

# How to use ClustersAnalysis Package

## Contents

<b>Introduction</b>	<b>1</b>
Short Descriptions of datasets . . . . .	1
Iris : . . . . .	1
Infert : . . . . .	2
Esoph : . . . . .	2
<b>Import ClustersAnalysis from Github (using devtools)</b>	<b>2</b>
How to access to help . . . . .	3
Univariate Analysis for qualitatives variariables . . . . .	3
Contingency table and size effect . . . . .	5

## Introduction

This is a demonstration of using the R package Clusters. You will see how to analyze classes according to one or more variables. The group variable must be of the type factor or character and the exploratory variables can be quantitative or qualitative. In this demonstration we are going to use natives dataset from R such as “iris”, “infert” or “esoph”.

## Short Descriptions of datasets

### Iris :

The data set consists of 50 samples from each of three species of Iris (Iris setosa, Iris virginica and Iris versicolor). Four features were measured from each sample: the length and the width of the sepals and petals, in centimeters. Based on the combination of these four features, Fisher developed a linear discriminant model to distinguish the species from each other.

```
summary(iris)
```

```
##   Sepal.Length   Sepal.Width   Petal.Length   Petal.Width
##   Min.    :4.300   Min.    :2.000   Min.    :1.000   Min.    :0.100
##   1st Qu.:5.100   1st Qu.:2.800   1st Qu.:1.600   1st Qu.:0.300
##   Median :5.800   Median :3.000   Median :4.350   Median :1.300
##   Mean    :5.843   Mean    :3.057   Mean    :3.758   Mean    :1.199
##   3rd Qu.:6.400   3rd Qu.:3.300   3rd Qu.:5.100   3rd Qu.:1.800
##   Max.    :7.900   Max.    :4.400   Max.    :6.900   Max.    :2.500
```

```
##      Species
## setosa      :50
## versicolor:50
## virginica   :50
##
##
##
```

**Infert :**

This is a matched case-control study dating from before the availability of conditional logistic regression.

```
summary(infert)
```

```
##      education      age      parity      induced
## 0-5yrs : 12   Min.   :21.00   Min.   :1.000   Min.   :0.0000
## 6-11yrs:120   1st Qu.:28.00   1st Qu.:1.000   1st Qu.:0.0000
## 12+ yrs:116   Median :31.00   Median :2.000   Median :0.0000
##              Mean    :31.50   Mean    :2.093   Mean    :0.5726
##              3rd Qu.:35.25   3rd Qu.:3.000   3rd Qu.:1.0000
##              Max.    :44.00   Max.    :6.000   Max.    :2.0000
##      case      spontaneous      stratum      pooled.stratum
## Min.   :0.0000   Min.   :0.0000   Min.   : 1.00   Min.   : 1.00
## 1st Qu.:0.0000   1st Qu.:0.0000   1st Qu.:21.00   1st Qu.:19.00
## Median :0.0000   Median :0.0000   Median :42.00   Median :36.00
## Mean    :0.3347   Mean    :0.5766   Mean    :41.87   Mean    :33.58
## 3rd Qu.:1.0000   3rd Qu.:1.0000   3rd Qu.:62.25   3rd Qu.:48.25
## Max.    :1.0000   Max.    :2.0000   Max.    :83.00   Max.    :63.00
```

**Esoph :**

Data from a case-control study of (o)esophageal cancer in Ille-et-Vilaine, France. This is a data frame with records for 88 age/alcohol/tobacco combinations.

```
summary(esoph)
```

```
##      agegp      alcgp      tobgp      ncases      ncontrols
## 25-34:15  0-39g/day:23  0-9g/day:24   Min.   : 0.000   Min.   : 1.00
## 35-44:15  40-79      :23  10-19      :24   1st Qu.: 0.000   1st Qu.: 3.00
## 45-54:16  80-119     :21  20-29      :20   Median : 1.000   Median : 6.00
## 55-64:16  120+       :21  30+       :20   Mean    : 2.273   Mean    :11.08
## 65-74:15                      3rd Qu.: 4.000   3rd Qu.:14.00
## 75+      :11                      Max.    :17.000   Max.    :60.00
```

## Import ClustersAnalysis from Github (using devtools)

```
#Use the below line to install devtools if necessary
#install.packages("devtools")
#library(devtools)
```

```
#install package from github
#Use the below line to install ClustersAnalysis if necessary
#devtools::install_github("clepadellec/ClustersAnalysis")

#load package
library(ClustersAnalysis)
```

## How to access to help

You can just use the function `help(function name)` to see all the documentation about your function.

```
help("u_plot_size_effect")
```

## Univariate Analysis for qualitative variables

To begin we will try to understand, for each qualitative explanatory variable, if it affects the group variable. It's necessary to create an object of univariate type. You can use the constructor **Univariate\_object**.

```
#Creation of univariate object using esoph dataframe and "agegp" (first column) as the group variable
u_esoph<-Univariate_object(esoph,1)
```

```
#detail of attributes associated with the object
print(u_esoph)
```

```
## $ind.qual
##      agegp      alcgp      tobgp      ncases ncontrols
##      TRUE       TRUE       TRUE       FALSE      FALSE
##
## $ind.quan
##      agegp      alcgp      tobgp      ncases ncontrols
##      FALSE      FALSE      FALSE       TRUE       TRUE
##
## $df
##      agegp      alcgp      tobgp      ncases ncontrols
## 1  25-34  0-39g/day  0-9g/day      0         40
## 2  25-34  0-39g/day  10-19      0         10
## 3  25-34  0-39g/day  20-29      0          6
## 4  25-34  0-39g/day  30+        0          5
## 5  25-34      40-79  0-9g/day      0         27
## 6  25-34      40-79  10-19      0          7
## 7  25-34      40-79  20-29      0          4
## 8  25-34      40-79  30+        0          7
## 9  25-34      80-119 0-9g/day      0          2
## 10 25-34      80-119 10-19      0          1
## 11 25-34      80-119 30+        0          2
## 12 25-34      120+  0-9g/day      0          1
## 13 25-34      120+  10-19      1          1
## 14 25-34      120+  20-29      0          1
## 15 25-34      120+  30+        0          2
## 16 35-44  0-39g/day  0-9g/day      0         60
```

## 17	35-44	0-39g/day	10-19	1	14
## 18	35-44	0-39g/day	20-29	0	7
## 19	35-44	0-39g/day	30+	0	8
## 20	35-44	40-79	0-9g/day	0	35
## 21	35-44	40-79	10-19	3	23
## 22	35-44	40-79	20-29	1	14
## 23	35-44	40-79	30+	0	8
## 24	35-44	80-119	0-9g/day	0	11
## 25	35-44	80-119	10-19	0	6
## 26	35-44	80-119	20-29	0	2
## 27	35-44	80-119	30+	0	1
## 28	35-44	120+	0-9g/day	2	3
## 29	35-44	120+	10-19	0	3
## 30	35-44	120+	20-29	2	4
## 31	45-54	0-39g/day	0-9g/day	1	46
## 32	45-54	0-39g/day	10-19	0	18
## 33	45-54	0-39g/day	20-29	0	10
## 34	45-54	0-39g/day	30+	0	4
## 35	45-54	40-79	0-9g/day	6	38
## 36	45-54	40-79	10-19	4	21
## 37	45-54	40-79	20-29	5	15
## 38	45-54	40-79	30+	5	7
## 39	45-54	80-119	0-9g/day	3	16
## 40	45-54	80-119	10-19	6	14
## 41	45-54	80-119	20-29	1	5
## 42	45-54	80-119	30+	2	4
## 43	45-54	120+	0-9g/day	4	4
## 44	45-54	120+	10-19	3	4
## 45	45-54	120+	20-29	2	3
## 46	45-54	120+	30+	4	4
## 47	55-64	0-39g/day	0-9g/day	2	49
## 48	55-64	0-39g/day	10-19	3	22
## 49	55-64	0-39g/day	20-29	3	12
## 50	55-64	0-39g/day	30+	4	6
## 51	55-64	40-79	0-9g/day	9	40
## 52	55-64	40-79	10-19	6	21
## 53	55-64	40-79	20-29	4	17
## 54	55-64	40-79	30+	3	6
## 55	55-64	80-119	0-9g/day	9	18
## 56	55-64	80-119	10-19	8	15
## 57	55-64	80-119	20-29	3	6
## 58	55-64	80-119	30+	4	4
## 59	55-64	120+	0-9g/day	5	10
## 60	55-64	120+	10-19	6	7
## 61	55-64	120+	20-29	2	3
## 62	55-64	120+	30+	5	6
## 63	65-74	0-39g/day	0-9g/day	5	48
## 64	65-74	0-39g/day	10-19	4	14
## 65	65-74	0-39g/day	20-29	2	7
## 66	65-74	0-39g/day	30+	0	2
## 67	65-74	40-79	0-9g/day	17	34
## 68	65-74	40-79	10-19	3	10
## 69	65-74	40-79	20-29	5	9
## 70	65-74	80-119	0-9g/day	6	13

```

## 71 65-74      80-119      10-19      4      12
## 72 65-74      80-119      20-29      2      3
## 73 65-74      80-119      30+      1      1
## 74 65-74      120+ 0-9g/day      3      4
## 75 65-74      120+      10-19      1      2
## 76 65-74      120+      20-29      1      1
## 77 65-74      120+      30+      1      1
## 78 75+ 0-39g/day 0-9g/day      1      18
## 79 75+ 0-39g/day      10-19      2      6
## 80 75+ 0-39g/day      30+      1      3
## 81 75+      40-79 0-9g/day      2      5
## 82 75+      40-79      10-19      1      3
## 83 75+      40-79      20-29      0      3
## 84 75+      40-79      30+      1      1
## 85 75+      80-119 0-9g/day      1      1
## 86 75+      80-119      10-19      1      1
## 87 75+      120+ 0-9g/day      2      2
## 88 75+      120+      10-19      1      1
##
## $group
## [1] 1
##
## $name_group
## [1] "agegp"
##
## $lst_quali
## [1] "agegp" "alcgp" "tobgp"
##
## $lst_quanti
## [1] "ncases"      "ncontrols"
##
## $multiple_var
## [1] TRUE

```

## Contingency table and size effect

The first thing we can do is to create a contingency table between the explanatory variable and the group variable and then visualize lines/columns profiles. In our case the explanatory variable is “tobgp” which is the tobacco consumption (gm/day). To do this you can use the `u_desc_profiles`

```

#use interact=TRUE to show an interactive graphique with widgets like zoom, comparisons...
ClustersAnalysis::u_desc_profiles(u_esoph,3,interact=FALSE)

```

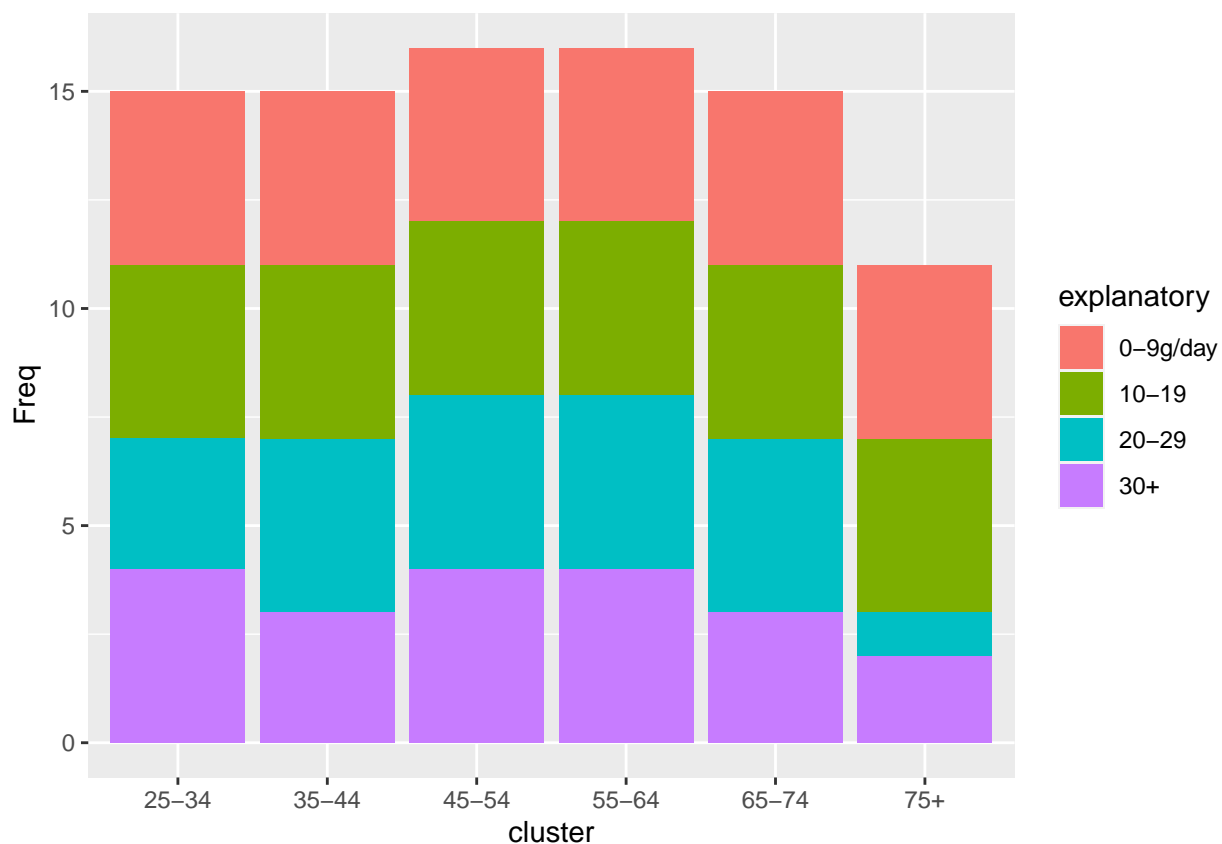
```

## [1] "Tableau de contingence : "
##
##           0-9g/day 10-19 20-29 30+
## 25-34           4      4      3      4
## 35-44           4      4      4      3
## 45-54           4      4      4      4
## 55-64           4      4      4      4
## 65-74           4      4      4      3
## 75+             4      4      1      2
## [1] "Profils lignes : "

```

```
##
##           0-9g/day 10-19 20-29 30+   Total
## 25-34      26.7      26.7  20.0  26.7 100.0
## 35-44      26.7      26.7  26.7  20.0 100.0
## 45-54      25.0      25.0  25.0  25.0 100.0
## 55-64      25.0      25.0  25.0  25.0 100.0
## 65-74      26.7      26.7  26.7  20.0 100.0
## 75+        36.4      36.4   9.1  18.2 100.0
## Ensemble  27.3      27.3  22.7  22.7 100.0
## [1] "Profils colonnes : "
```

```
##
##           0-9g/day 10-19 20-29 30+   Ensemble
## 25-34  16.67      16.67  15.00  20.00  17.05
## 35-44  16.67      16.67  20.00  15.00  17.05
## 45-54  16.67      16.67  20.00  20.00  18.18
## 55-64  16.67      16.67  20.00  20.00  18.18
## 65-74  16.67      16.67  20.00  15.00  17.05
## 75+    16.67      16.67   5.00  10.00  12.50
## Total 100.00     100.00 100.00 100.00 100.00
```



The distributions don't show any particular phenomena. The most represented classes are the 35-44 and 55-64 years. Then we can see that there are more than a half that smokes less than 20 g/days. The only fact that we can see is that the 75+ people are close to 75% to don't smoke a lot.

Now we are going to see in details if there is a size effect between these two variables. To do this we can use `u_desc_size_effect` which return the test statistic `vt` (comparison between proportions). Then we can also use `u_plot_size_effect` which create a mosaic plot.

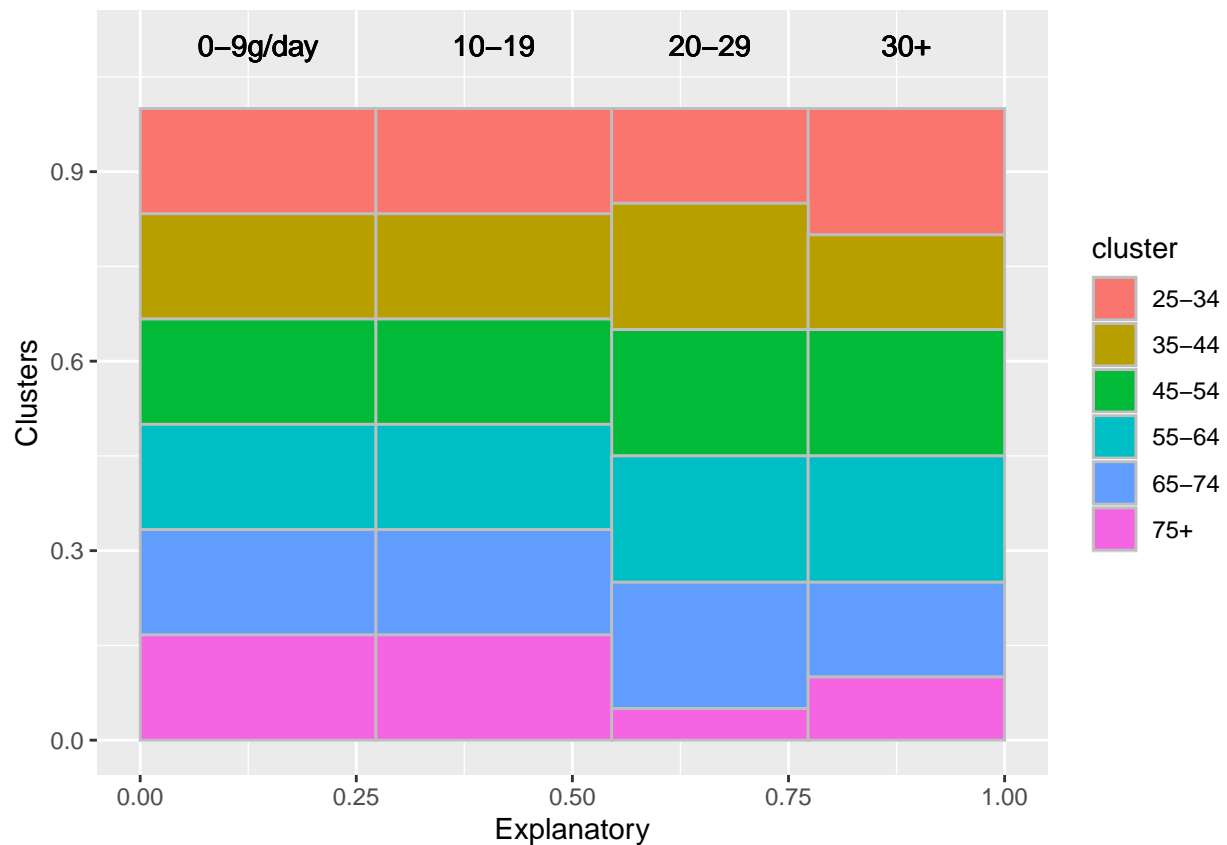
```
u_desc_size_effect(u_esoph,3)
```

```
##
##      0-9g/day 10-19 20-29 30+ Sum
## 25-34      4      4      3      4 15
## 35-44      4      4      4      3 15
## 45-54      4      4      4      4 16
## 55-64      4      4      4      4 16
## 65-74      4      4      4      3 15
## 75+       4      4      1      2 11
## Sum       24     24     20     20 88
```

```
##
##      0-9g/day      10-19      20-29      30+
## 25-34 -0.1291915 -0.1291915 -0.6565969 0.9484178
## 35-44 -0.1291915 -0.1291915 0.9484178 -0.6565969
## 45-54 -0.5038193 -0.5038193 0.5690195 0.5690195
## 55-64 -0.5038193 -0.5038193 0.5690195 0.5690195
## 65-74 -0.1291915 -0.1291915 0.9484178 -0.6565969
## 75+    1.6158165  1.6158165 -2.7373833 -0.9124611
```

```
#use interact=TRUE to show an interactive graphique with widgets like zoom, comparisons...
u_plot_size_effect(u_esoph,3,interact=FALSE)
```

```
## Warning: Ignoring unknown aesthetics: width
```



If we refer to the results we can see that the biggest “vt” values are for 75+ peoples who smokes 0-9g/day or 10-19 g/day. So this is for these two combinations that one can most easily conclude that there is an over-representation. We can use the mosaic plot to confirm our comment.