#### Initiation to R software Session VI

Pierre Michel

Master AMSE 1st year, 2019

### Statistics in R

Introduction: statistics in R

#### Functions available for statistical analysis:

- Package stats: functions for classical statistical analysis (hypotheses tests, linear models, distributions, statistical summary, time series, multivariate analyses). It is loaded when R starts. USe library(help = "stats")
- Specific packages: other packages available for statistical methods, distributed with R or downloadable: class (for classification), boot (for bootstrap), survival (for survival analysis)...

## Introduction: statistical modelling

#### Functions from stats for statistical modelling:

function_name	description
SSlogis	Logistic model
aov	Analysis of variance
glm	Generalized linear models
lm	Linear models
manova	Mutivariate analysis of variance
ksmooth	Kernel smoother regression
step	Model selection with Stepwise algorithm

#### Introduction: statistical tests

#### Functions from stats for statistical tests:

function_name	description
bartlett.test	Bartlett's variances homogeneity test
binom.test	Exact binomial test
chisq.test	Chi2 test
cor.test	Correlation test
t.test	Student's means comparison t-Test
var.test	Fisher's variances comparison F-test
wilcox.test	Wilcoxon's rank test
ks.test	Kolmogorov-Smirnov's test

### Introduction: estimation and data analysis

#### Functions for estimation:

function_name	description
density ecdf	Kernel density estimation Empirical cumulative distribution function

#### Functions for data analysis:

function_name	description		
hclust	Hierarchical clustering		
kmeans	Partitioning clustering		
princomp	Principal components analysis		
dist	Computes distance matrices		

### Formulas

Formulas: syntax

Many statistical functions use formulas: function(formula, data, ...)

- formula of type response ~ predictors
- data data table containing the variables in formula

response is the target variable, predictors is the set of predictor variables, separated by arithmetic symbols.

#### Formulas: examples

- y~a: predictor a
- y~a+b: predictors a and b
- ▶ y~M: as many models as predictors in M (matrix)
- ▶ y~x-1: model without intercept
- y~log(b): predictor log(b)
- y~a+I(b+c): predictors a and b+c
- y~a:b: interaction of a and b
- y~a\*b:a+b+a:b: main effects and interaction of a and b
- y~(a+b)^2: main effects a and b and second-order interaction of a+b+a:b
- y~a\*b-a:b: a+b

### Formulas: examples with formulas

Example 1: Simple linear regression of y on x (x and y quantitative).

```
x=sample(1:10,200,replace=TRUE)
y=3+7*x+rnorm(200,0,100)
linreg=lm(y~x)
```

Example 2: Multivariate linear regression of fertility on education and infant mortality (all variables quantitative).

```
data(swiss)
mlinreg = lm(Fertility~Education+Infant.Mortality,swiss)
```

Example 3: 1-factor analysis of variance: analyze the effect of a bug spray (6 different types) on the number of insects (count), based on the observation of 12 cultures sprayed succesively with the bug sprays. The response is qualitative.

```
data(InsectSprays)
anova = aov(sqrt(count) ~ spray, data = InsectSprays)
```

#### Formulas: examples without formulas

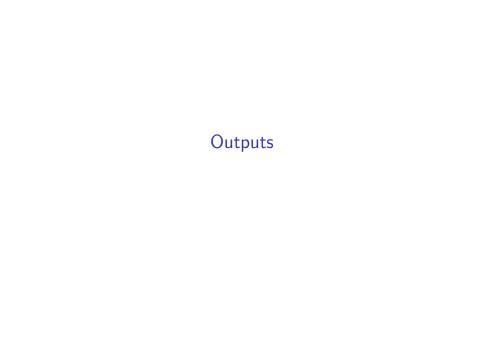
Example 4:  $\chi^2$ -test to test the relation between qualitative variables X and Y, based on a sample in the form of a contengency table.

```
0 = matrix(c(442,514,38,6),nrow=2,byrow=TRUE)
colnames(0) = c("male","female");
rownames(0) = c("sighted","blind")
X2 = chisq.test(0,correct=FALSE)
```

Example 5: Student's *t*-test to test the equality between two means, based on a sample of each subpopulation  $X = (X_1, ..., X_{n_1})$  and  $Y = (Y_1, ..., Y_{n_2})$ 

```
x = rnorm(100,1,1); y = rexp(200,1)
st = t.test(x,y)

x = rnorm(100,1,1); z = rep(c(T,F),50)
st2 = t.test(x~z)
```



# Outputs: Linear regression (lm())

```
##
## Call:
## lm(formula = y \sim x)
##
   Coefficients:
## (Intercept)
```

29.610

X

4.592

linreg

##

# Outputs: Multivariate linear regression (lm())

```
##
## Call:
## lm(formula = Fertility ~ Education + Infant.Mortality, of
##
## Coefficients:
## (Intercept) Education Infant.Mortality
## 48.8213 -0.8167 1.5187
```

# Outputs: Analysis of variance (aov())

```
anova
## Call:
      aov(formula = sqrt(count) ~ spray, data = InsectSpray
##
##
## Terms:
##
                      spray Residuals
## Sum of Squares 88.43787 26.05798
## Deg. of Freedom
                          5
                                   66
##
## Residual standard error: 0.6283453
## Estimated effects may be unbalanced
```

# Outputs: $\chi^2$ -test (chisq.test())

```
%
##
## Pearson's Chi-squared test
##
## data: 0
## X-squared = 27.139, df = 1, p-value = 1.894e-07
```

# Outputs: Student's *t*-test (t.test())

```
st
##
##
    Welch Two Sample t-test
##
## data: x and y
## t = -0.16849, df = 197.7, p-value = 0.8664
## alternative hypothesis: true difference in means is not
## 95 percent confidence interval:
## -0.2606861 0.2196456
## sample estimates:
## mean of x mean of y
## 0.9942497 1.0147699
```

# Outputs: Student's *t*-test (t.test())

```
st2
##
##
   Welch Two Sample t-test
##
## data: x by z
## t = -0.77103, df = 92.686, p-value = 0.4427
## alternative hypothesis: true difference in means is not
## 95 percent confidence interval:
## -0.5958588 0.2625704
## sample estimates:
## mean in group FALSE mean in group TRUE
             0.8617696
##
                                 1.0284138
```

#### Outputs

In R, statistical functions return an object of class inherited from the name of the function (e.g aov() return an object of class aov, lm() returns an object of class lm, tests return objects of class lm, tests lm, lm

The returned object contains the results of the analysis. This is generally a list, whose displayed results depend on its class (e.g aov() and lm() do not provide the same lists).

The elements in the list-object can be printed using attributes() and names().

## Outputs: Linear regression and ANOVA

```
names(linreg)
    [1] "coefficients"
                         "residuals"
                                          "effects"
##
##
    [5] "fitted.values" "assign"
                                           "qr"
##
    [9] "xlevels"
                         "call"
                                           "terms"
                                                            "mo
names (anova)
##
    [1] "coefficients"
                         "residuals"
                                          "effects"
    [5] "fitted.values" "assign"
                                          "qr"
##
                         "xlevels"
##
    [9] "contrasts"
                                          "call"
                                                            "te
##
   [13] "model"
```

### Outputs: $\chi^2$ -test

```
attributes(X2)
```

"da

# Outputs: Simple/Multivariate linear regression

(Intercent)

118.3705958

146.3313526

## -155.6741978

16

21

26

##

##

##

##

##

##

##

coefficients(mlinreg) # or mlinreg\$coefficients

##	(Interce	pt) Ec	iucation iniai	ic.Morcaricy	
##	48.8212	736 -0.	8166573	1.5187190	
res	<pre>iduals(linreg</pre>	) # or linreg	g\$residuals		
##	1	2	3	4	
##	-197.9339498	-95.3936445	-43.2006639	-61.2268520	43
##	6	7	8	9	
##	-90.4280540	-35.8655301	19.5964264	-140.8449593	-7
##	11	12	13	14	

17

22

27

-20.4186140 -177.7767247

134.3449689 -19.6699709

65.1376230 -115.2425095 -133.4402101

33.5042972 -84.9596497

Education Infant Mortality

18

23

28

-65.3849881

160.9829140

127.5716138

19

24

29

5

116.0808205 -103

# Outputs: $\chi^2$ -test

```
X2$expected # theoretical counts
## male female
## sighted 458.88 497.12
## blind 21.12 22.88
X2$residuals # residuals (theoretical - observed)
##
               male female
## sighted -0.7879939 0.7570801
## blind 3.6730385 -3.5289413
sum(X2$residuals^2) # Khi2 statistic value
## [1] 27.13874
```

# Generic functions

#### Generic functions

Some functions are also used to extract the desired results: the **generic functions**.

They work specifically, according to the object class.

Generic functions have a **single syntax** for all cases.

#### Generic functions

- print() returns a short summary of the analysis.
- summary() returns a detailed summary of the analysis.
- df.residuals() returns the number of degrees of freedom of residual.
- coef() returns the estimated coefficients (sometimes with standard errors).
- residuals() returns the residuals.
- fitted() returns values predicted by the model.
- logLik() computes the log-likelihood and the number of parameters of a model.
- AIC() computes Akaike's information criterion.
- anova() table of analysis of variance.
- plot() returns a plot adapted to the analysis.

# Generic functions: summary()

summary() prints a detailed summary of the analysis, specific to the object class.

#### apropos("^summary")

```
##
    [1] "summary"
                                    "Summary"
##
    [3]
       "summary.aov"
                                    "summary.connection"
##
    [5]
       "summary.data.frame"
                                    "Summary.data.frame"
    [7]
       "summary.Date"
                                    "Summary.Date"
##
       "summary.default"
                                    "Summary.difftime"
##
   [11] "summary.factor"
                                    "Summary.factor"
   [13] "summary.glm"
                                    "summary.lm"
   [15] "summary.manova"
##
                                    "summary.matrix"
##
   [17] "Summary.numeric_version"
                                   "Summary.ordered"
   [19]
                                    "Summary.POSIXct"
       "summary.POSIXct"
   [21]
       "summary.POSIX1t"
                                    "Summary.POSIX1t"
   [23]
       "summary.proc time"
                                    "summary.srcfile"
                                    "summary.stepfun"
   [25]
       "summary.srcref"
                                    "summary.warnings"
   [27] "summary.table"
```

# Generic functions: summary() summary(linreg)

##

## Call:

```
## lm(formula = y \sim x)
##
## Residuals:
      Min 1Q Median
                             3Q
                                   Max
##
## -332.41 -69.66 2.23 66.15 287.10
##
## Coefficients:
             Estimate Std. Error t value Pr(>|t|)
##
## (Intercept) 29.610 16.083 1.841 0.0671.
## x
              4.592 2.448 1.876 0.0621 .
## Signif. codes: 0 '***' 0.001 '**' 0.05 '.' 0.3
##
## Residual standard error: 99.43 on 198 degrees of freedom
```

# Generic functions: summary()

```
Df Sum Sq Mean Sq F value Pr(>F)
##
## spray 5 88.44 17.688 44.8 <2e-16 ***
```

## Residuals 66 26.06 0.395

## ---

summary(anova)

## Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.05 '.' 0.3

## Generic functions: plot()

plot() returns graphics adapted to the current analysis.

```
apropos("^plot")
```

```
## [1] "plot" "plot.default" "p
```

### Generic functions: plot()

par(mfrow = c(2,2))
plot(anova)

