

Programación Estadística: Gráficas

Jocelyn Simmonds (jsimmond@dcc.uchile.cl)

Departamento de Ciencias de la Computación

Cambiar el directorio de trabajo

Pueden usar `getwd()` y `setwd()` para cambiar el directorio en el que están trabajando. En este ejemplo, asuman esta estructura de directorios:

R-workspace

```
|--- clase01  
|--- clase02  
|--- clase03  
|--- clase04  
|--- .....
```

```
1 > getwd()  
2 [1] "/home/jsimmond/R-workspace/clase04"  
3 > setwd("../clase03")  
4 > getwd()  
5 [1] "/home/jsimmond/R-workspace/clase03"  
6 > setwd("../")  
7 > getwd()  
8 [1] "/home/jsimmond/R-workspace"  
9 > setwd("../clase04")  
10 > getwd()  
11 [1] "/home/jsimmond/R-workspace/clase04"
```

Dos nombres de directorios especiales:

- `.` corresponde al directorio actual
- `..` corresponde al directorio “padre” del directorio actual

Graficando datos cualitativos

Datos cualitativos

Hemos visto que podemos usar factores para categorizar datos:

```
1 > head(mtcars, 5)
2           mpg  cyl  disp  hp  drat    wt  qsec vs  am  gear  carb
3 Mazda RX4      21.0   6  160  110 3.90 2.620 16.46 0  1   4    4
4 Mazda RX4 Wag  21.0   6  160  110 3.90 2.875 17.02 0  1   4    4
5 Datsun 710     22.8   4  108   93 3.85 2.320 18.61 1  1   4    1
6 Hornet 4 Drive 21.4   6  258  110 3.08 3.215 19.44 1  0   3    1
7 Hornet Sportabout 18.7   8  360  175 3.15 3.440 17.02 0  0   3    2
8 > mtcars$cyl
9  [1] 6 6 4 6 8 6 8 4 4 6 6 8 8 8 8 8 8 4 4 4 4 8 8 8 8 4 4 4 8 6 8 4
10 >
11 > fact_cyl <- factor(mtcars$cyl)
12 > fact_cyl
13  [1] 6 6 4 6 8 6 8 4 4 6 6 8 8 8 8 8 8 4 4 4 4 8 8 8 8 4 4 4 8 6 8 4
14 Levels: 4 6 8
```

Usen `help(nombre_data_set)` para ver la descripción de un dataset disponible en R.

Pruébenlo

El paquete MASS contiene varios datasets, como `painters` y `crabs`. Para usar estos datasets, primero deben cargar el paquete MASS:

```
1 > library(MASS)
```

Exploren los datasets `painters` y `crabs` del paquete MASS. ¿Cuáles columnas son factores, o candidatos a factores?

Distribución de frecuencia

Pueden usar la función `table()` para calcular la frecuencia de ocurrencia de cada categoría de un factor:

```
1 > table(painters$School)
2  A  B  C  D  E  F  G  H
3 10  6  6 10  7  4  7  4
4 >
5 > table(fact_cyl)
6 fact_cyl
7  4  6  8
8 11  7 14
9 > table(crabs$sex, crabs$sp)
10
11      B  O
12  F 50 50
13  M 50 50
```

Estadísticas sobre categorías

`table()` solo cuenta ocurrencias, deben usar `tapply()` para calcular estadísticas sobre cruces de categorías:

```
1 > tapply(crabs$FL, list(crabs$sp, crabs$sex), mean)
2           F           M
3 B 13.270 14.842
4 O 17.594 16.626
5 >
6 > tapply(painters$Composition, painters$School, mean)
7           A           B           C           D           E           F           G           H
8 10.400000 12.166667 13.166667  9.100000 13.571429  7.250000 13.857143 14.000000
```

Distribución de frecuencia relativa

$$\text{Frecuencia relativa} = \frac{\text{Frecuencia}}{\text{Tamaño muestra}}$$

Calculamos la distribución de frecuencia relativa usando `table()` y `nrow()`:

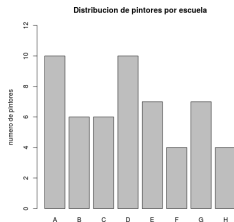
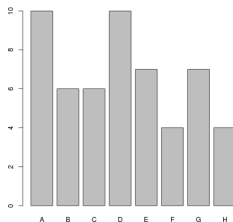
```
1 > freq_painters <- table(painters$School)
2 > freq_painters
3
4   A  B  C  D  E  F  G  H
5 10  6  6 10  7  4  7  4
6 > rel_freq <- freq_painters/nrow(painters)
7 > rel_freq
8
9           A           B           C           D           E           F           G
10 0.18518519 0.11111111 0.11111111 0.18518519 0.12962963 0.07407407 0.12962963
11           H
12 0.07407407
```

Cambiando las opciones de R

```
1 > rel_freq
2
3           A           B           C           D           E           F           G
4 0.18518519 0.11111111 0.11111111 0.18518519 0.12962963 0.07407407 0.12962963
5           H
6 0.07407407
7 > options(digits=1)
8 > rel_freq
9
10          A    B    C    D    E    F    G    H
11 0.19 0.11 0.11 0.19 0.13 0.07 0.13 0.07
12 > options(digits=3)
13 > rel_freq
14
15          A    B    C    D    E    F    G    H
16 0.1852 0.1111 0.1111 0.1852 0.1296 0.0741 0.1296 0.0741
```

Graficando la distribución de frecuencia

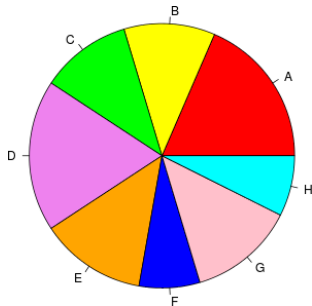
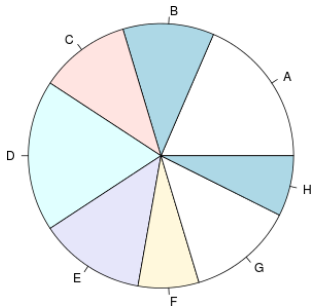
```
1 > freq_painters <- table(painters$School)
2 > barplot(freq_painters)           # grafica en pantalla
3 >
4 > png("mi_archivo.png")          # dejara el grafico en mi_archivo.png
5 > barplot(freq_painters)
6 > dev.off()                       # cierra el "device", copiando los
7 >                                 # contenidos al archivo
8 >
9 > x11()                            # o windows() si estan en Windows
10 > barplot(freq_painters, ylab = "numero de pintores", ylim = c(0,12),
11 + main="Distribucion de pintores por escuela")
```



Graficando la distribución de frecuencia

También podemos generar gráficos de torta:

```
1 > pie(freq_painters)
2 >
3 > # con otros colores
4 > colores = c("red", "yellow", "green", "violet", "orange", "blue",
5 + "pink", "cyan")
6 > pie(freq_painters, col = colores)
```



Otra forma de guardar gráficos

Pueden guardar un gráfico generado en pantalla en forma directa, usando `dev.copy()`:

```
1 > pie(freq_painters, col = colores)
2 > dev.copy(png, "piechart.png")
3 png
4 4
5 > dev.off()
6 png
7 2
```

Al hacer `dev.copy()`, cambian el valor del device actual al indicado.

Graficando datos cuantitativos

Distribución de frecuencia

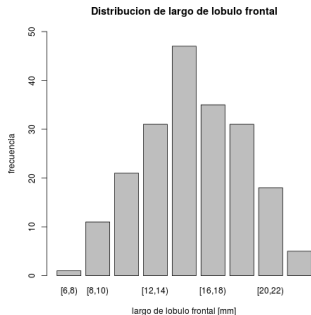
Primero debemos generar una partición del rango de los datos:

```
1 > head(crabs$FL, 20)      # tamaño del lobulo frontal [mm]
2 [1]  8.1  8.8  9.2  9.6  9.8 10.8 11.1 11.6 11.8 11.8 12.2 12.3 12.6 12.8 12.8
3 [16] 12.9 13.1 13.1 13.3 13.9
4 >
5 > range(crabs$FL)
6 [1]  7.2 23.1
7 >
8 > cortes <- seq(6.0, 24.0, by=2)
9 > length(cortes)
10 [1] 10
11 > crabs.cut.FL <- cut(crabs$FL, cortes, right=FALSE)
12 > crabs.cut.FL
13 [1] [8,10) [8,10) [8,10) [8,10) [8,10) [10,12) [10,12) [10,12) [10,12)
14 [10] [10,12) [12,14) [12,14) [12,14) [12,14) [12,14) [12,14) [12,14) [12,14)
15 [19] [12,14) [12,14) [14,16) [14,16) [14,16) [14,16) [14,16) [14,16) [14,16)
16 ...
17 [190] [20,22) [20,22) [20,22) [20,22) [20,22) [20,22) [20,22) [20,22) [20,22)
18 [199] [22,24) [22,24)
19 9 Levels: [6,8) [8,10) [10,12) [12,14) [14,16) [16,18) [18,20) ... [22,24)
```

Distribución de frecuencia

Ahora podemos usar `table()` y `barplot()`:

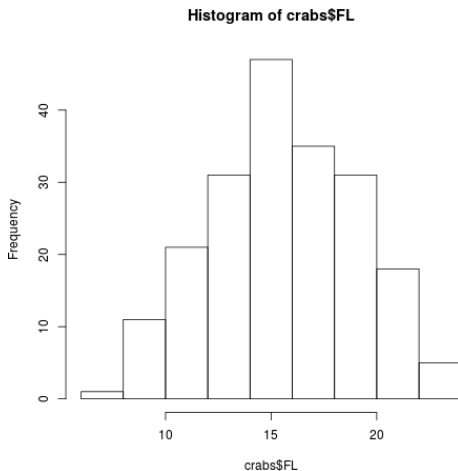
```
1 > freq.FL <- table(crabs.cut.FL)
2 > freq.FL
3 crabs.cut.FL
4   [6,8) [8,10) [10,12) [12,14) [14,16) [16,18) [18,20) [20,22) [22,24)
5         1      11      21      31      47      35      31      18       5
6 > barplot(freq.FL, ylim=c(0, 50), main="Distribucion de largo de lobulo
7 + frontal", ylab="frecuencia", xlab="largo de lobulo frontal [mm]")
```



Distribución de frecuencia

También podemos usar la función `hist()` para generar histogramas:

```
1 > hist(crabs$FL, right=FALSE)
```



Distribución de frecuencia relativa y acumulada

Podemos reutilizar los intervalos que calculamos antes usando `cut()` para ahora calcular las distribuciones de frecuencia relativa y acumulada:

```
1 > freq.FL <- table(crabs.cut.FL)
2 >
3 > relfreq.FL <- freq.FL / nrow(crabs)
4 > relfreq.FL
5 crabs.cut.FL
6 [6,8) [8,10) [10,12) [12,14) [14,16) [16,18) [18,20) [20,22) [22,24)
7 0.005 0.055 0.105 0.155 0.235 0.175 0.155 0.090 0.025
8 >
9 > acumfreq.FL <- cumsum(freq.FL)
10 > acumfreq.FL
11 [6,8) [8,10) [10,12) [12,14) [14,16) [16,18) [18,20) [20,22) [22,24)
12 1 12 33 64 111 146 177 195 200
```

Generando un gráfico de la frecuencia acumulada

```
1 > acumfreq0.FL <- c(0, cumsum(freq.FL))
2 > plot(cortes, acumfreq0.FL,
3 +   main="Lobulo frontal del cangrejo Leptograpsus variegatus",
4 +   xlab="tamano [mm]", ylab="frecuencia acumulada")
5 > lines(cortes, acumfreq0.FL)
```

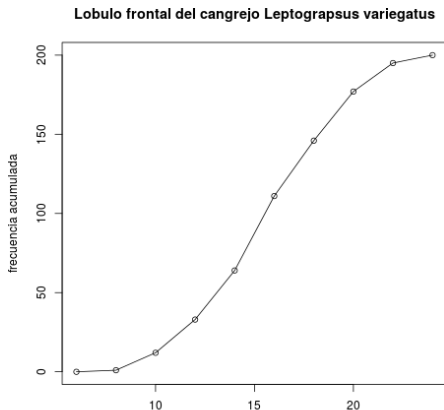


Gráfico de dispersión (scatter plot)

Para ver la relación entre dos variables numéricas, podemos usar un gráfico de dispersión:

```
1 > plot(crabs$FL, crabs$CL, xlab="tamano del lobulo frontal [mm]",  
2 + ylab="tamano del caparazon [mm]", main="cangrejo Leptograpsus variegatus")
```

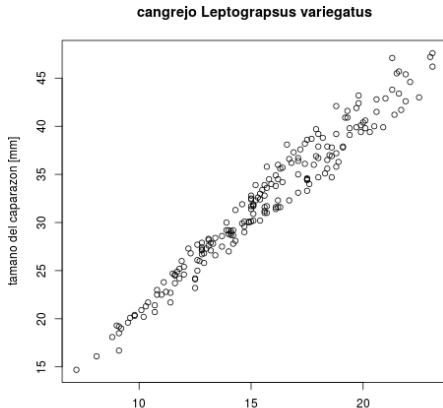
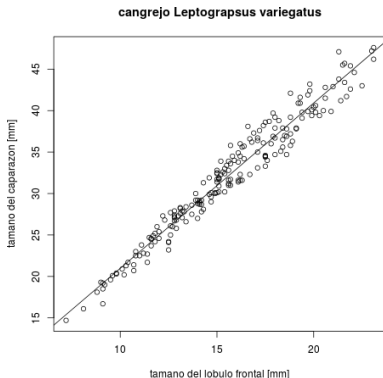


Gráfico de dispersión (scatter plot)

En este caso, hay una clara relación lineal positiva entre las dos variables. Podemos usar `abline()` y `lm()` para agregar una línea de tendencia:

```
1 > plot(crabs$FL, crabs$CL, xlab="tamano del lobulo frontal [mm]",  
2 + ylab="tamano del caparazon [mm]", main="cangrejo Leptograpsus variegatus")  
3 > abline(lm(crabs$CL ~ crabs$FL))
```



Algunos estadísticas adicionales

Cuartiles

Podemos dividir los datos en cuartiles:

- primer cuartil (cuartil inferior o Q1): la mediana de la primera mitad de los datos
- segundo cuartil (el cuartil medio o Q2): es la mediana
- tercer cuartil (cuartil superior o Q3): la mediana de la segunda mitad de los datos

```
1 > quantile(crabs$FL)
2   0%   25%   50%   75%  100%
3  7.20 12.90 15.55 18.05 23.10
```

El rango intercuartil se define como $Q3 - Q1$:

```
1 > IQR(crabs$FL)
2 [1] 5.15
```

Percentiles

Los percentiles son los 99 valores que dividen la serie de datos en 100 partes iguales. Por ejemplo:

- P_{10} es el valor que separa el primer 10% de los valores de los datos, cuando se ordenan en forma ascendente
- P_{50} es la mediana

```
1 > quantile(crabs$FL, c(.10, .25, .67, .89))
2   10%   25%   67%   89%
3 10.980 12.900 17.400 20.011
```

Gráfico de caja (box-and-whiskers)

```
1 > boxplot(crabs$FL, horizontal=TRUE, xlab="distribucion de tamaño de  
2 + lobulo frontal [mm]", main="cangrejo Leptograpsus variegatus")
```

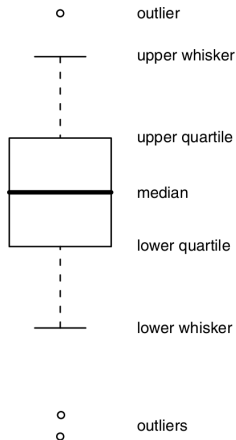
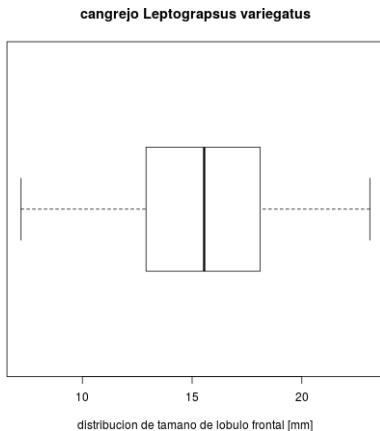
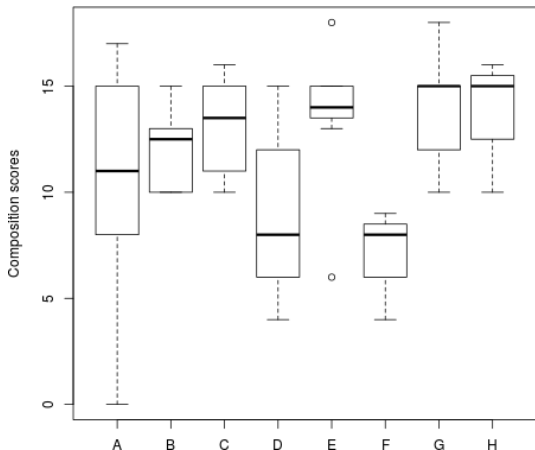


Gráfico de caja (box-and-whiskers)

También pueden mostrar distribuciones por alguna categoría:

```
1 > boxplot(Composition ~ School, data=painters, ylab="Composition scores")
```



Covarianza y el coeficiente de correlación

- La covarianza de dos variables da una idea de cuanto varían estas variables en forma conjunta.
- El coeficiente de correlación indica la magnitud de la correlación: positiva (valor cercano a 1), negativa (valor cercano a -1) o sin relación (valor cercano a 0).

```
1 > cov(crabs$FL, crabs$CL)
2 [1] 24.356677
3 >
4 > cor(crabs$FL, crabs$CL)
5 [1] 0.97884179
```

Podemos concluir que hay una correlación positiva entre el tamaño del caparazón y del lóbulo frontal en el cangrejo *Leptograpsus variegatus*.