

María Isabel Borrajo García, Mercedes Conde Amboage e Rosa María Crujeiras Casais (Departamento de Estatística, Análise Matemática e Optimización)

**R** é unha linguaxe de programación orientada ao ámbito da estatística. Trátase dunha linguaxe interpretada que se emprega maioritariamente de xeito interactivo, distribuída baixo licenza GNU, e que dispón dunha serie de características no módulo base que se pode ir ampliando por medio de librarías. O proceso de instalación e toda a información relativa ao programa pode atoparse en <http://www.r-project.org/>.

Para poder empregar calquera das librarías ou paquetes os pasos que cómpre seguir son:

- Instalar a librería dende o repositorio CRAN (precísase de conexión á Rede), para o que se emprega a orde `install.packages("nome_da_librería")`.
- Cargar a librería na sesión de **R** actual para poder empregar as funcións que contén, isto farase coa sentenza `library(nome_do_paquete)`.

Para calquera dúbida sobre estas ou outras funcións do programa, pode consultarse a súa páxina de axuda tecleando na consola `?nome_da_función` ou `help(nome_da_función)`.

O programa estatístico **R** dispón de diversas interfaces gráficas máis amigables para as persoas usuarias non expertas. Unha das máis habitualmente empregadas é RStudio, que pode obterse de xeito libre dende <https://rstudio.com/>, outra opción dispoñible é *R Commander*, unha interface incluída dentro do propio **R** mediante a librería `Rcmdr`.

### BÁSICOS

O primeiro que debemos de saber para manexar **R** é que a asignación (gardar un valor nunha variable) faise mediante o símbolo `=` ou `<-`, por exemplo, `x=5` ou `x<-5` para gardar o valor 5 nunha variable denominada x.

O seguinte que se debe coñecer son as distintas clases de obxectos que **R** permite manexar:

- **logical**: consta unicamente de dous valores de tipo lóxico que son TRUE (denotado tamén por T) ou FALSE (denotado tamén por F).
- **integer**: almacena valores de tipo enteiro (números naturais e opostos, sen decimais).
- **numeric**: almacena valores de tipo real.
- **complex**: almacena valores de tipo complexo (parte real e parte imaxinaria).
- **character**: almacena valores de tipo carácter, é dicir, secuencias de letras ou texto.
- **factor**: almacena variables de tipo categórico, onde as categorías constitúen os niveis do factor.

```
> x
[1] 1 2 3 4 5
>
> A=matrix(1:6,nr=2,nc=3)
> A
      [,1] [,2] [,3]
[1,]    1    3    5
[2,]    2    4    6
> |
```

A continuación verase unha colección de comandos básicos para traballar cos distintos tipos de obxectos, así como coas diferentes estruturas de datos que se poden manexar en **R**: `vector`, `matrix` (matriz), `list` (lista), `data.frame` (base de datos),...

#### Creación de datos

<code>c(...)</code>	función para combinar argumentos separados por comas e formar un vector.
<code>dende:ata</code>	xera unha secuencia que comeza en <code>dende</code> e remata en <code>ata</code> con paso unha unidade.
<code>factor(x,levels=lev)</code>	crea un obxecto de tipo factor con categorías dadas en <code>lev</code> .
<code>matrix(datos,nrow=n,ncol=m,byrow=F)</code>	crea unha matriz a partir de datos de dimensión <code>n x m</code> onde os elementos se colocan por columnas; colocaranse enchendo filas se o argumento <code>byrow=T</code> .
<code>rbind(x1,x2,...); cbind(x1,x2,...)</code>	combina os elementos <code>x1,x2,...</code> (vector, matriz ou base de datos) concatenándoos por filas (r) ou columnas (c).
<code>rep(x,times=v)</code>	constrúe un vector repetindo <code>v</code> veces o elemento <code>x</code> , que pode ser un valor, un vector, unha matriz, unha lista,...
<code>seq(dende,ata,by=p,length=1)</code>	xera unha secuencia que comeza en <code>dende</code> e remata en <code>ata</code> , ben con paso <code>p</code> ou ben de lonxitude 1.

## BÁSICOS (cont.)

### Extraer valores

<code>x[k]</code>	accede ao elemento k-ésimo do vector <code>x</code> .
<code>x[-k]</code>	accede a todos os elementos do vector <code>x</code> salvo ao k-ésimo.
<code>x[c(1,3,7)]</code>	accede ao primeiro, terceiro e sétimo elementos do vector <code>x</code> .
<code>x[x&lt;2]</code>	accede a todos os elementos de <code>x</code> que cumpran a condición de seren menores que 2.
<code>x[i,k]</code>	accede ao elemento da fila <code>i</code> e columna <code>k</code> da matriz ou base de datos <code>x</code> .
<code>x[i, ]</code>	accede á fila <code>i</code> -ésima da matriz ou base de datos <code>x</code> .
<code>x[ ,k]</code>	accede á columna <code>k</code> -ésima da matriz ou base de datos <code>x</code> .
<code>x[ ,c(2,4)]</code>	accede ás columnas segunda e cuarta da matriz ou base de datos <code>x</code> .
<code>x[[k]]</code>	accede ao elemento k-ésimo da lista <code>x</code> .
<code>x[["nome"]]</code>	accede ao elemento da lista <code>x</code> denominado <code>nome</code> .
<code>x\$nome</code>	accede ao elemento da lista ou base de datos denominado <code>nome</code> .



### Información sobre variables

<code>class(x)</code>	devolve o tipo de obxecto que é <code>x</code> .
<code>dim(x)</code>	devolve a dimensión (nº de filas e nº de columnas) da matriz ou base de datos <code>x</code> .
<code>is.numeric(x)</code>	devolve cun valor lóxico indicando se o obxecto <code>x</code> é (T) ou non é (F) de tipo numérico. Análogo para <code>is.character</code> , <code>is.na</code> , <code>is.complex</code> , <code>is.data.frame</code> ,...
<code>length(x)</code>	devolve o número de elementos do obxecto <code>x</code> .
<code>nrow(x); ncol(x)</code>	devolve o número de filas ( <code>nr</code> ) ou columnas ( <code>nc</code> ) da matriz ou base de datos <code>x</code> .
<code>summary(x)</code>	resume a información contida no obxecto <code>x</code> ; dependendo da clase de <code>x</code> proporciona diferente información.

### Selección e manipulación de datos

<code>choose(x,k)</code>	calcula as posibles combinacións de coller <code>k</code> elementos en <code>n</code> repeticións.
<code>cut(x,breaks=b)</code>	crea a partir do obxecto <code>x</code> un novo obxecto de tipo factor mediante a división de <code>x</code> en intervalos (categorías); <code>b</code> pode ser o número de intervalos, ou os propios puntos de corte que delimitan os intervalos.
<code>det(x)</code>	calcula o determinante da matriz <code>x</code> .
<code>diag(x)</code>	accede á diagonal da matriz <code>x</code> .
<code>eigen(x)</code>	calcula os autovalores e autovectores da matriz <code>x</code> .
<code>match(x,y)</code>	define un vector da mesma lonxitude de <code>x</code> , cos elementos de <code>x</code> que están en <code>y</code> e NA (valor co que <b>R</b> codifica os datos ausentes) noutro caso.
<code>round(x,digits=nd)</code>	redondea os valores numéricos de <code>x</code> ao número de decimais dado por <code>nd</code> .
<code>rowsum(x); colsum(x)</code>	calcula a suma por filas ( <code>row</code> ) ou columnas ( <code>col</code> ) da matriz ou base de datos <code>x</code> .
<code>rowMeans(x); colMeans(x)</code>	calcula a media por filas ( <code>row</code> ) ou columnas ( <code>col</code> ) da matriz ou base de datos <code>x</code> .
<code>solve(x)</code>	calcula a matriz inversa da matriz <code>x</code> .
<code>sort(x, decreasing=F)</code>	ordena de xeito crecente (ou decrecente se <code>decreasing=T</code> ) os elementos de <code>x</code> . Empregando o argumento <code>index.return=T</code> pódese obter o valor da posición dos elementos segundo a súa orde (crecente ou decrecente).
<code>t(x)</code>	calcula a matriz trasposta da matriz <code>x</code> .
<code>x%*%y</code>	calcula o produto escalar ou matricial dos vectores ou matrices <code>x</code> e <code>y</code> .
<code>which(x&lt;2)</code>	determina a/a(s) posición/s dos elementos de <code>x</code> que cumpran a condición de ser menores que 2.
<code>which.min(x); which.max(x)</code>	determina a posición que ocupa en <code>x</code> o seu elemento de menor ( <code>min</code> ) ou maior ( <code>max</code> ) valor.

$$A = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{pmatrix}$$

### Lectura e escritura de datos

<code>data(nome)</code>	carga o conxunto de datos <code>nome</code> gardado en <b>R</b> ; se non se especifica ningún argumento, lista os conxuntos de datos dispoñibles na sesión actual de <b>R</b> .
<code>load(...)</code>	carga datos dende unha área de traballo previamente gardada (en xeral con <code>save</code> ).
<code>read.table(file="nome.ext", sep=";", dec=".", header=T)</code>	carga os datos almacenados no arquivo <code>nome</code> , de extensión <code>.ext</code> ( <code>.txt</code> , <code>.dat</code> ). Co argumento <code>sep</code> determínase o elemento que serve para separar valores no arquivo, con <code>dec</code> indícase o símbolo empregado para o separador decimal e <code>header</code> é un argumento lóxico que indica se a primeira fila da base de datos contén os nomes das variables (T) ou contén os valores do primeiro individuo (F). Análogo <code>read.csv</code> para extensións <code>.csv</code> .

## BÁSICOS (cont.)

### Lectura e escritura de datos (cont.)

<code>attach(dat)</code>	carga no contorno de traballo as variables da base de datos <code>dat</code> .
<code>complete.cases(x,y)</code>	indica as posicións comúns sen NA nos elementos <code>x</code> e <code>y</code> .
<code>head(dat)</code>	amosa as seis primeiras liñas do elemento <code>dat</code> .
<code>names(dat)</code>	indica os nomes dos elementos almacenados en <code>dat</code> .
<code>save(nome, file=arquivo)</code>	garda o elemento <code>nome</code> en <code>arquivo</code> .
<code>save.image()</code>	garda todos os elementos que se atopen na memoria da sesión nese momento nun arquivo <code>.RData</code> .
<code>View(dat)</code>	permite visualizar a base de datos <code>dat</code> nun formato de táboa ao estilo dunha folla de cálculo.
<code>write.table(x, file="nome.ext")</code>	garda o elemento <code>x</code> nun arquivo <code>.ext</code> ( <code>.txt</code> , <code>.csv</code> ) tras convertelo en base de datos.



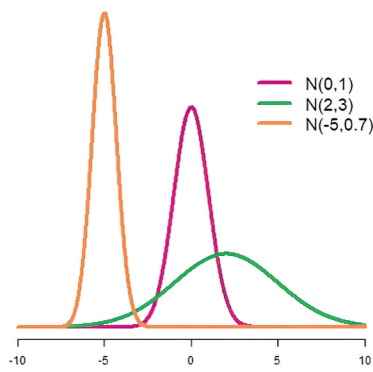
## ESTADÍSTICA DESCRIPTIVA

<code>addmargins(tab)</code>	engade a fila e columna dos totais (marxinais) a unha táboa de dobre entrada <code>tab</code> .
<code>barplot(tab)</code>	representa o diagrama de barras asociado ás frecuencias <code>tab</code> , que se poden obter, por exemplo, coa función <code>table</code> .
<code>boxplot(x)</code>	representa un diagrama de caixa para os valores de <code>x</code> .
<code>boxplot(x~y)</code>	representa un diagrama de caixa para os valores de <code>x</code> agrupados segundo as categorías do factor <code>y</code> .
<code>cov(x,y); cor(x,y)</code>	calcula a cuasicovarianza ( <code>cov</code> ) ou correlación ( <code>cor</code> ) entre os valores de <code>x</code> e de <code>y</code> .
<code>cumsum(x); cumprod(x)</code>	vector no que o elemento <code>i</code> -ésimo é a suma ( <code>sum</code> ) ou produto ( <code>prod</code> ) acumulados dos elementos <code>x[1:i]</code> .
<code>describe(x)</code>	calcula numerosas características dos elementos de <code>x</code> : tamaño da mostra, número de valores perdidos, media, cuantís, táboa de frecuencias,... Esta función está dispoñible na librería <code>Hmisc</code> .
<code>diff(x)</code>	calcula un vector de lonxitude unha unidade menor que a lonxitude de <code>x</code> coas diferenzas iteradas dos elementos de <code>x</code> .
<code>hist(x, breaks=b, freq=T)</code>	representa un histograma con <code>b+1</code> intervalos (ou cos intervalos dados en <code>b</code> ); se <code>freq=T</code> representa as frecuencias absolutas, en caso contrario representa densidades.
<code>IQR(x)</code>	calcula o rango intercuartílico (terceiro menos primeiro cuartil) dos elementos de <code>x</code> .
<code>kurtosis(x)</code>	calcula o coeficiente de curtose para os valores de <code>x</code> . Esta función está dispoñible na librería <code>moments</code> .
<code>max(x)</code>	calcula o valor máximo dos elementos de <code>x</code> .
<code>mean(x)</code>	calcula a media da mostra <code>x</code> . Engadindo o argumento <code>trim=p</code> , calcula a media recortada coa proporción de datos eliminados dada por <code>p</code> .
<code>median(x)</code>	calcula a mediana da mostra <code>x</code> .
<code>min(x)</code>	calcula o valor mínimo dos elementos de <code>x</code> .
<code>pie(tab)</code>	representa o diagrama de sectores asociado ás frecuencias <code>tab</code> , que se poden obter, por exemplo, coa función <code>table</code> .
<code>pie3D(tab)</code>	representa o diagrama de sectores con volume (3D) asociado ás frecuencias <code>tab</code> , que se poden obter, por exemplo, coa función <code>table</code> .
<code>proportions(tab, margin=m)</code>	calcula as proporcións condicionadas da táboa <code>tab</code> dadas as marxinais <code>m</code> .
<code>qqplot(x,y)</code>	representa os cuantís de <code>x</code> fronte aos cuantís de <code>y</code> .
<code>quantile(x, probs=p)</code>	calcula o cuantil dos elementos de <code>x</code> definido pola probabilidade en <code>p</code> . Tamén é válido que <code>p</code> sexa un vector de probabilidades.
<code>range(x)</code>	calcula simultaneamente os valores mínimo e máximo dos elementos de <code>x</code> .
<code>skewness(x)</code>	calcula o coeficiente de asimetría para os valores de <code>x</code> . Esta función está dispoñible na librería <code>moments</code> .
<code>sum(x); prod(x)</code>	calcula a suma ( <code>sum</code> ) ou produto ( <code>prod</code> ) dos elementos de <code>x</code> .
<code>table(x)</code>	conta os elementos dos distintos valores de <code>x</code> .
<code>table(x,y)</code>	realiza unha táboa de dobre entrada dos elementos <code>x</code> e <code>y</code> .
<code>unique(x)</code>	vector cos elementos de <code>x</code> eliminando os elementos repetidos.
<code>var(x); sd(x)</code>	calcula a cuasivarianza ( <code>var</code> ) ou cuasidesviación típica ( <code>sd</code> ) dos elementos de <code>x</code> .
<code>weighted.mean(x,w)</code>	calcula a media da mostra <code>x</code> ponderada polo vector de pesos <code>w</code> .



## VARIABLES ALEATORIAS

dDISTR(x)



calcula a función de masa de probabilidade (función de densidade) nun punto  $x$  asociada a unha variable DISTR discreta (variable continua). Amais, os argumentos desta función dependen da DISTR considerada. As principais variables que podemos empregar con **R** son:

DISTRIBUCIÓN	DISTR	ARGUMENTOS
Binomial ( $n, p$ )	binom	size= $n$ , prob= $p$
Binomial negativa ( $n, p$ )	nbinom	size= $n$ , prob= $p$
Poisson ( $\lambda$ )	pois	lambda= $\lambda$
Uniforme ( $a, b$ )	unif	min= $a$ , max= $b$
Normal ( $\mu, \sigma$ )	norm	mean= $\mu$ , sd= $\sigma$
Exponencial ( $\lambda$ )	exp	rate= $\lambda$
Gamma ( $n, \lambda$ )	gamma	shaper= $n$ , rate= $\lambda$
T-Student con $gl$ graos de liberdade	t	df= $gl$
$\chi^2$ con $gl$ graos de liberdade	chisq	df= $gl$
F de Snedecor con $gl_1$ e $gl_2$ graos de liberdade	f	df1= $gl_1$ , df2= $gl_2$

density(x)

estima a función de densidade asociada a unha mostra  $x$  empregando métodos non paramétricos.

ecdf(x)

calcula a función de distribución empírica asociada a unha mostra  $x$ .

fitdistr(x, densfun="NOME")

emprega o método de máxima verosimellanza para axustar os parámetros asociados a unha distribución NOME para a mostra  $x$ . As distintas distribucións compatibles con esta función poden verse en ?fitdistr. Esta función está dispoñible na librería MASS.

pDISTR(x)

calcula a función de distribución nun punto  $x$  asociada a unha variable DISTR discreta ou continua. Esta función pódese aplicar para as mesmas distribucións que dDISTR.

qDISTR( $\tau$ )

calcula o cuantil de orde  $\tau$  asociado a unha variable DISTR discreta ou continua. Esta función pódese aplicar para as mesmas distribucións que dDISTR.

qqPlot(x)

representa os cuantís asociados á mostra  $x$  fronte aos cuantís dunha distribución normal. Esta función está dispoñible na librería car.

rDISTR(n)

xera  $n$  datos asociados a unha variable DISTR discreta ou continua. Esta función pódese aplicar para as mesmas distribucións que dDISTR.

## INFERENCIA ESTADÍSTICA

ks.test(x, "pDISTR")

realiza un contraste de Kolmogorov-Smirnov para contrastar se a mostra  $x$  (que procede dunha variable aleatoria continua) segue unha distribución DISTR.

chisq.test(x,y)

realiza un contraste de independencia sobre dúas variables categóricas das que se coñecen as mostras  $x$  e  $y$ . Tamén se podería realizar o contraste introducindo simplemente unha táboa de continxencia.

prop.test

(x, n, conf.level= $\beta$ )

calcula un intervalo de confianza de nivel  $\beta$  para a probabilidade de éxito asociado a un experimento de Bernoulli do que se observaron  $n$  realizacións das cales  $x$  foron éxitos. Se  $n$  e  $x$  fosen vectores de dimensión 2, calcularía un intervalo de confianza de nivel  $\beta$  para a diferenza das probabilidades de éxito asociadas a dous experimentos de Bernoulli.

(x, n, p=p, alternative=alt)

realiza un contraste de hipóteses para a probabilidade de éxito asociado a un experimento de Bernoulli do que se observaron  $n$  realizacións, das cales  $x$  foron éxitos. O valor da probabilidade de éxito que se quere contrastar será  $p$  mentres que o tipo de contraste vén determinado polo argumento alternative\*. Amais, se  $n$  e  $x$  fosen vectores de dimensión 2, realizaría un contraste para a diferenza das probabilidades de éxito asociadas a dous experimentos de Bernoulli (neste caso o argumento  $p$  non se emprega).

shapiro.test(x)

realiza un contraste de normalidade para unha variable  $X$  da que se ten a mostra  $x$ .

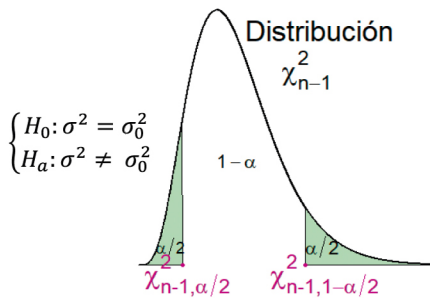
\* Nas funcións prop.test, sigma.test, t.test, var.test e z.test o argumento alternative pode ser "two.sided" (contraste bilateral), "less" (contraste unilateral esquerdo) ou "greater" (contraste unilateral dereito).

## INFERENCIA ESTADÍSTICA (cont.)

sigma.test

(x, conf.level= $\beta$ )

(x, sigmasq =  $\sigma^2$ , alternative=alt)



calcula un intervalo de confianza de nivel  $\beta$  para a varianza poboacional dunha distribución normal da que se coñece unha mostra x.

realiza un contraste de hipóteses para a varianza poboacional dunha distribución normal da que se coñece unha mostra x. O valor da varianza poboacional que se quere contrastar será  $\sigma^2$  mentres que o tipo de contraste vén determinado polo argumento `alternative**`.

Esta función está dispoñible na librería TeachingDemos.

t.test

(x, conf.level= $\beta$ )

(x, y, conf.level= $\beta$ , paired=T)

(x, mu= $\mu_0$ , alternative=alt)

(x, y, mu= $\mu_0$ , paired=F)

calcula un intervalo de confianza de nivel  $\beta$  para a media poboacional dunha distribución normal da que se coñece unha mostra x.

calcula un intervalo de confianza de nivel  $\beta$  para a diferenza de medias poboacionais de dúas distribucións normais dependentes (`paired=T`) ou independentes (`paired=F`) das que se coñecen as mostras x e y. No caso de que `paired=F`, emprégase o argumento `var.equal=T` ou `F` para indicar que as varianzas poboacionais son iguais ou non.

realiza un contraste de hipóteses para a media poboacional dunha distribución normal da que se coñece unha mostra x. O valor da media poboacional que se quere contrastar será  $\mu_0$  mentres que o tipo de contraste vén determinado polo argumento `alternative**`.

calcula un contraste de hipóteses para a diferenza de medias poboacionais de dúas distribucións normais dependentes (`paired=T`) ou independentes (`paired=F`) con varianzas coñecidas ou descoñecidas (`var.equal=T` ou `F`) das que se coñecen as mostras x e y. O valor da diferenza de medias poboacionais que se quere contrastar será  $\mu_0$  mentres que o tipo de contraste vén determinado polo argumento `alternative**`.

var.test

(x, y, conf.level= $\beta$ )

(x, y, ratio=r, alternative=alt)

calcula un intervalo de confianza de nivel  $\beta$  para o cociente das varianzas poboacionais de dúas distribucións normais das que se coñecen as mostras x e y.

realiza un contraste de hipóteses para o cociente das varianzas poboacionais de dúas distribucións normais das que se coñecen as mostras x e y. O valor do cociente de varianzas poboacionais que se quere contrastar será r mentres que o tipo de contraste vén determinado polo argumento `alternative**`.

z.test

(x, stdev =  $\sigma$ , conf.level= $\beta$ )

(x, mu= $\mu_0$ , stdev =  $\sigma$ , alternative=alt)

calcula un intervalo de confianza de nivel  $\beta$  para a media poboacional dunha distribución normal con desviación típica  $\sigma$  coñecida da que se coñece unha mostra x.

realiza un contraste de hipóteses para a media poboacional dunha distribución normal con desviación típica  $\sigma$  coñecida da que se coñece unha mostra x. O valor da media poboacional que se quere contrastar será  $\mu_0$  mentres que o tipo de contraste vén determinado polo argumento `alternative**`.

Esta función está dispoñible na librería TeachingDemos.

\*\* Nas funcións `prop.test`, `sigma.test`, `t.test`, `var.test` e `z.test` o argumento `alternative` pode ser “two.sided” (contraste bilateral), “less” (contraste unilateral esquerdo) ou “greater” (contraste unilateral dereito).

## MODELOS DE REGRESIÓN

`aov(y~grupos)`

realiza un contraste de igualdade de medias poboacionais asociadas a unha distribución normal (da que coñecemos unha mostra  $y$ ) medida en diferentes grupos (determinados polo vector `grupos`).

`boxcox(y~x1+x2+...+xd)`

calcula o parámetro  $\lambda$  óptimo asociado a unha transformación tipo Box-Cox que permite corrixir a falta de cumprimento das hipóteses básicas do modelo de regresión linear múltiple con variable resposta  $y$  e variables explicativas  $x_1, x_2, \dots, x_d$ . Esta función está dispoñible na librería MASS.

`Box.test(x, lag =  $\gamma$ ,  
type="Ljung-Box")`

realiza o contraste de Ljung-Box que comproba se as observacións da mostra  $x$  son non correladas fronte á hipótese alternativa de que exista correlación secuencial de orde  $\gamma$ .

`confint(mod, level =  $\beta$ )`

calcula intervalos de confianza de nivel  $\beta$  asociados aos parámetros do modelo de regresión `mod`.

`glm(y~x1+x2+...+xd, family = NOME)`

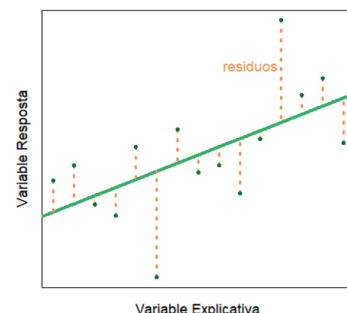
estima un modelo de regresión linear xeneralizado con variable resposta  $y$  e variables explicativas  $x_1, x_2, \dots, x_d$ . O argumento `family` permite seleccionar a función de enlace adecuada en función da natureza da variable resposta.

`hmctest(mod)`

contrasta a homocedasticidade do modelo de regresión linear múltiple `mod`. Esta función está dispoñible na librería `lmtest`.

`lm(y~x1+x2+...+xd)`

estima un modelo de regresión linear múltiple con variable resposta  $y$  e variables explicativas  $x_1, x_2, \dots, x_d$  empregando o método de mínimos cadrados.



`names(mod)`

proporciona os nomes das estimacións asociadas a un modelo de regresión, como os residuos (`mod$residuals`) ou os valores axustados (`mod$fitted.values`).

`nls(y~fun( $\theta_1, \theta_2, \dots, \theta_k$ ,  
start=list( $(\theta_1=n_1, \theta_2=n_2, \dots, \theta_k=n_k)$ ))`

estima un modelo de regresión non linear cuxa expresión vén determinada por `fun` (coñecida) salvo polos parámetros  $\theta_1, \theta_2, \dots, \theta_k$ . Co argumento `start` introdúcese o punto de arranque do algoritmo iterativo de estimación.

`plot(mod)`

realiza catro representacións gráficas que permiten levar a cabo unha validación e diagnose do modelo de regresión múltiple `mod`.

`predict(mod, data.frame(x1= $\delta_1$ ,  
x2= $\delta_2, \dots, x_d=\delta_d$ ))`

calcula predicións puntuais da variable resposta dun modelo de regresión `mod` con covariables  $x_1, x_2, \dots, x_d$  dada unha nova observación  $(\delta_1, \delta_2, \dots, \delta_d)$  das covariables. Tamén permite calcular intervalos de confianza para a esperanza condicional (`interval="confidence"`) e intervalos de predición (`interval="prediction"`).

`resettest(y ~ x, power=2/3)`

realiza o contraste do modelo linear simple con variable resposta  $y$  e variable explicativa  $x$  fronte a modelos polinómicos de orde 2/3. Esta función está dispoñible na librería `lmtest`.

`sm.regression(x,y,  
model="linear")`

contrasta a linearidade dun modelo de regresión linear múltiple con covariables  $x$  (como máximo de dimensión 2) e variable resposta  $y$ . Esta función está dispoñible na librería `sm`.

`summary(mod)`

proporciona unha lista de estatísticos resumo asociados a un modelo de regresión `mod`, como por exemplo, contrastes de significación dos coeficientes, a estimación da desviación típica do erro ou o coeficiente de determinación.