

# Programación con R. PRACTICA 6

```

1 ##PRACTICA 6
2
3 #BLOQUE 4 TEMA 1
4
5
6 a<-c(1,1,6); b<-c(4,5,1)
7 a-b
8 sum((a-b)^2)
9 sqrt(sum((a-b)^2))
10
11 #pag 7
12 #Lo primero que hacemos es descargar la libreria.
13
14 source("http://www.bioconductor.org/biocLite.R")
15 biocLite()
16 biocLite("multtest")
17 library(multtest);
18 data(golub)
19 golub.gnames[1042,]
20
21 #pag 8
22
23 #Queremos buscar el patron Cyclin en el vector. Las posiciones nos las da INDEX.
24 #grep: search for matches to argument pattern within each element of a character vector:
25 #they differ in the format of an amount of detail in the results.
26 index<-grep("Cyclin", golub.gnames[,2]) #esto es un ejemplo de funcion de distribucion de secuencias
27
28 golub.gnames[index,2]
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

Console

```

> golub.gnames[1042,]
[1] "2354"
> a<-c(1,1,6); b<-c(4,5,1)
> a-b
[1] -3 -4 5
> sum((a-b)^2)
[1] 50
> sqrt(sum((a-b)^2))
[1] 7.071068
> golub.gnames[1042,]
[1] "2354"
>

```

```

9
10 sqrt(sum((a-b)^2))
11
12 #pag 7
13 #Lo primero que hacemos es descargar la libreria.
14
15 source("http://www.bioconductor.org/biocLite.R")
16 biocLite()
17 biocLite("multtest")
18 library(multtest);
19 data(golub)
20 golub.gnames[1042,]
21
22 #pag 8
23
24 #Queremos buscar el patron Cyclin en el vector. Las posiciones nos las da INDEX.
25 #grep: busca coincidencias con el patrón de argumento dentro de cada elemento de un vector de caracteres:
26 #diffieren en el formato de una cantidad de detalles en los resultados.
27 index<-grep("Cyclin", golub.gnames[,2]) #esto es un ejemplo de funcion de distribucion de secuencias
28
29 golub.gnames[index,2]
30
31 #pag9
32
33 #ahora vamos a calcular la distancia euclidea (se pueden usar otras distancias) entre todos los genes cuyo
34 #descriptor contenga la palabra Cyclin, hacemos:
35 dist.cyclin<-dist(golub[index,],method="euclidian") #esto nos da las distancias entre datos
36 diam<-as.matrix(dist.cyclin) #aunque ya nos lo da como una matriz, lo organizamos como matriz (de otro modo)
37 colnames(diam)<-golub.gnames[index,3] #el nombre de las columnas va a ser...
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

Console

```

> gol.true <- factor(golub.c1,levels=0:1,labels= c("ALL","not ALL"))
> gol.pred <- factor(golub[1042,],levels=c("TRUE","FALSE"), #si el factor sale con valor >1.27, TRUE, que se corresponde con ALL
+ labels=c("ALL","notALL"))
> table(gol.pred,gol.true)
gol.true
gol.pred ALL not ALL
ALL      25      1
notALL   2      10
> gol.true<-factor(golub.c1,levels=0:1,labels=c("TRUE","FALSE")) #aquí medimos si se parece al gen "ciclina" el resto de registros
> pred<-prediction(golub[1042,],gol.true)
> perf<-performance(pred,"tpr","fpr")
> plot(perf)
+

```

True positive rate

False positive rate

RStudio

```

25 #they differ in the format of an amount of detail in the results.
26 index<-grep("Cyclin", golub.gnames[,2]) #esto es un ejemplo de funcion de distribucion de secuencias
27 index
28 golub.gnames[index,2]
29
30 #pag9
31
32 #ahora vamos a calcular la distancia euclidea (se pueden usar otras distancias) entre todos los genes cuyo
33 #descriptor conengan la palabra Cyclin, hacemos:
34 dist.cyclin<-dist(golub[index,],method="euclidian") #esto nos da las distancias entre datos
35 diam<-as.matrix(dist.cyclin) #aunque ya nos lo da como una matriz, lo organizamos como matriz (de otro modo)
36 colnames(diam)<-golub.gnames[index,3] #el nombre de las columnas va a ser...
37 rownames(diam)<-colnames(diam) #el nombre de las filas va a ser... (los mismos)
38 diam[1:5,1:5] #muestre de fila 1 a fila 5 y de col 1 a col 5 (porque es muy larga)
39
40 #pag10
41 #vamos a calcular la dist. euclidea del gen 1389_at a todos los genes
42 #para quedarnos con los 10 mas cercanos
43
44 biocLite("genefilter")
45 library("genefilter")
46 biocLite("ALL")
47 library("ALL")
48 data(ALL)
49 c1389_at<-genefilter(ALL, "1389_at", 10, method="euc") #nos da los 10 mas parecidos
50 c1389_at[[1]]$indices
51 round(c1389_at[[1]]$dists,1)
52
53 <
54 (Top Level)

```

```

> dist.cyclin<-dist(golub[index,],method="euclidian") #esto nos da las distancias entre datos
> diam<-as.matrix(dist.cyclin) #aunque ya nos lo da como una matriz, lo organizamos como matriz (de otro modo)
> colnames(diam)<-golub.gnames[index,3] #el nombre de las columnas va a ser...
> rownames(diam)<-colnames(diam) #el nombre de las filas va a ser... (los mismos)
> diam[1:5,1:5]
      D13639_at M68520_at M92287_at U09579_at U11791_at
D13639_at  0.000000  8.821806  11.53349  10.056814  8.669112
M68520_at  8.821806  0.000000  11.70156  5.931260  2.934802
M92287_at 11.533494 11.701562  0.00000  11.991333  11.900558
U09579_at 10.056814  5.931260  11.99133  0.000000  5.698232
U11791_at  8.669112  2.934802  11.90056  5.698232  0.000000

```

Environment History Connections

Global Environment

Object	Class	Attributes
d	num	[1:3] 4 3 1
colours	chr	[1:11] "red" "blue" "yellow" "green"
cont.inue	"S"	
corAF	dbl	0.461520020598407
corHA	dbl	0.74983493804404
corHF	dbl	0.500680424892131
corR	dbl	0.058818470543691
covar	dbl	0.302628282828283

User Library

Package	Description	Version
abind	Combine Multidimensional Arrays	1.4-5
adabag	Applies Multiclass AdaBoostM1, SAMME and Bagging	4.1
akima	Interpolation of Irregularly and Regularly Spaced Data	0.6-2
ALL	A data package	1.18.0
annotate	Annotation for microarrays	1.54.0
AnnotationDbi	Annotation Database Interface	1.38.2
assertthat	Easy Pre and Post Assertions	0.2.0
backports	Reimplementations of Functions Introduced Since R-3.0.0	1.1.1
base64enc	Tools for base64 encoding	0.1-3
bestglm	Best Subset GLM	0.36
BH	Boost C++ Header Files	1.65.0-1
bindr	Parameterized Active Bindings	0.1
bindrcpp	An "Rcpp" Interface to Active Bindings	0.2
Biobase	Biobase: Base functions for Bioconductor	2.36.2
BiocGenerics	S4 generic functions for Bioconductor	0.22.1
BiocInstaller	Install/Update Bioconductor CRAN, and github Packages	1.26.1
bit	A class for vectors of 1-bit booleans	1.1-12
bit64	A S3 Class for Vectors of 64-bit Integers	0.9-7

RStudio

```

37 rownames(diam)<-colnames(diam) #el nombre de las filas va a ser... (los mismos)
38 diam[1:5,1:5] #muestre de fila 1 a fila 5 y de col 1 a col 5 (porque es muy larga)
39
40 #pag10
41
42 # Ahora vamos a encontrar los diez genes más cercanos a uno que ya tenemos.
43 #vamos a calcular la dist. euclidea del gen 1389_at a todos los genes
44 #Esto se puede hacer con la funcion genefinder especificando un índice o un nombre.
45
46 biocLite("genefilter")
47 library("genefilter")
48 biocLite("ALL")
49 library("ALL")
50 data(ALL)
51 c1389_at<-genefilter(ALL, "1389_at", 10, method="euc") #nos da los 10 mas parecidos
52 c1389_at[[1]]$indices
53 round(c1389_at[[1]]$dists,1)
54 featureNames(ALL)[c1389_at[[1]]$indices]
55 str(c1389_at)
56
57 <
58 (Top Level)

```

```

> #biocLite("ALL")
> library("ALL")
> data(ALL)
> c1389_at<-genefilter(ALL, "1389_at", 10, method="euc") #nos da los 10 mas parecidos
> c1389_at[[1]]$indices
[1] 2653 1096 6634 9255 6639 11402 9849 8518 10736
> round(c1389_at[[1]]$dists,1)
[1] 12.6 12.8 12.8 12.8 13.0 13.0 13.1 13.2 13.3 13.4
> featureNames(ALL)[c1389_at[[1]]$indices]
[1] "32629_fat" "1988_at" "36571_at" "39168_at" "36576_at" "41295_at" "39756_g_at" "32254_at" "38438_at"
[10] "40635_at"
> str(c1389_at)
List of 1
 $ 1389_at:List of 2
 .. indices: num [1:10] 2653 1096 6634 9255 6639 ...
 .. dists : num [1:10] 12.6 12.8 12.8 12.8 13 ...
> #pag 13

```

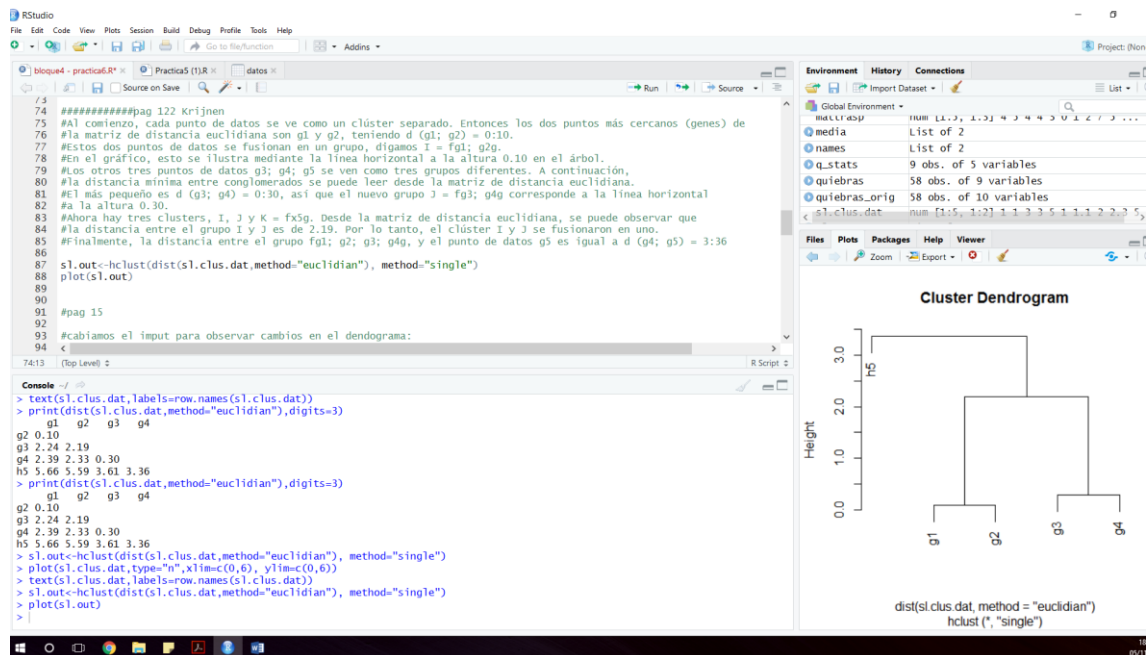
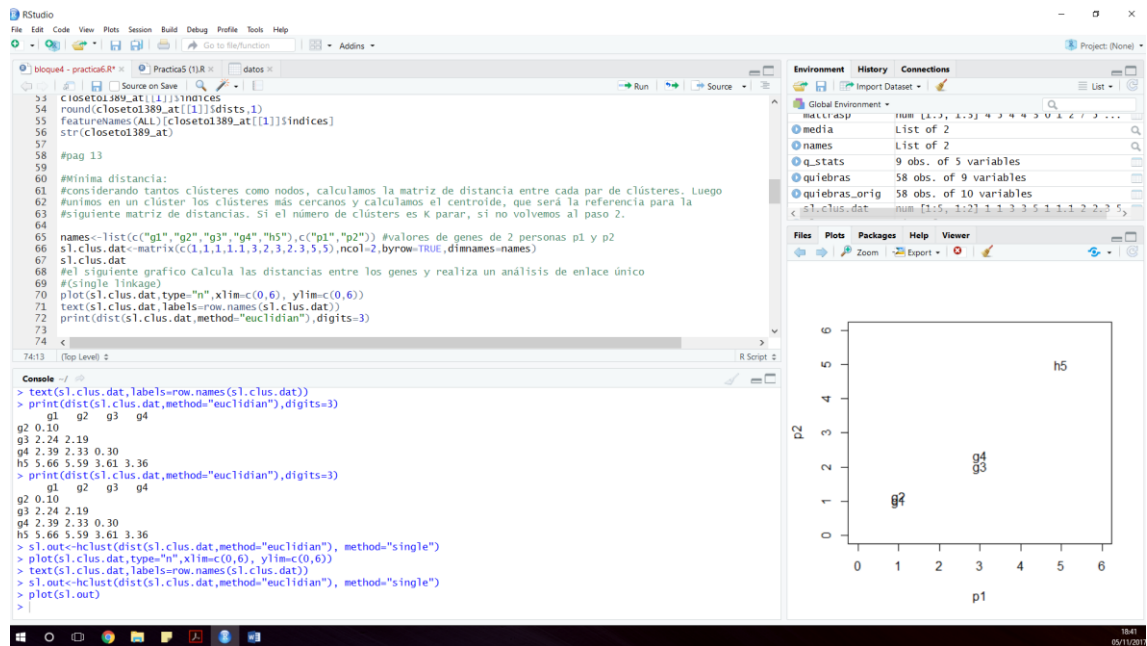
Environment History Connections

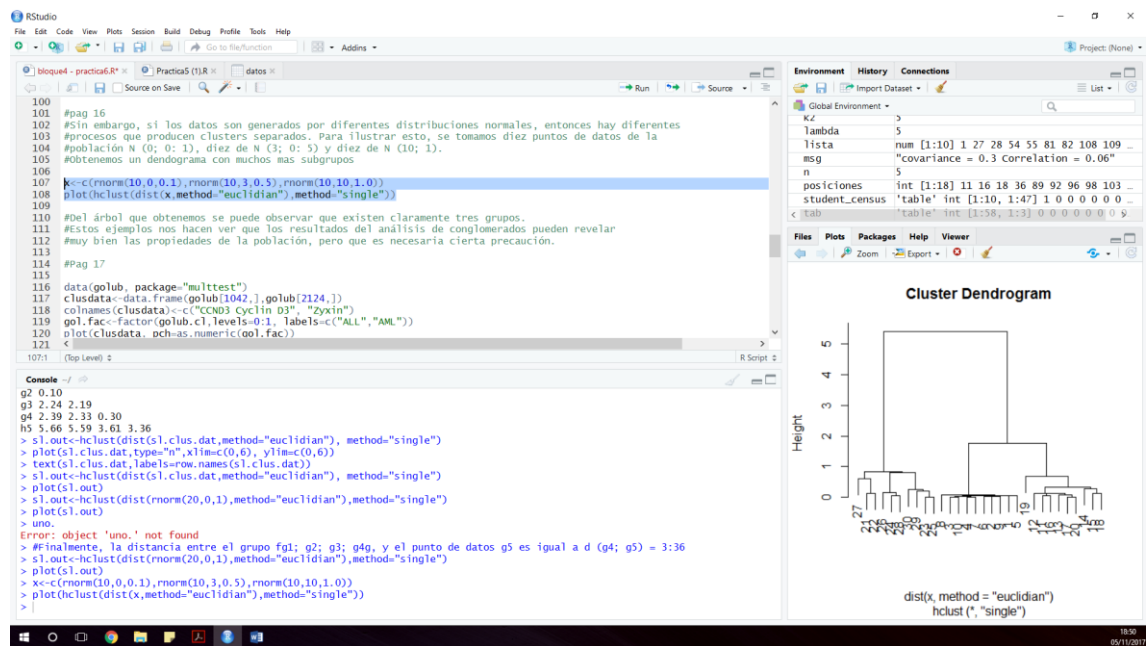
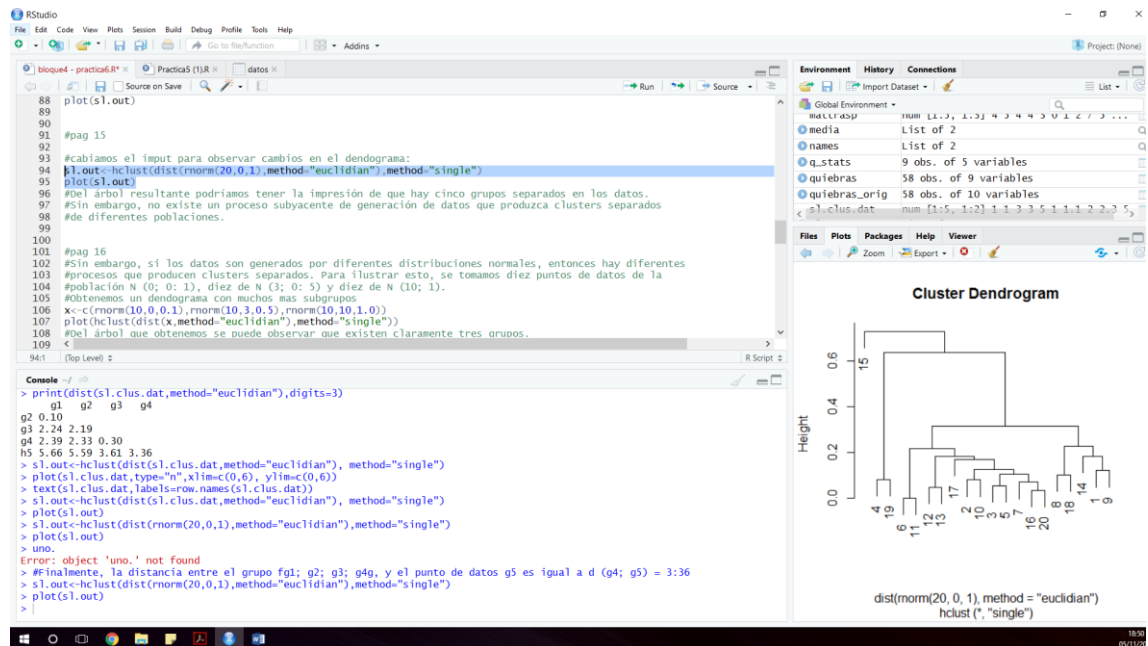
Global Environment

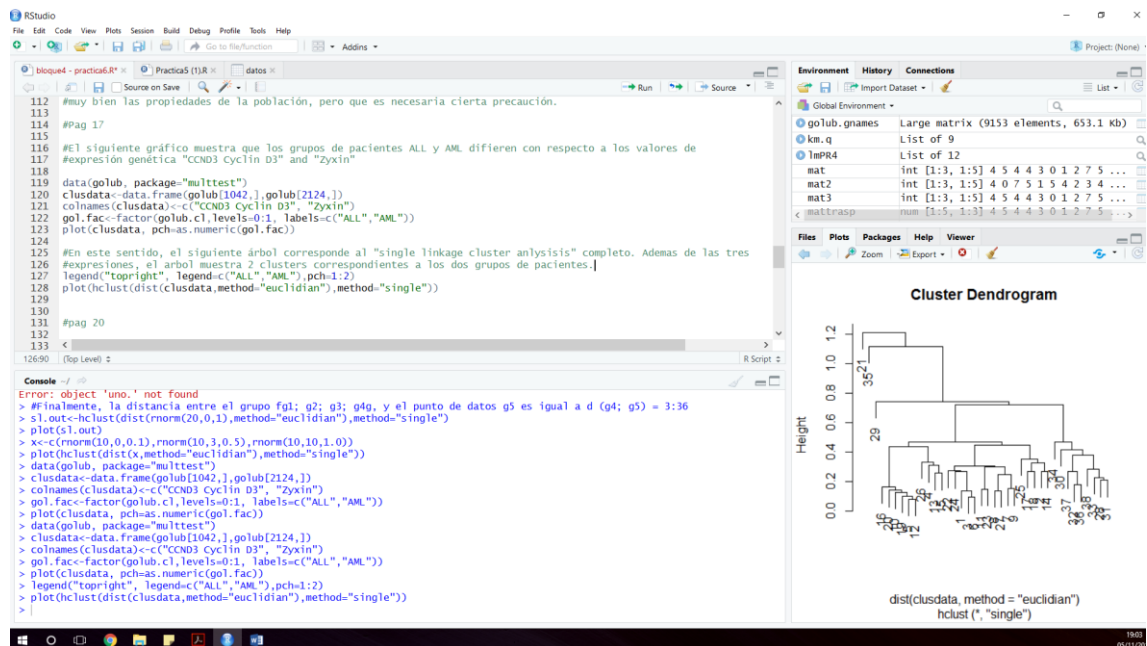
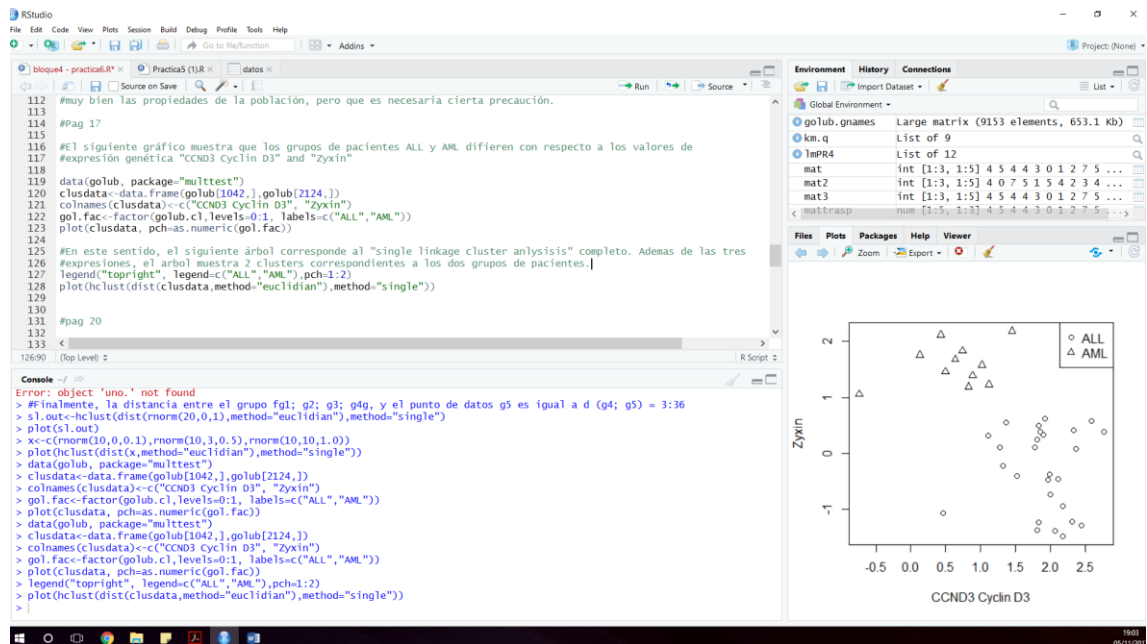
Data

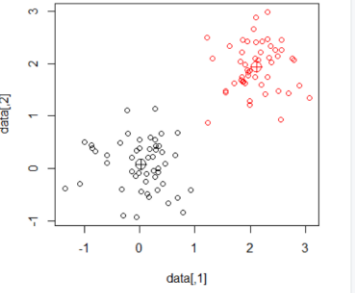
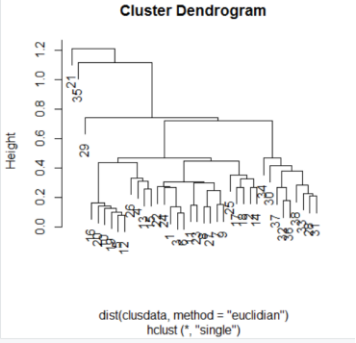
- ALL: Large ExpressionSet (858.1 kb)
- c1: List of 9
- c1389\_at: List of 1
- clusdata: 38 obs. of 2 variables
- data: num [1:100, 1:2] 0.09438 -0.67112 0.00...
- datos: 200 obs. of 42 variables

Files Plots Packages Help Viewer

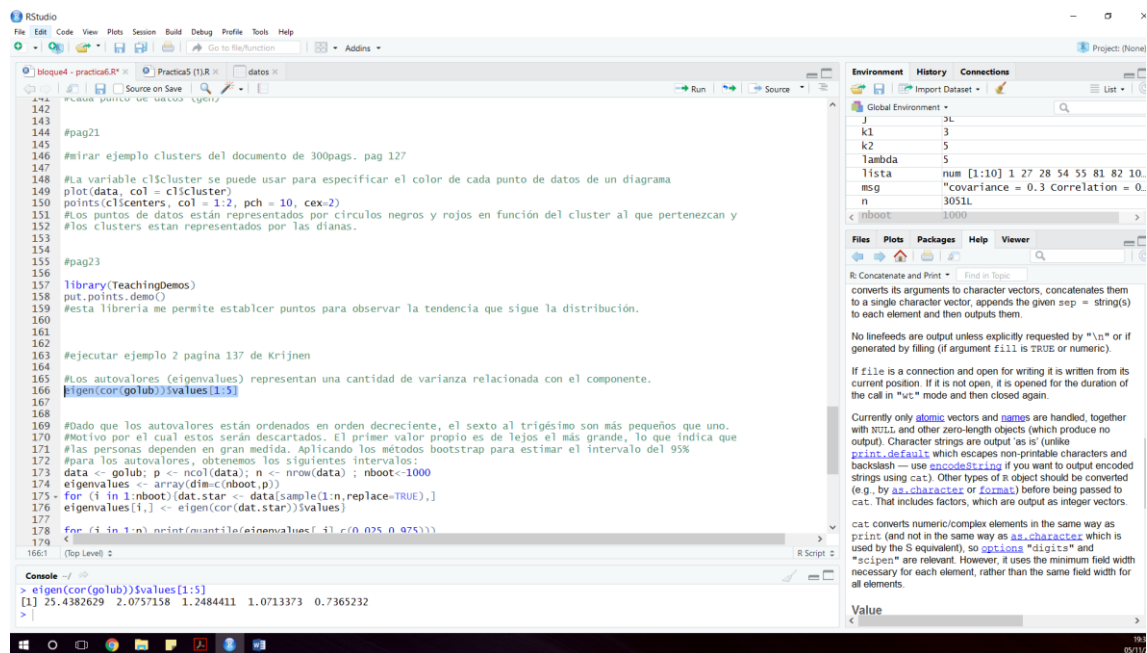
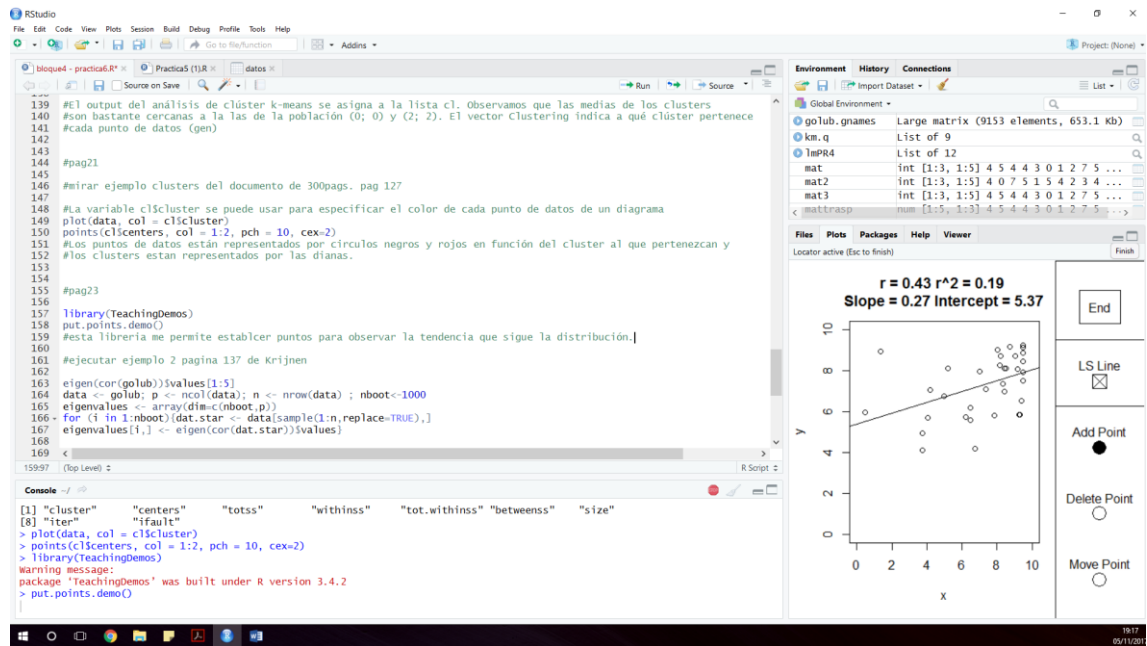












RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

bloque4 - practica6.R | Practica5 (1)R | datos

```

158 put.points(demo())
159 #esta libreria me permite establecer puntos para observar la tendencia que sigue la distribución.
160
161
162
163 #ejecutar ejemplo 2 pagina 137 de Krijnen
164
165 #los autovalores (eigenvalues) representan una cantidad de varianza relacionada con el componente.
166 eigen(cor(golub))$values[1:5]
167
168
169 #Dado que los autovalores están ordenados en orden decreciente, el sexto al trigésimo son más pequeños que uno.
170 #Motivo por el cual estos serán descartados. El primer valor propio es de lejos el más grande, lo que indica que
171 #las personas dependen en gran medida. Aplicando los métodos bootstrap para estimar el intervalo del 95%
172 #para los autovalores, obtenemos los siguientes intervalos:
173 data <- golub; p <- ncol(data); R <- nrow(data); nboot<-1000
174 eigenvalues <- array(dim=c(nboot,p))
175 for (i in 1:nboot){dat.star <- data[sample(1:n,replace=TRUE),]}
176 eigenvalues[i,] <- eigen(cor(dat.star))$values
177
178 for (j in 1:p) print(quantile(eigenvalues[,j],c(0.025,0.975)))
179
180 for (j in 1:5) cat(j,as.numeric(quantile(eigenvalues[,j], #cat da salida a los objetos y concatena las representaciones
181 + c(0.025,0.975))),"\n")
182
183
184
173:1 (Top Level)

```

Console

```

2.5% 97.5%
0.09279025 0.10724545
2.5% 97.5%
0.08250041 0.09815758
>
+
1 24.83642 26.06715
2 1.925606 2.244856
3 1.138249 1.388624
4 0.987678 1.144923
5 0.6865587 0.801799
>

```

Environment History Connections

Global Environment

J	3L
k1	3
k2	5
lambda	5
lista	num [1:10] 1 27 28 54 55 81 82 10
msg	"covariance = 0.3 Correlation = 0.
n	3051L
nboot	1000

Files Plots Packages Help Viewer

R: Concatenate and Print

converts its arguments to character vectors, concatenates them to a single character vector, appends the given sep = string(s) to each element and then outputs them.

No linefeeds are output unless explicitly requested by "\n" or if generated by filling (if argument fill is TRUE or numeric).

If file is a connection and open for writing it is written from its current position. If it is not open, it is opened for the duration of the call in "wt" mode and then closed again.

Currently only **atomic** vectors and **names** are handled, together with **NULL** and other zero-length objects (which produce no output). Character strings are output 'as is' (unlike **print.default** which escapes non-printable characters and backslash — use **encodeString** if you want to output encoded strings using **cat**). Other types of R object should be converted (e.g. by **as.character** or **format**) before being passed to **cat**. That includes factors, which are output as integer vectors.

**cat** converts numeric/complex elements in the same way as **print** (and not in the same way as **as.character** which is used by the S equivalent), so **options** "digits" and "scientific" are relevant. However, it uses the minimum field width necessary for each element, rather than the same field width for all elements.

Value

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

bloque4 - practica6.R | Practica5 (1)R | datos

```

176 eigenvalues[i,] <- eigen(cor(dat.star))$values
177
178 for (j in 1:p) print(quantile(eigenvalues[,j],c(0.025,0.975)))
179
180 for (j in 1:5) cat(j,as.numeric(quantile(eigenvalues[,j], #cat da salida a los objetos y concatena las representaciones
181 + c(0.025,0.975))),"\n")
182
183
184
185 #Bloque 4 tema 2
186 #pagina 148 de Krijnen, chapter 8
187
188 library(ROCR)
189 gol.true <- factor(golub.c1,levels=0:1,labels=c("TRUE","FALSE")) #aquí medimos si se parece al gen "ciclina" el resto de registros
190 pred <- prediction(golub[1042,],gol.true)
191 perf <- performance(pred,"tpr","fpr")
192 plot(perf) #nos sale que el test es de muy buena calidad, cuanto mas pegado a la zona izquierda y superior, mejor.
193
194 data(golub, package = "multtest")
195 gol.true <- factor(golub.c1,levels=0:1,labels= c(" ALL","not ALL"))
196 gol.pred <- factor(golub[1042,]>1.27,levels=c("TRUE","FALSE"), #si el factor sale con valor >1.27, TRUE, que se corresponde con ALL
197 labels=c("ALL","notALL"))
198 table(gol.pred,gol.true) #en el resultado: 27 pacientes de leucemia ALL, 25 bien clasificados, 2 falsos negativos.
199 #De los que no tienen ALL, 10 se han clasificado bien y uno es un falso negativo
200
201 #En este sentido, la true positive rate es 0.93 (25/27) y la false positive rate es 0.09 (1/11)
202
192:116 (Top Level)

```

Console

```

> gol.true <- factor(golub.c1,levels=0:1,labels= c(" ALL","not ALL"))
> gol.pred <- factor(golub[1042,]>1.27,levels=c("TRUE","FALSE"), #si el factor sale con valor >1.27, TRUE, que se corresponde con ALL
+ labels=c("ALL","notALL"))
> table(gol.pred,gol.true)
gol.true
gol.pred ALL not ALL
ALL 25 1
notALL 2 10
> gol.true <- factor(golub.c1,levels=0:1,labels=c("TRUE","FALSE")) #aquí medimos si se parece al gen "ciclina" el resto de registros
> pred <- prediction(golub[1042,],gol.true)
> perf <- performance(pred,"tpr","fpr")
> plot(perf)
+

```

Environment History Connections

Global Environment

perf	Formal class performance
pred	Formal class prediction
q.stats	9 obs. of 5 variables
quiebras	58 obs. of 9 variables
quiebras	58 obs. of 10 variables
s1.clus.d	num [1:5, 1:2] 1 1 3 3 5
s1.out	List of 7

Files Plots Packages Help Viewer