura Elena Raileanu

Fabien Dutoit

# Table des matières

<b>Introduction</b>	<b>1</b>
<b>1 Installation et configuration</b>	<b>2</b>
1.1 Installation . . . . .	2
1.2 Configuration . . . . .	2
<b>2 Analyse de données</b>	<b>4</b>
2.1 Arbre de décisions . . . . .	4
2.2 Clustering . . . . .	8
2.2.1 2 clusters . . . . .	9
2.2.2 8 clusters . . . . .	11
2.3 Règles d'associations . . . . .	15
2.3.1 Apriori . . . . .	15
2.3.2 FilteredAssociator . . . . .	16
<b>3 Conclusion</b>	<b>18</b>

# Introduction

Pour le second laboratoire du cours de Data Management, nous allons mettre en pratique les concepts d'analyse de donnée vus en cours :

- Arbre de décisions
- Clustering
- Règles d'associations

Pour faire cela, nous utiliserons le logiciel Weka car il dispose d'une grande collection d'outils pour l'analyse de donnée et est également libre et multi-plateforme, ce qui en fait un outil idéal pour le cadre académique dans lequel ce laboratoire est réalisé.

# 1 | Installation et configuration

Dans notre situation nous avons décidé d'installer Weka sur une machine virtuelle Ubuntu 16.10.

## 1.1 Installation

Nous avons installé :

- Java JDK 1.8\_111
- Weka 3.8.0
- Xampp x64 7.0.13
- MySQL connector for Java 5.1.39

Un élément à mentionner est que nous avons utilisé le connecteur dans sa version 5.1.39 alors que sur le site de MySQL<sup>1</sup> il est disponible en 5.1.40. Ceci est dû au fait que nous avons utilisé la version disponible dans un repository ppa. L'installation a pu être faite avec la commande `sudo apt-get install libmysql-java`.

Hormis le point cité au paragraphe ci dessus, cette étape n'a pas posé de problème particulier qu'il serait nécessaire de mentionner. Il a suffit de suivre les différents tutoriaux relatifs à chacun des logiciels à installer.

## 1.2 Configuration

Pour configurer la machine virtuelle, un seul point a posé problème : la connexion à la base de donnée MySQL depuis Weka en utilisant le connecteur MySQL Java.

Le problème rencontré a été la prise en compte de la variable d'environnement `CLASSPATH` lors des tentatives de connexions. Sans cette variable, Weka ne trouve pas le connecteur.

Après avoir effectué une recherche, nous nous sommes rendu compte qu'avec Ubuntu, il fallait préciser le `CLASSPATH` dans la commande pour démarrer Weka à l'aide du paramètre `-cp`. Le listing Figure 1.1 présente le `.sh` qui a été utilisé pour démarrer Weka et le connecter à MySQL avec succès.

---

1. <https://dev.mysql.com/downloads/connector/j/>

## 1.2. Configuration

---

```
1 # Variables
2 java8=/home/fredmontet/Desktop/lab02-weka/asset/java/jdk1.8.0_111/jre/bin/
   ↪ java
3 weka=/home/fredmontet/Desktop/lab02-weka/asset/weka-3-8-0/weka.jar
4 connector=/usr/share/java/mysql-connector-java.jar
5
6 # Command
7 export CLASSPATH=$connector:$CLASSPATH
8 $java8 -Xmx300m -cp "$weka:$connector" weka.gui.GUIChooser
```

FIGURE 1.1 – Contenu du fichier bash à exécuter pour lancer Weka

## 2 | Analyse de données

Les données à dispositions dans la base de donnée qui a été importée à partir du fichier *aventure2014.sql.gz* n'ont pas été utilisées complètement pour ce laboratoire. Uniquement les attributs suivants ont été utilisés :

- MaritalStatus
- Gender
- YearlyIncome
- TotalChildren ;
- EnglishEducation
- EnglishOccupation
- HouseOwnerFlag
- NumberCarsOwned

### 2.1 Arbre de décisions

L'algorithme choisi pour cette méthode de classification a été REPTree. Après avoir effectué plusieurs runs aux paramètres différents, nous avons choisi un des meilleurs de nos tests.

Sur le listing Figure 2.1, on observe les informations du run choisis ainsi que les attributs sur la base desquels on veut trouver les classes de l'attributs **EnglishOccupation** qui sont :

- Professional
- Management
- Skilled Manual
- Clerical
- Manual

## 2.1. Arbre de décisions

```
1 Scheme:      weka.classifiers.trees.REPTree -M 2 -V 0.001 -N 10 -S 42 -L -1
   ↳ -P -I 0.0 -batch-size 200
2 Relation:    QueryResult-weka.filters.unsupervised.attribute.Remove-R1
   ↳ -7,9,11,14,16,18,21-25
3 Instances:   18484
4 Attributes:  8
5              MaritalStatus
6              Gender
7              YearlyIncome
8              TotalChildren
9              EnglishEducation
10             EnglishOccupation
11             HouseOwnerFlag
12             NumberCarsOwned
13 Test mode:   split 80.0% train, remainder test
```

FIGURE 2.1 – Informations du run REPTree

Le résultat de ce run a montré que la moyenne pondérée des instances correctement classifiées est 91.29%. La valeur de ce résultat est bonne et indique que notre modèle est bon. Cependant, ce seul résultat ne permet pas d'affirmer que le modèle est bon. Il est nécessaire d'observer les différentes métriques dans la partie **Detailed Accuracy By Class** du listing Figure 2.2 pour en savoir d'avantage. Sur ce dernier, on constate que 0.913 est la moyenne des métriques TP Rate, Precision, Recall et F-Measure pour les différentes classes. Notre analyse est donc confortée par la stabilité du résultat moyen de ces différentes métriques.

Maintenant, il est intéressant de savoir quelles sont les forces et faiblesses de notre modèle. Pour cela, nous baserons nos constats sur la F-Measure. Cette dernière étant robuste car elle prend en compte precision et recall. On y voit que la classe au meilleur score est **Manual** avec 0.970 et la moins bonne est **Skilled Manual** avec 0.882.

Afin d'améliorer notre modèle, on se concentrera donc sur la classe **Skilled Manual**. La matrice de confusion nous donne le détail de classification des instances de **Skilled Manual**. On y voit que **Skilled Manual** est souvent pris pour **Professional** et vice-versa. Ceci peut nous donner une piste pour de plus amples analyse ou on se questionnera sur les similarités et différences entre 2 instances de chacune des classes pour comprendre quel sont les facteurs qui rentrent en jeux.

## 2.1. Arbre de décisions

```

1 === Summary ===
2
3 Correctly Classified Instances      3375      91.2902 %
4 Incorrectly Classified Instances    322      8.7098 %
5 Kappa statistic                     0.8883
6 Mean absolute error                 0.0434
7 Root mean squared error            0.1549
8 Relative absolute error             13.8991 %
9 Root relative squared error         39.2271 %
10 Total Number of Instances          3697
11
12 === Detailed Accuracy By Class ===
13
14 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC      ROC Area  PRC Area
15  ↳ Class
16 0.904    0.043    0.899      0.904    0.902      0.860    0.982    0.963
17  ↳ Professional
18 0.912    0.019    0.909      0.912    0.910      0.892    0.992    0.970
19  ↳ Management
20 0.873    0.035    0.892      0.873    0.882      0.844    0.981    0.945
21  ↳ Skilled Manual
22 0.946    0.013    0.930      0.946    0.938      0.927    0.991    0.971
23  ↳ Clerical
24 0.972    0.005    0.968      0.972    0.970      0.966    0.997    0.981
25  ↳ Manual
26
27 Weighted Avg.
28
29 0.913    0.027    0.913      0.913    0.913      0.885    0.987    0.963
30
31 === Confusion Matrix ===
32
33      a      b      c      d      e  <-- classified as
34 1002    25    69    12     0 | a = Professional
35  46   568     9     0     0 | b = Management
36  62    32   802    20     3 | c = Skilled Manual
37   4     0    15   545    12 | d = Clerical
38   0     0     4     9   458 | e = Manual

```

FIGURE 2.2 – Résultats du run REPTree

Concernant l'arbre de décision, son début est visible dans le listing Figure 2.3. Sur cet exemple, on voit que la division commence par l'apport financier, puis l'éducation, les enfants, l'état civil et la possession d'une maison.

Il est intéressant de noter que juste après la division homme-femme, l'état civil et la possession d'une maison ont été interverti dans la branche M et F. Ceci met en évidence que la meilleure manière de diviser les classes est prise en compte en fonction de chacune des instances.



## 2.1. Arbre de décisions

```
1 YearlyIncome < 45000
2 |   YearlyIncome < 25000
3 |   |   EnglishEducation = Bachelors
4 |   |   |   YearlyIncome < 15000
5 |   |   |   |   TotalChildren < 1.5
6 |   |   |   |   |   Gender = M
7 |   |   |   |   |   |   MaritalStatus = M
8 |   |   |   |   |   |   |   HouseOwnerFlag = 1 : Manual (27/1) [0/0]
9 |   |   |   |   |   |   |   HouseOwnerFlag = 0 : Manual (3/0) [0/0]
10 |   |   |   |   |   |   |   MaritalStatus = S : Manual (11/0) [0/0]
11 |   |   |   |   |   |   |   Gender = F
12 |   |   |   |   |   |   |   HouseOwnerFlag = 1
13 |   |   |   |   |   |   |   |   MaritalStatus = M : Manual (21/1) [0/0]
14 |   |   |   |   |   |   |   |   MaritalStatus = S : Manual (7/0) [0/0]
15 |   |   |   |   |   |   |   |   HouseOwnerFlag = 0 : Manual (8/1) [0/0]
```

FIGURE 2.3 – Début du REPTree

## 2.2 Clustering

Première observation pour le clustering : Plusieurs des méthodes proposées par Weka reposent sur d'autres méthodes.

Ainsi, `FilteredClusterer` et `MakeDensityBasedClusterer` demandent de choisir une méthode de clustering parmi celles mises à disposition par Weka. Petite particularité de Weka : Celui-ci n'empêche pas d'indiquer à (par exemple) `FilteredClusterer` d'utiliser `FilteredClusterer` pour fonctionner, permettant ainsi de "chaîner" indéfiniment ceux deux méthodes de clustering. Après plusieurs tests, ceci ne change rien aux valeurs de résultats (mais augmente le temps de génération de ceux-ci).

Deuxième observation : D'autres méthodes permettent de définir un algorithme de calcul de distance parmi ceux proposés par Weka.

Ainsi, `HierarchicalClusterer` et `SimpleKMeans` sont tous les deux configurable à ce niveau.

Dans notre cas nous avons choisi de détailler les résultats de la méthode `FarthestFirst` car sa simplicité permet d'en combiner l'usage avec les cluster `FilteredClusterer` et `MakeDensityBasedClusterer` et bien comprendre son fonctionnement facilite ainsi la compréhension d'autres méthodes que nous pourrions combiner avec lors de cas extra-académiques.

La méthode `FarthestFirst` est dite de clustering non-hiérarchique (ou Flat Clustering) car elle ne crée pas d'arbre, et donc pas de descendance / hiérarchie entre les cluster qu'elle génère. Deux paramètres numériques sont définissables :

- Le nombre de cluster à générer
- La graine à utiliser pour tirer aléatoirement des exemples

Contrairement à d'autres méthodes essayées et évaluées, `FarthestFirst` ne tire pas une graine aléatoirement quand elle est définie à zéro. Il est donc du devoir de l'utilisateur de Weka d'être attentif à cette particularité s'il souhaite lancer plusieurs générations afin de déterminer les meilleurs clusters possible.

Les options ne permettent pas de définir un nombre maximum d'itérations, on peut donc en déduire que la méthode évaluée repose sur la convergence des données dans les clusters pour déterminer quand s'arrêter.

Pour nos tests nous avons essayé avec les nombres de clusters suivants :

- 2 : C'est le nombre minimum de clusters assignable et cela permet de voir quels données sont dissociées
- 8 : C'est le nombre de type de données différentes que nous traitons dans cet exercice

### 2.2.1 2 clusters

Les listes déroulantes du haut de l'interface permettent de définir les éléments à afficher sur l'axe des X et ceux à afficher sur l'axe des Y pour la visualisation. Si les types de données sont discrètes, un paramètre (jitter) permet de visualiser plus facilement l'étendue des répartitions en ajoutant une forme de bruit en X et en Y aux graphes. Dans les exemples ci-dessous, les graphes sont faits avec le revenu annuel en abscisse et la possession ou non de maison en ordonnée. En rouge et en bleu on voit les deux clusters déterminés par le run, dont le résultat est également ci-joint.

## 2.2. Clustering

```
1 === Run information ===
2
3 Scheme:          weka.clusterers.FarthestFirst -N 2 -S 0
4 Relation:        QueryResult-weka.filters.unsupervised.attribute.Remove-R2
5                   ↪ -7,9,11,14,16,18,21-25-weka.filters.unsupervised.attribute.Remove-R1
6 Instances:       18484
7 Attributes:      8
8                   MaritalStatus
9                   Gender
10                  YearlyIncome
11                  TotalChildren
12                  EnglishEducation
13                  EnglishOccupation
14                  HouseOwnerFlag
15                  NumberCarsOwned
16 Test mode:       evaluate on training data
17
18 === Clustering model (full training set) ===
19
20
21 FarthestFirst
22 =====
23
24 Cluster centroids:
25
26 Cluster 0
27     M M 60000.0 2.0 Partial College Professional 1 1.0
28 Cluster 1
29     S F 170000.0 4.0 Bachelors Management 0 4.0
30
31
32
33 Time taken to build model (full training data) : 0.07 seconds
34
35 === Model and evaluation on training set ===
36
37 Clustered Instances
38
39 0      14166 ( 77%)
40 1      4318 ( 23%)
```

FIGURE 2.4 – Output de l'exécution avec 2 clusters

## 2.2. Clustering

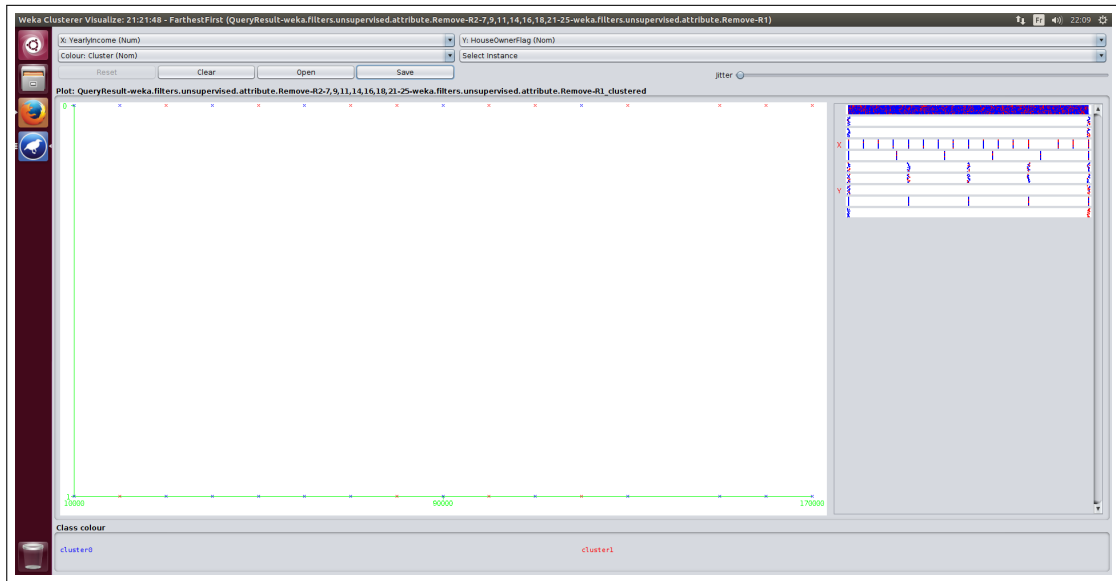


FIGURE 2.5 – Clustering sans jitter et avec deux clusters

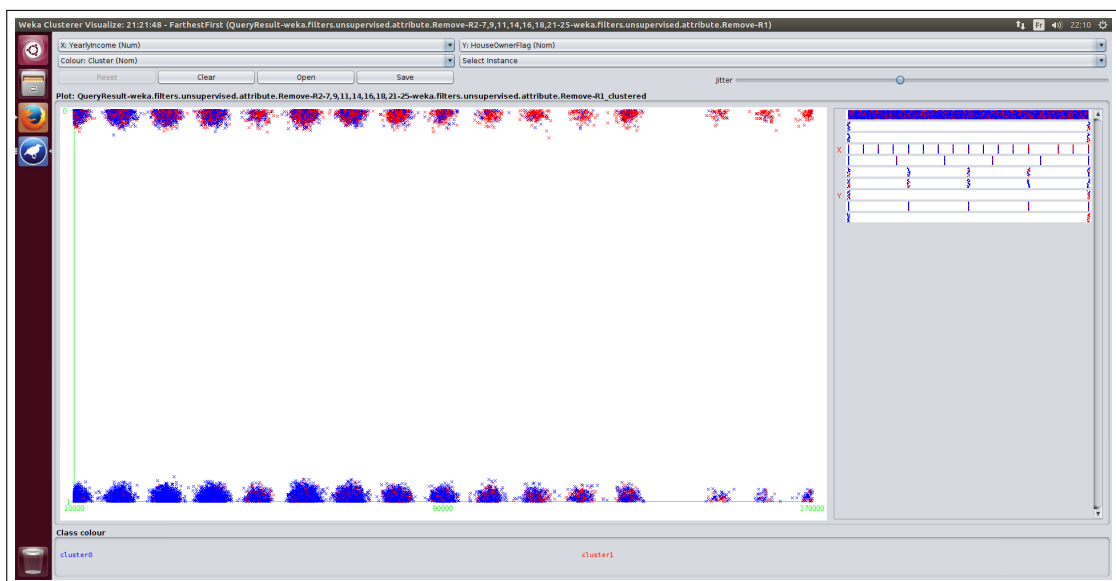


FIGURE 2.6 – Clustering avec jitter et avec deux clusters

### 2.2.2 8 clusters

Pour cette génération (le jitter a été volontairement poussé pour permettre une meilleure visualisation des résultats), on constate que les revenus influencent certains clusters quand à l'état civil (marié ou non). Il ne nous est pas en l'état possible de savoir pourquoi et de plus amples recherches seraient nécessaires pour déterminer les facteurs

communs de ces résultats.

```

1 === Run information ===
2
3 Scheme:          weka.clusterers.FarthestFirst -N 8 -S 1
4 Relation:        QueryResult-weka.filters.unsupervised.attribute.Remove-R2
5                   ↪ -7,9,11,14,16,18,21-25-weka.filters.unsupervised.attribute.Remove-R1
6 Instances:       18484
7 Attributes:      8
8                   MaritalStatus
9                   Gender
10                  YearlyIncome
11                  TotalChildren
12                  EnglishEducation
13                  EnglishOccupation
14                  HouseOwnerFlag
15                  NumberCarsOwned
16 Test mode:       evaluate on training data
17
18 === Clustering model (full training set) ===
19
20
21 FarthestFirst
22 =====
23
24 Cluster centroids:
25
26 Cluster 0
27     M M 30000.0 0.0 Partial College Skilled Manual 0 1.0
28 Cluster 1
29     S F 120000.0 5.0 Partial High School Professional 1 4.0
30 Cluster 2
31     S F 160000.0 0.0 Graduate Degree Management 0 1.0
32 Cluster 3
33     S M 10000.0 3.0 High School Manual 1 0.0
34 Cluster 4
35     M M 170000.0 1.0 Bachelors Management 1 4.0
36 Cluster 5
37     M F 30000.0 5.0 Graduate Degree Clerical 1 0.0
38 Cluster 6
39     M F 120000.0 1.0 High School Professional 0 4.0
40 Cluster 7
41     S M 130000.0 5.0 High School Management 0 4.0
42
43
44
45 Time taken to build model (full training data) : 0.07 seconds
46
47 === Model and evaluation on training set ===
48
49 Clustered Instances
50
51 0          3967 ( 21%)
52 1          1534 (  8%)
53 2          1825 ( 10%)
54 3          3031 ( 16%)
55 4          2374 ( 13%)
56 5          3405 ( 18%)
57 6          1736 (  9%)
58 7           612 (  3%)

```

FIGURE 2.7 – Output de l'exécution avec 8 clusters

## 2.2. Clustering

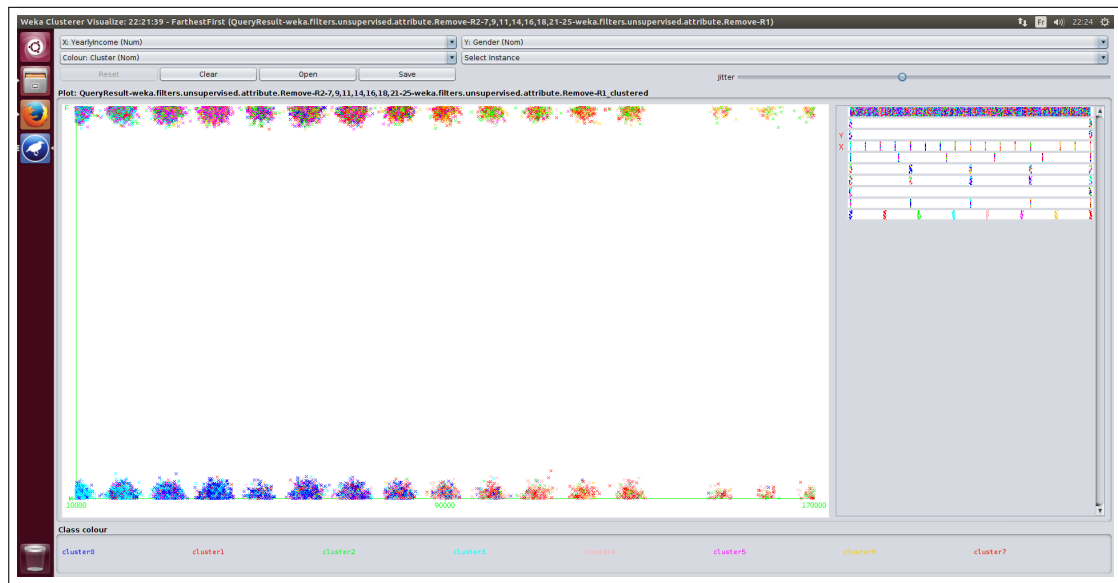


FIGURE 2.8 – Clustering avec jitter et avec huit clusters



## 2.3 Règles d'associations

Pour cette section, le filtre *NumericToNominal* a été appliqué aux attributs *NumberCarsOwned*, *YearlyIncome* et *TotalChildren*.

Seul les algorithmes *Apriori* et *FilteredAssociator* ont été testé. Seul ceux-ci étaient disponible (voir 2.9).

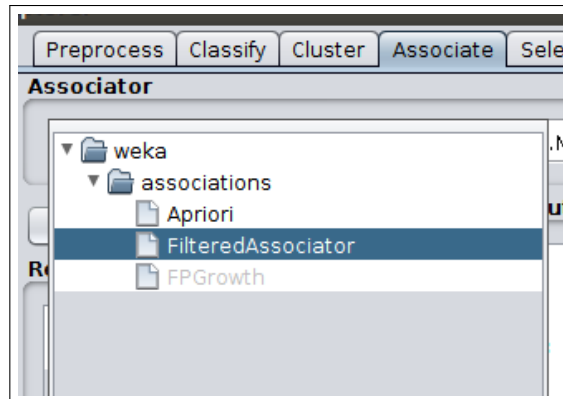


FIGURE 2.9 – Règles d'associations : Algos disponible

### 2.3.1 Apriori

Le seuil minimum de confiance a été modifié pour obtenir plus de règles.

Voici les paramètres :

```
Apriori -N 10 -T 0 -C 0.8 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c 1
```

Et le résultat (figure 2.10).

## 2.3. Règles d'associations

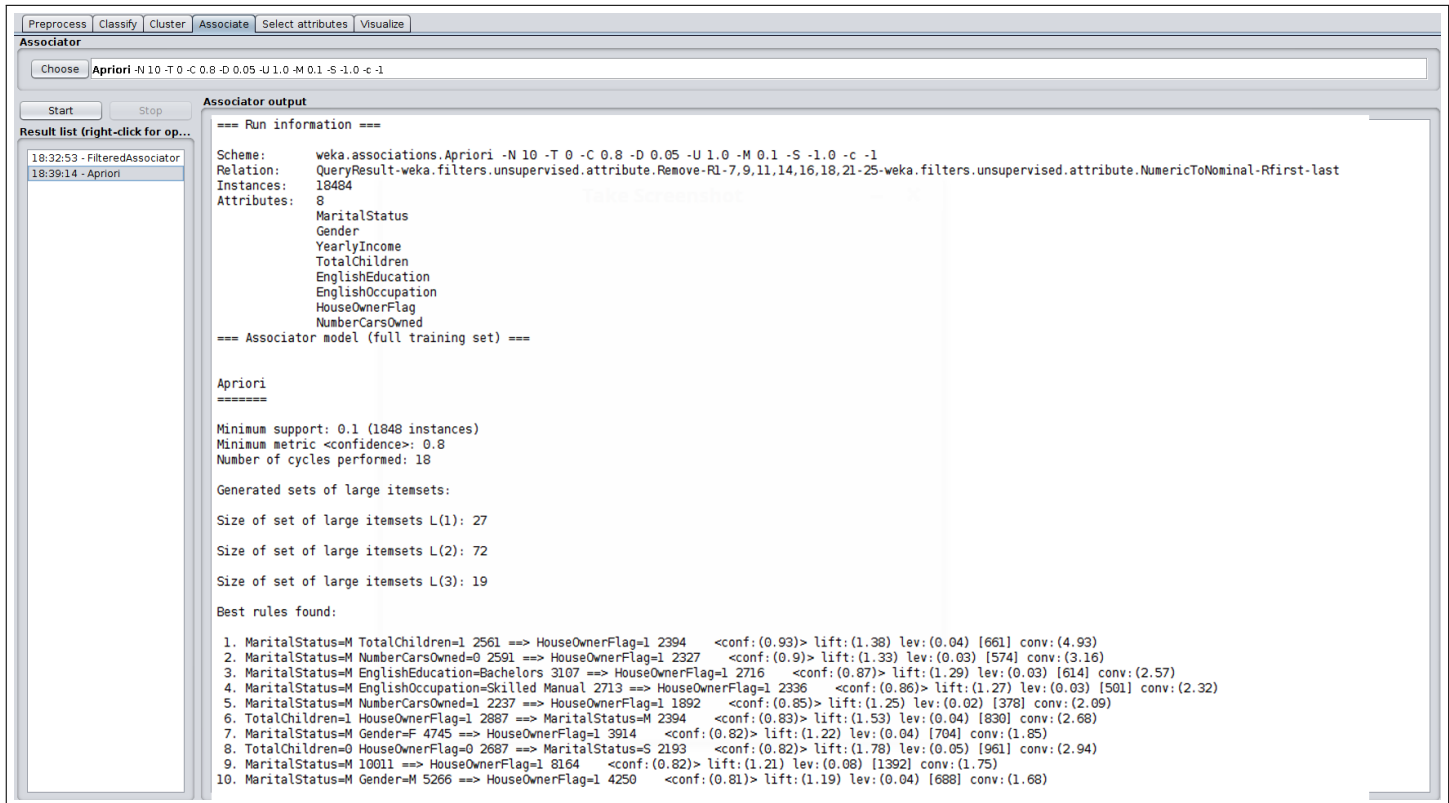


FIGURE 2.10 – Règles d'associations : Apriori résultat

### 2.3.2 FilteredAssociator

Cet algorithme permet de créer des règles d'association filtrées.

Les 10 règles obtenues (figure 2.11) sont identiques au résultat de l'*Apriori*.

## 2.3. Règles d'associations

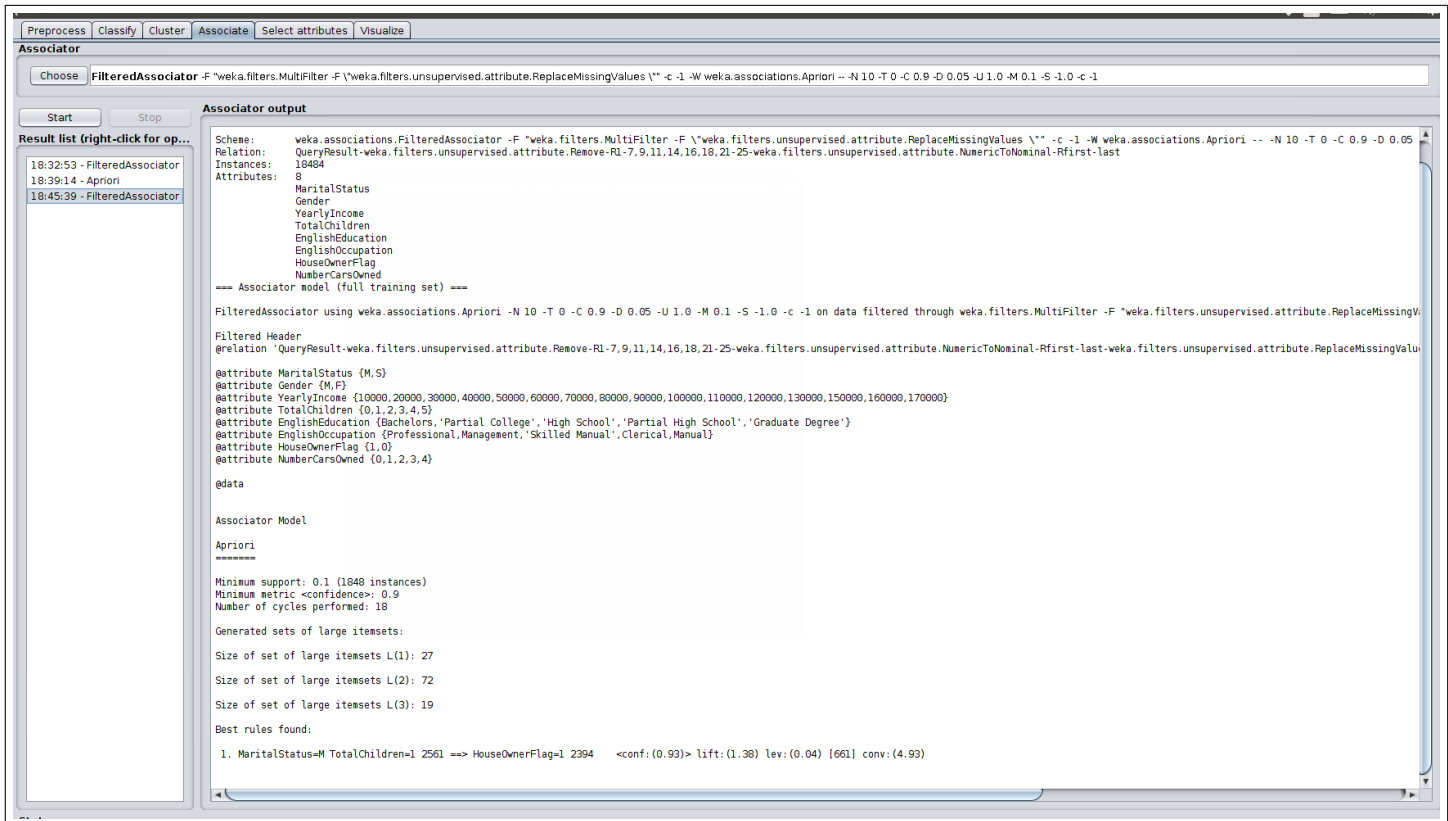


FIGURE 2.11 – Règles d'associations : FilteredAssociator résultat

## 3 | Conclusion

En faisant ce laboratoire, nous avons découvert la puissance du logiciel Weka. Malgré une interface qui pourrait être plus intuitive, la possibilité de se connecter directement à une base de donnée MySQL permet de faire des analyses très rapidement. Ainsi, Weka semble être un logiciel très utile lorsqu'il s'agit de faire une analyse de donnée préliminaire ou, lorsqu'une contrainte de coût et/ou délai est présente. Un exemple d'utilisation relativement réaliste serait celui où Weka est connecté à une base de donnée en production pour ajuster les différents public cible d'une gamme de produit.

En testant les différents types d'algorithmes que nous avons vu en cours, notre compréhension de ces derniers s'est améliorée par le biais d'un exemple pratique. Nous avons compris que certains algorithmes nécessitent un paramétrage particulier pour donner de bon résultats de classification. Désormais, notre vision critique vis-à-vis des algorithmes à disposition est plus claire.