
Table of Contents

- 1 Import data
- 2 Compute sum scores
- 3 Graphics
- 4 Esport graphs
- 5 R for statistical computing
- 6 Classical hypothesis testing in R
- 7 Generalized Linear Models (GLMs)

Table of Contents

- 1 Import data
- 2 Compute sum scores
- 3 Graphics
- 4 Esport graphs
- 5 R for statistical computing
- 6 Classical hypothesis testing in R
- 7 Generalized Linear Models (GLMs)

Conditioning according to variable labels

They must present a regular expression (common root):

```
colnames(benessere)
```

```
[1] "ID" "età" "genere" "frat" "item1" "item2"  
"item3" "item4"  
[9] "item5" "au1" "au2" "au3" "au4" "au5" "au6"  
"au7"  
[17] "au8" "au9" "au10"
```

Columns with `item` → well-being items

Columns with `au` → self-esteem items

10 / 102


```
summary(babies)
```

id	genere	peso	al
Length:10	Length:10	Min. : 5.411	Min.
Class :character	Class :character	1st Qu.: 7.809	1st Q
Mode :character	Mode :character	Median : 9.429	Media
		Mean : 9.970	Mean
		3rd Qu.:11.268	3rd Q
		Max. :17.343	Max.

table()

```
table(babies$genere)
```

f m

6 4

```

      genere
new_benessere  f  m
      alto   22 18
      basso  32 28

```


Aggregate a response variable according to grouping variable(s) (e.g., averaging per experimental conditions):

```
# one dependent variable (y) and single grouping variable
aggregate(y ~ x, data = data, FUN, ...)
```

```
# Multiple response variables, multiple grouping variables
aggregate(cbind(y1, y2) ~ x1 + x2, data = data, FUN, ...)
```

Aggregating: Examples

```
benessere = read.csv("data/benessereScores.csv", header = T,  
head(benessere)
```

	ID	età	genere	frat	item1	item2	item3	item4	item5	au1	au2	au3
1	1	16	1	0	1	2	4	3	4	5	4	5
2	2	21	1	1	2	2	3	4	3	5	4	5
3	3	28	2	4	2	3	5	1	2	4	4	5
4	4	15	2	2	3	4	2	2	2	4	5	5
5	5	23	1	3	4	5	1	5	2	3	2	5
6	6	31	1	2	2	3	2	1	2	5	1	5

	au8	au9	au10	score_ben	score_au
1	4	2	4	14	33
2	5	1	5	14	40
3	4	2	4	13	38
4	4	2	4	13	38
5	4	1	3	17	26

```
aggregate(score_au ~ genere, data = benessere, mean)
```

	genere	score_au
1	1	33.08750
2	2	32.62857

Compute the mean according of self esteem and well being according to gender

```
aggregate(cbind(score_ben, score_au) ~ genere, data = benesse
```

	genere	score_ben	score_au
1	1	14.46250	33.08750
2	2	14.41429	32.62857

Your turn!



- Recode `frat` and assign it to a new var into the data frame (`siblings: > 0 siblings → no > > 1+ sibilings → yes`)
- Compute the mean of `score_ben` according to `siblings`
- Compute the mean of `score_ben` and `score_au` according to `sibilingsand gender` (assign it to the object `mean_dep`)
- Compute the standard deviation of `score_ben` and `score_au` according to `sibilingsand gender` (assign it to the object `sd_dep`)
- merge `mean_dep` and `sd_dep` and assign the resulting object to `descr`

WARNING!

When using `merge` the column names must be different

Result

```
descr
```

	siblings	genere	mean_score_ben	mean_score_au	sd_score_ben	sd_score_au
1	no	1	13.90909	33.27273	3.250042	4.452734
2	no	2	14.51852	32.11111	3.309315	3.004270
3	yes	1	14.67241	33.01724	3.347625	3.743961
4	yes	2	14.34884	32.95349	3.228472	4.396717

```
benessere$siblings = ifelse(benessere$frat == 0, "no", "yes")

mean_dep = aggregate(cbind(score_ben, score_au) ~ siblings + genere,
                      data = benessere,
                      mean)

colnames(mean_dep)[3:4] = paste("mean",
                                colnames(mean_dep)[3:4],
                                sep = "_")

sd_dep = aggregate(cbind(score_ben, score_au) ~ siblings + genere,
                   data = benessere,
                   sd)

colnames(sd_dep)[3:4] = paste("sd", colnames(sd_dep)[3:4],
                              sep = "_")

descr = merge(mean_dep, sd_dep)
```

tidyverse()

```
install.packages("tidyverse")
library(tidyverse)
```


Descriptive statistica

```
benessere %>% # object
  group_by(siblings, genere) %>% # groupings
  summarise(m_benessere = mean(score_ben), # functions
            sd_benessere = sd(score_ben),
            m_au = mean(score_au),
            sd_au = sd(score_au))
```

A tibble: 4 x 6

Groups: siblings [2]

	siblings	genere	m_benessere	sd_benessere	m_au	sd_au
	<chr>	<int>	<dbl>	<dbl>	<dbl>	<dbl>
1	no	1	13.9	3.25	33.3	4.45
2	no	2	14.5	3.31	32.1	3.00
3	yes	1	14.7	3.35	33.0	3.74
4	yes	2	14.3	3.23	33.0	4.40



- Compute minimum, maximum, median of `score_au` and `score_ben` with `tidyverse`
- Import the `babies` data set and compute the descriptive stats of `weight` and `height` with `tidyverse`

Table of Contents

- 1 Import data
- 2 Compute sum scores
- 3 Graphics**
- 4 Esport graphs
- 5 R for statistical computing
- 6 Classical hypothesis testing in R
- 7 Generalized Linear Models (GLMs)

- Traditional graphics
- Grid graphics & ggplot2

For both:

- High level functions → actually produce the plot
- Low level functions → make it looks better =>

Traditional graphics I

High level functions

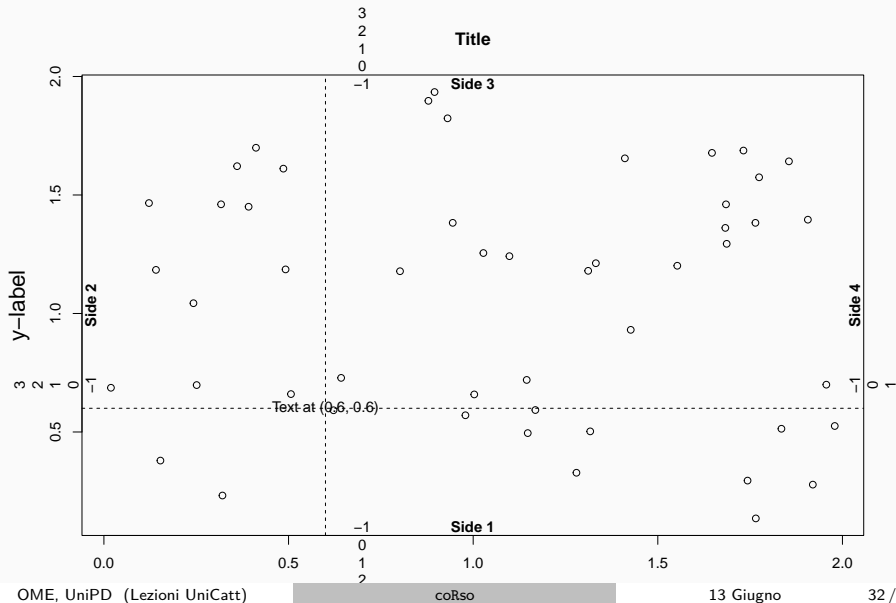
```
plot()           # scatter plot, specialized plot methods
boxplot()
hist()          # histogram
qqnorm()        # quantile-quantile plot
barplot()
pie()           # pie chart
pairs()          # scatter plot matrix
persp()         # 3d plot
contour()        # contour plot
coplot()         # conditional plot
interaction.plot()
```

demo(graphics) for a guided tour of base graphics!


```
x <- runif(50, 0, 2) # 50 uniform random numbers
y <- runif(50, 0, 2)
plot(x, y, main="Title",
      sub="Subtitle", xlab="x-label",
      ylab="y-label") # produce plotting window
```

```
text(0.6, 0.6, "Text at (0.6, 0.6)")
abline(h=.6, v=.6, lty=2) # horizont. and vertic.
                           # lines
```

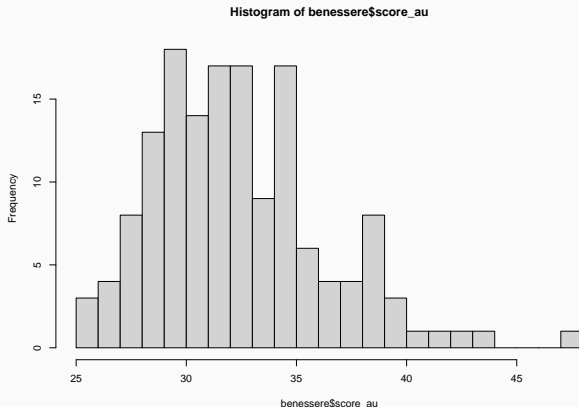

Margins region




```
plot(x) # solo una variabile
plot(x, y) # due variabile (scatter plot)
plot(y ~ x) # due variabile, y in funzione di x
```


hist(): Frequenze

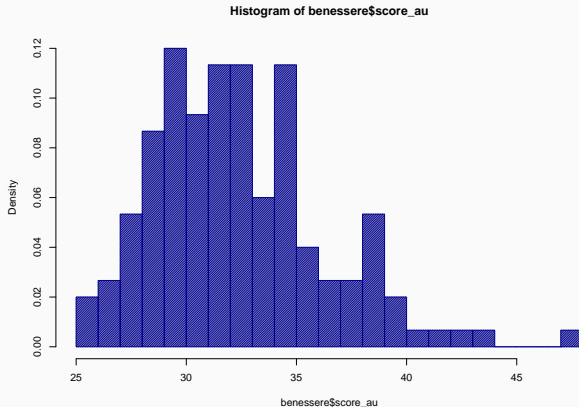
```
hist(benessere$score_au, breaks = 20)
```



hist(): Densità

Densità

```
hist(benessere$score_au,
     density=50, breaks=20,
     prob=TRUE, col = "darkblue")
```

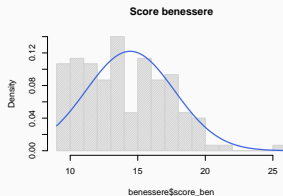
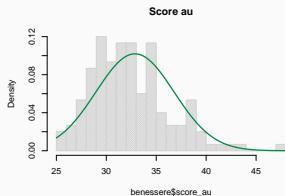


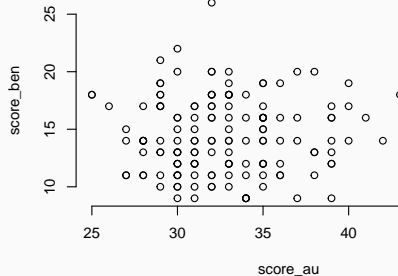
Multi plot (in riga)

```
par(mfrow=c(1, 2))

hist(benessere$score_au,density=50, breaks=20, prob=TRUE,
     main = "Score au")
curve(dnorm(x, mean=mean(benessere$score_au),
                  sd=sd(benessere$score_au)),
      col="springgreen4", lwd=2, add=TRUE, yaxt="n")

[...]
```





Multiplot (in colonna), codice

```
par(mfcol= c(2,2))
hist(benessere$score_au,density=50, breaks=20, prob=TRUE,
     main = "Score au")
curve(dnorm(x, mean=mean(benessere$score_au),
                    sd=sd(benessere$score_au)),
      col="springgreen4", lwd=2, add=TRUE, yaxt="n")

hist(benessere$score_ben,density=50, breaks=20, prob=TRUE,
     main = "Score benessere")
curve(dnorm(x, mean=mean(benessere$score_ben),
                    sd=sd(benessere$score_ben)),
      col="royalblue", lwd=2, add=TRUE, yaxt="n")

with(benessere,
     plot(score_au, score_ben, frame = FALSE))
```



```
item_ben = benessere[, grep("item", colnames(benessere))]  
  
par(mfrow = c(2, round(ncol(item_ben)/2 + 0.2)))  
temp = NULL  
for (i in 1:ncol(item_ben)) {  
  temp = table(benessere$genere, item_ben[,i])  
  for (j in 1:nrow(temp)) {  
    temp[j,] = temp[j,]/table(item_ben[,i])  
  }  
  barplot(temp, ylim=c(0,1), legend = rownames(temp),  
          main = colnames(item_ben)[i])  
  abline(h = .5, lty = 2, col = "red")  
}
```

```
interaction.plot()
```

Permette di vedere l'interazione tra due variabili a seconda di una terza variabile:

```
interaction.plot(x, v. categoriale, y)
```

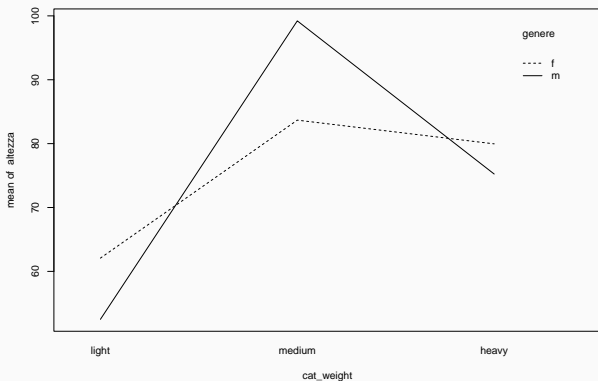


```
babies$cat_weight = with(babies,
  ifelse(
    peso <= quantile(babies$peso)[2],
    "light",
    ifelse(peso > quantile(babies$peso)[2] &
      peso > quantile(babies$peso)[4],
      "medium",
      "heavy")))
babies$cat_weight = factor(babies$cat_weight,
  levels = c("light",
    "medium",
    "heavy"))
```

```
babies
```

	id	genere	peso	altezza	cat_weight
1	baby1	f	7.424646	62.07722	light
2	baby2	m	7.442727	58.18877	light
3	baby3	f	9.512598	84.52737	heavy
4	baby4	f	11.306349	85.13573	medium
5	baby5	m	9.345165	75.23783	heavy
6	babv6	m	5.411290	46.80163	light

```
interaction.plot(cat_weight, genere, altezza))
```

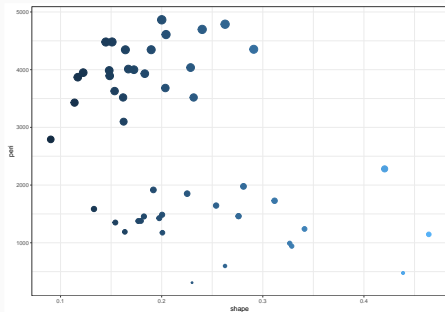


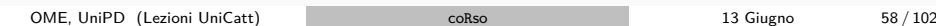
ggplot2 (Grammar of Graphics plot, Wickman, 2016) is one of the best packages for plotting raw data and results:

```
ggplot(data,
  ars(x = var.x,
      y = var.y,
      col = var.color, # factor or character
      fill = var.filling, # factor or character
      shape = var.shape, # actor or character
      size = var.size, # numeric
      ...)) + geom_graph.type() + ...
```

Raw data

```
ggplot(rock,
      aes(y=peri,x=shape, color =shape,
          size = peri)) + geom_point() +
  theme_bw() + theme(legend.position = "none")
```



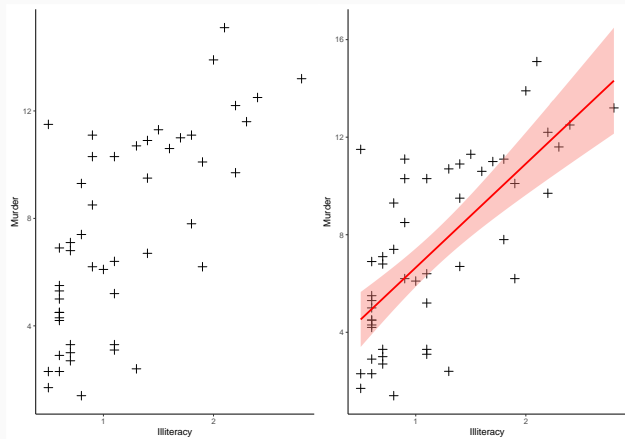



```
install.packages("gridExtra") ; library(gridExtra)
```

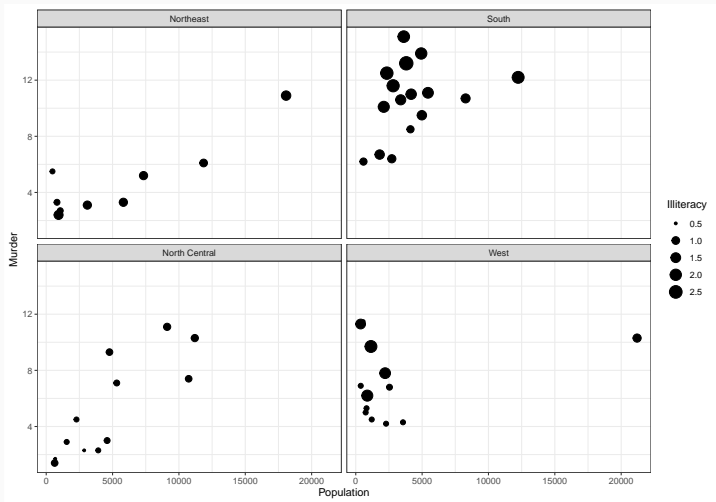
```
murder_lm = ggplot(states, # lm
  aes(x = Illiteracy, y = Murder)) +
  ....
```

```
grid.arrange(murder_raw, murder_lm,
              nrow=1) # plots forced to be the same row
```

Combine the plots together



Multi Panel



Multi panel code

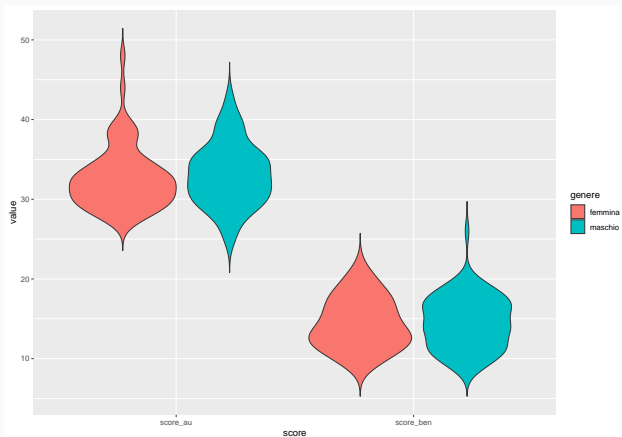
```
states = data.frame(state.x77, state.name = state.name,
                     state.region = state.region)

ggplot(states,
        aes(x = Population, y = Murder,
             size = Illiteracy)) + geom_point() +
  facet_wrap(~state.region) + theme_bw()
```


	ID	score	value
1.score_au	1	score_au	33
2.score_au	2	score_au	40
3.score_au	3	score_au	38
4.score_au	4	score_au	38
5.score_au	5	score_au	26
6.score_au	6	score_au	33


```
score_long = merge(score_long, benessere[, c("ID", "genere")])
```

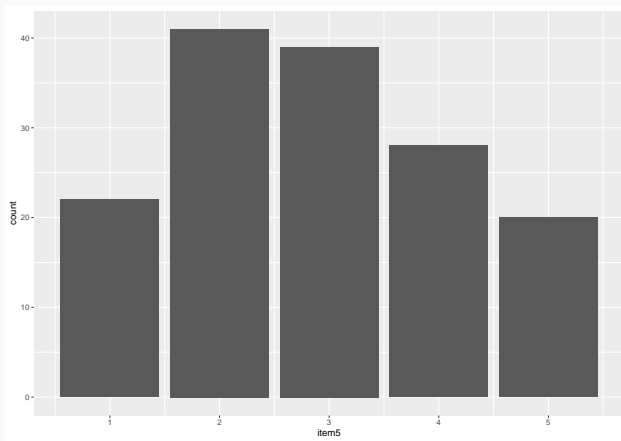
```
ggplot(score_long,
       aes(x = score, y = value,
           fill = genere)) + geom_violin(trim = FALSE)
```



69 / 102

```
geom_bar(stat = "count"):
```

```
ggplot(benessere,
       aes(x = item5)) + geom_bar(stat = "count")
```




```
geom_bar(stat = "identity")
```

```
item_5 = data.frame(table(benessere$item5)/nrow(benessere))
item_5
```

	Var1	Freq
1	1	0.1466667
2	2	0.2733333
3	3	0.2600000
4	4	0.1866667
5	5	0.1333333

Now the plot

```
ggplot(item_5,
       aes(x = Var1, y = Freq)) + geom_bar(stat = "identity")
```

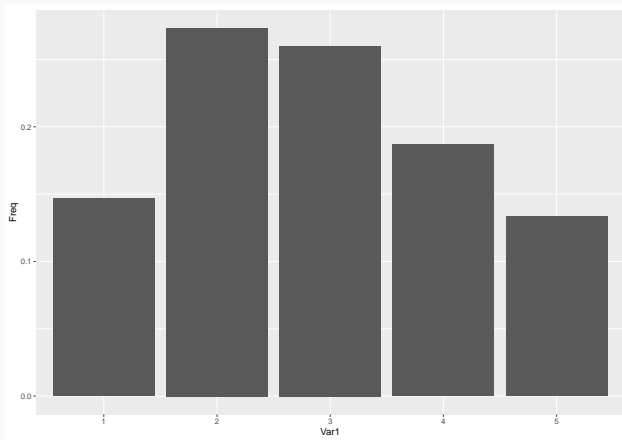


Table of Contents

- 1 Import data
- 2 Compute sum scores
- 3 Graphics
- 4 Esport graphs**
- 5 R for statistical computing
- 6 Classical hypothesis testing in R
- 7 Generalized Linear Models (GLMs)

In automatico

```
pdf("nome_grafico.pdf")
png("nome_grafico.png")
tiff("nome_grafico.tiff")
jpeg("nome_grafico.jpeg")
bmp("nome_grafico.bmp")
```

Remember to turn off the graphical device

```
dev.off()
```

```
pdf("il_mio_violin.pdf")
ggplot(score_long,
       aes(x = score, y = value,
           fill = score)) + geom_violin(trim = FALSE)
dev.off()
```


Manually

```
knitr::include_graphics("data/esporta.png")
```

Table of Contents

- 1 Import data
- 2 Compute sum scores
- 3 Graphics
- 4 Esport graphs
- 5 R for statistical computing**
- 6 Classical hypothesis testing in R
- 7 Generalized Linear Models (GLMs)

```
see library(help=stats)
```

Table of Contents

- 1 Import data
- 2 Compute sum scores
- 3 Graphics
- 4 Esport graphs
- 5 R for statistical computing
- 6 Classical hypothesis testing in R**
- 7 Generalized Linear Models (GLMs)

- `chisq.test()`: contingency table χ^2 tests
- `cor.test()`: association between paired samples
- `t.test()`: one- and two-sample t tests (also for paired data)
- `lm(y ~ x1 + x2 + x3 ...)`: fit a linear model
- `glm(y ~ x1 + x2 + x3 ...)`: fit a generalize linear model

p -value:

If $p < \alpha$ (i.e., the probability of rejecting the null hypothesis when it is true) \rightarrow the null hypothesis is rejected



χ^2 test

Independence of observations

Hypothesis:

- $H_0: P(X_{ij} = k) = p_k$ for all $i = 1, \dots, r$ and $j = 1, \dots, c$
- $H_0: P(X_{ij} = k) \neq P(X_{i'j} = k)$ for *at least* one $i \in \{1, \dots, r\}$ and $j \in \{1, \dots, c\}$

Test statistic:

$$\chi^2 = \sum_{i=1}^n \frac{(x_{ij} - \hat{x}_{ij})^2}{\hat{x}_{ij}}, \quad \chi^2 \sim \chi^2(r-1)(c-1)$$

In R:

```
tab <- xtabs(~ education + induced, infert)
chisq.test(tab)
```

Correlation test

Hypothesis:

- $H_0: \rho_{XY} = 0, H_1: \rho_{xy} \neq 0$
- $H_0: \rho_{XY} = 0, H_1: \rho_{xy} < 0$
- $H_0: \rho_{XY} = 0, H_1: \rho_{xy} > 0$

Test statistic:

$$T = \frac{r_{xy}}{\sqrt{1 - r_{xy}^2}} \sqrt{n - 2}, \quad T \sim t(n - 2)$$

In R:

```
cor.test(~ speed + dist, cars,
         alternative = "two.sided")
```


Two (indepdent) sample t test

Independent samples from normally distributions where σ^2 are unknown but homogeneous

- $H_0: \mu_{x_1-x_2} = 0, H_1: \mu_{x_1-x_2} \neq 0$
- $H_0: \mu_{x_1-x_2} = 0, H_1: \mu_{x_1-x_2} < 0$
- $H_0: \mu_{x_1-x_2} = 0, H_1: \mu_{x_1-x_2} > 0$

Test statistic:

$$T = \frac{\bar{x}_1 - \bar{x}_2}{\sigma_{\bar{x}_1 - \bar{y}_2}}, \quad T \sim t(n_1 + n_2 - 2)$$

R function:

```
t.test(len ~ supp, data = ToothGrowth,
       var.equal = FALSE)
```

Two (depedent) sample t test

Observations on the same sample

Hypothesis:

- $H_0: \mu_D = 0, H_1: \mu_D \neq 0$
- $H_0: \mu_D = 0, H_1: \mu_D < 0$
- $H_0: \mu_D = 0, H_1: \mu_D > 0$

Test statistic:

$$T = \frac{d}{\sigma_d}, \quad T \sim t(m-1)$$

R function:

```
with(sleep,
      t.test(extra[group == 1],
              extra[group == 2], paired = TRUE))
```

Table of Contents

- 1 Import data
- 2 Compute sum scores
- 3 Graphics
- 4 Esport graphs
- 5 R for statistical computing
- 6 Classical hypothesis testing in R
- 7 Generalized Linear Models (GLMs)**

Statistical models are represented by formulae which are extremely close to the actual statistical notation:

91 / 102

In R:

```
coef()      # Extract the regression coefficients
summary()   # Print a comprehensive summary of the results of
            # the regression analysis
anova()     # Compare nested models and produce an analysis
resid()     # Extract the (matrix of) residuals
plot()      # Produce four plots, showing residuals, fitted
            # values and some diagnostics
model.matrix()
            # Return the design matrix
```

```
vcov()      # Return the variance-covariance matrix of the
            # main parameters of a fitted model object
predict()   # A new data frame must be supplied having the
            # same variables specified with the same labels
            # as the original. The value is a vector or
            # matrix of predicted values corresponding to
            # the determining variable values in data frame
step()      # Select a suitable model by adding or dropping
            # terms and preserving hierarchies. The model
            # with the smallest value of AIC (Akaike's
            # Information Criterion) discovered in the
            # stepwise search is returned
```


$g()$ is the link functions that connects the mean to the linear combination of predictors.

A GLM is composed of three elements: The response distribution, the link function, and the linear combination of predictors

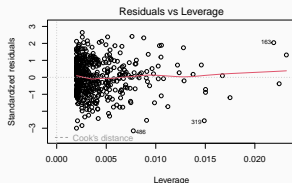
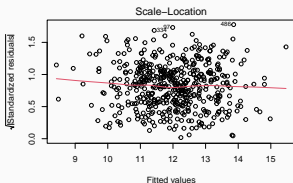
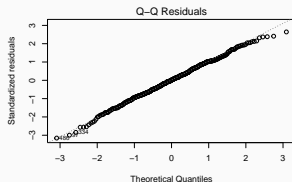
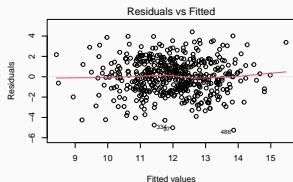
In R:

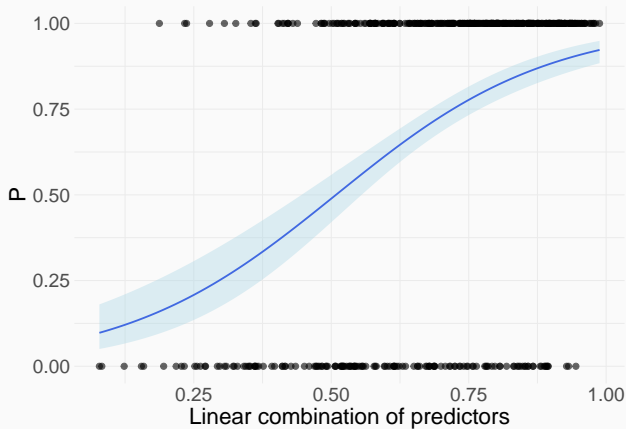
```
glm(y ~ x1 + x2 + ... + xk, family(link), data)
```



Diagnostic

```
par(mfrow = c(2,2))
plot(model)
```





101 / 102

A special link function to each response variable. In R some different link functions are available by default:

102 / 102