

MÉTODOS DE REMUESTREO

Tema 1. Introducción a los Métodos de Remuestreo

basado en

- B. Efron, R. Tibshirani (1993). An Introduction to the bootstrap.
O. Kirchkamp (2017). Resampling methods.

Curso 2018/19

Introducción

- ▶ La Estadística es una ciencia puramente experimental que aprende de la experiencia, especialmente de la que aparece de manera *secuencial*.
- ▶ Se puede considerar que comienza alrededor del siglo XVII cuando **Graunt** (UK) empezó a usar **tablas de mortalidad**.
- ▶ Actualmente se aplica en todas las ciencias que requieren tratamiento de la información, desde ciencias biomédicas, psicología, educación y economía, hasta el estudio de las partículas en física cuántica, o de galaxias extremadamente distantes.

Introducción

- ▶ Pero la mayoría de las personas y los dispositivos no son muy eficientes para encontrar patrones en un mar de datos con *ruido*.
- ▶ Es decir, tenemos tendencia a contemplar patrones inexistentes fuera que suceden para satisfacer nuestros propósitos: astrología...
- ▶ La teoría estadística intenta proporcionar métodos óptimos para la búsqueda de señales reales en ambientes ruidosos y también proporciona controles estrictos contra la *sobre-interpretación* de patrones al azar.

Introducción

La teoría estadística trata de responder a tres preguntas básicas:

1. ¿Cómo debo recoger o muestrear mis datos?
2. ¿Cómo debo analizar y resumir los datos que he recogido?
3. ¿Cómo son de exactos los resúmenes de mis datos?

El punto 3 constituye parte del proceso conocido como *Inferencia Estadística*.

Introducción

- ▶ Las técnicas de remuestreo son técnicas desarrolladas hace pocos años para calcular estadísticos, basándose en técnicas computacionales intensivas que evitan los cálculos complejos de la teoría estadística tradicional.
- ▶ Aunque el uso de las técnicas de remuestreo implican el uso de conceptos tradicionales de inferencia estadística, cambia radicalmente su implementación.
- ▶ Mediante el uso de computación intensiva se aplican las técnicas de manera flexible, fácil y con un mínimo de aparato matemático.

Ejemplo

- ▶ Los tres conceptos básicos: recogida de datos, resumen de los mismos e inferencia se ilustran en la noticia aparecida en el *New York Times*.
- ▶ Se hizo un estudio sobre si pequeñas dosis de aspirina podrían prevenir los ataques al corazón en personas de mediana edad.
- ▶ Los datos del estudios se tomaron de manera eficiente mediante un estudio controlado, aleatorizado y *doble ciego*.

HEART ATTACK RISK FOUND TO BE CUT BY TAKING ASPIRIN

LIFESAVING EFFECTS SEEN

Study Finds Benefit of Tablet Every Other Day Is Much Greater Than Expected

By HAROLD M. SCHMECK Jr.

A major nationwide study shows that a single aspirin tablet every other day can sharply reduce a man's risk of heart attack and death from heart attack.

The lifesaving effects were so dramatic that the study was halted in mid-December so that the results could be reported as soon as possible to the participants and to the medical profession in general.

The magnitude of the beneficial effect was far greater than expected, Dr. Charles H. Hennekens of Harvard, principal investigator in the research, said in a telephone interview. The risk of myocardial infarction, the technical name for heart attack, was cut almost in half.

'Extreme Beneficial Effect'

A special report said the results showed "a statistically extreme beneficial effect" from the use of aspirin. The report is to be published Thursday in The New England Journal of Medicine.

In recent years smaller studies have demonstrated that a person who has had one heart attack can reduce the risk of a second by taking aspirin, but there had been no proof that the beneficial effect would extend to the general male population.

Dr. Claude Lenfant, the director of the National Heart Lung and Blood Institute, said the findings were "extremely important," but he said the general public should not take the report as an indication that everyone should start taking aspirin.

Ejemplo

- ▶ La mitad de las personas recibió una sustancia placebo y las personas se asignaron de manera aleatoria a los tratamientos.
- ▶ Tanto los sujetos como las personas que trabajaban en el estudio no sabían la identificación de tanto el tipo de las dosis como de los pacientes.
- ▶ Los estadísticos de resumen del artículo eran muy simples:

	Ataques Corazón	No Ataque Corazón
Aspirina	104	11037
Placebo	189	11034

Ejemplo

- ▶ A simple vista parece que hay *menos* ataques de corazón en el grupo de aspirina.
- ▶ La *razón de odds* (la razón o ratio entre ambos ratios de ataques la corazón) es

$$\hat{\theta} = \frac{104/11037}{189/11034} = 0,55$$

Ver

es.wikipedia.org/wiki/Razón_de_momios

- ▶ Según este estudio, las personas que toman aspirinas tienen casi la mitad de riesgo de sufrir un ataque al corazón.

Ejemplo

- ▶ Realmente $\hat{\theta}$ es solo un estimador del valor **poblacional** desconocido.
- ▶ La muestra parece suficientemente grande en el estudio: 22071, pero la conclusión de que la aspirina funciona bien se basa en **solo** 293 casos observados de ataques al corazón.
- ▶ ¿Se puede asegurar que se obtendría el mismo resultado si tomamos otra muestra distinta?
- ▶ Usando estadística clásica, aproximando mediante la distribución normal, se obtiene que un intervalo al 95 % para el verdadero valor θ es $0,43 < \theta < 0,70$
- ▶ Parece obvio que nunca excede el valor de 1, es decir que la aspirina en cualquier caso **no** es significativamente mala para la salud.

Razón de odds y distribución asintótica

- ▶ Dada una tabla de contingencia 2×2

	C	D
A	n_{11}	n_{12}
B	n_{21}	n_{22}

- ▶ El estimador que se utiliza para la razón de odds es

$$\hat{\theta} = \frac{n_{11}n_{22}}{n_{12}n_{21}}$$

- ▶ La distribución de este estimador es muy **asimétrica** por lo que para considerar una aproximación a la normal es mejor tomar la transformación $\log(\hat{\theta})$.

Razón de odds y distribución asintótica

- Una estimación del error estándar de $\log(\hat{\theta})$ (aplicando el método *delta*) es

$$\hat{\sigma}_{\log(\hat{\theta})} = \sqrt{\frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}}$$

de modo que el correspondiente intervalo de Wald es

$$\log(\hat{\theta}) \pm z_{\frac{\alpha}{2}} \hat{\sigma}_{\log(\hat{\theta})}$$

- Si se toman la función exponencial (*antilogaritmo*) de los extremos se obtiene el intervalo correspondiente para θ .
- El test es algo **conservador** (la probabilidad de cubrimiento es algo mayor que el nivel nominal).

Ejemplo sobre aspirina

```
(teta = (104/11037)/(189/11034))
```

```
[1] 0.550115
```

```
SE = sqrt(1/104 + 1/11037 + 1/189 + 1/11034)
```

```
LSup = log(teta) + qnorm(1-0.05/2)*SE
```

```
LInf = log(teta) - qnorm(1-0.05/2)*SE
```

```
exp(LInf)
```

```
exp(LSup)
```

```
[1] 0.4324113
```

```
[1] 0.699858
```

Ejemplo sobre aspirina

- ▶ Por otro lado, en un estudio sobre accidentes cerebrovasculares (*ictus*) se observó

	Ictus	No Ictus
Aspirina	119	11037
Placebo	98	11034

$$\hat{\theta} = \frac{119/11037}{98/11034} = 1,21$$

- ▶ de modo que un intervalo al 95 % para el verdadero valor es θ es
 $0,93 < \theta < 1,59$
- ▶ Este intervalo incluye al valor *neutro* 1 en el que se obtiene que la aspirina no es significativamente ni mejor ni peor que el placebo.

Bootstrap

- ▶ Se podría decir que la aspirina es *significativamente* beneficiosa para prevenir ataques al corazón pero no es significativamente perjudicial para causar accidentes cerebrovasculares. Pero es un tema que sigue siendo controvertido.
- ▶ El bootstrap es un método de simulación que se puede usar en casos como el anterior.
- ▶ El término *bootstrap* deriva de la frase *to pull oneself up by one's bootstrap*. Se basa en el libro del siglo XVIII *las aventuras del barón Munchausen* de Rudolph E. Raspe.
- ▶ El barón había caído en el fondo de un profundo lago y se le ocurrió escapar *tirando* de los cordones de sus propias botas...



Ideas de bootstrap en el ejemplo de aspirina

- ▶ En el caso del bootstrap, se consideran dos poblaciones:
 - ▶ la primera con 11037 observaciones contiene 119 *unos* y 10918 *ceros*
 - ▶ la segunda con 11034 observaciones contiene 98 *unos* y 10936 *ceros*.
- ▶ Se genera una muestra **con reemplazamiento** de 11037 elementos de la primera población, y una muestra de 11034 elementos de la segunda población.
- ▶ Se calcula la proporción respectiva de *unos* en ambas muestras:

$$\hat{\theta}^* = \frac{\text{Proporción de unos en la muestra 1}}{\text{Proporción de unos en la muestra 2}}$$

Bootstrap

- ▶ Se repite este proceso un número elevado de veces, digamos $N = 1000$, y se obtiene una muestra de 1000 valores de $\hat{\theta}^*$.
- ▶ Este proceso es fácil de implementar con cualquier software estadístico como R o MatLab, por ejemplo.
- ▶ Esta muestra de 1000 valores de $\hat{\theta}^*$ contiene información que se puede usar para realizar inferencia a partir de los datos reales.
- ▶ Por ejemplo, Efron y Tibshirani obtienen una desviación estándar muestral alrededor de 0,17 y un intervalo de confianza basado en los cuantiles muestrales alrededor de (0,93; 1,60).

Bootstrap

- ▶ Basta tomar el elemento de la muestra que ocupa el lugar 25 y el que ocupa el lugar 975 una vez ordenada la muestra de los 1000 valores de $\hat{\theta}^*$.
- ▶ Este resultado es similar al obtenido mediante métodos clásicos de inferencia (basados en la aproximación a la distribución normal).
- ▶ Pero aquí **NO** se ha usado ningún argumento de aproximación del tipo *Teorema Central del Límite* como en los métodos clásicos.

Ejemplo sobre aspirina

```
n1 = 11037 + 104      # tamaño muestra 1
s1 = 104               # numero de exitos

n2 = 11034 + 189      # tamaño muestra 2
s2 = 189               # numero de exitos

p1pre = c(rep(1,s1), rep(0,n1-s1))
p2pre = c(rep(1,s2), rep(0,n2-s2))

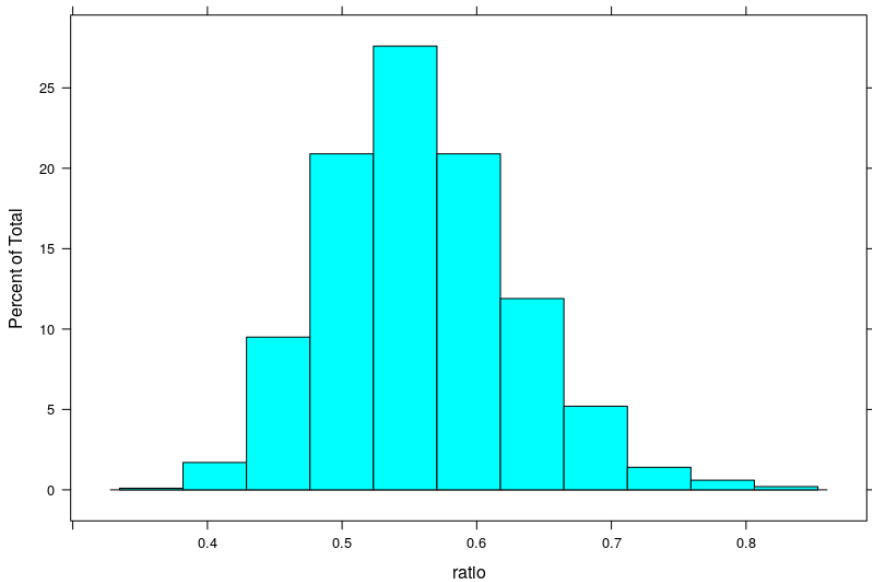
p1 = sample(p1pre, n1)  # muestra 1
p2 = sample(p2pre, n2)  # muestra 2

n.bs = 1000            # tomo n.bs muestras bootstrap

# reservo dos vectores de ceros
bs1 = rep(0, n.bs)
bs2 = rep(0, n.bs)
```

Ejemplo sobre aspirina

```
for (i in 1:n.bs) {  
  # proporcion de exitos en muestras bootstrap 1 y 2  
  bs1[i] = sum(sample(p1, n1, replace=TRUE))/n1  
  bs2[i] = sum(sample(p2, n2, replace=TRUE))/n2  
}  
  
# replicas de la estimacion bootstrap del ratio  
ratio = bs1/bs2  
  
# histograma de las estimaciones del ratio  
lattice::histogram(ratio)
```



Ejemplo de aspirina

```
mean(ratio)
median(ratio)
```

```
[1] 0.5571126
[1] 0.5539163
```

```
# IC bootstrap son los cuantiles del 0.025 y
# 0.975 de la muestra ordenada
quantile(ratio, probs=c(0.025,0.975))
```

```
      2.5%      97.5%
0.4330636 0.7086892
```

Ejemplo de aspirina

```
# 0 bien a mano:
# ordenas las estimas del ratio para obtener
# los IC bootstrap
  rats = sort(ratio)

# IC bootstrap son los cuantiles del 0.025 y
# 0.975 de la muestra ordenada
  CI.bs = c(rats[round(0.025*n.bs)],
            rats[round(0.975*n.bs)])
  CI.bs
```

```
[1] 0.4310098 0.7086662
```


Ejemplo con Rcpp

- Se escribe un programa en C++ y se graba e.j. en el fichero denominado `boot_ratio2prop.cpp`:

```
#include <Rcpp.h>
using namespace Rcpp ;
// [[Rcpp::export]]

NumericVector boot_ratio2prop(NumericVector p1,
NumericVector p2, int replicas=1000) {

    int n1 = p1.size();
    int n2 = p2.size();

    NumericVector bs1(replicas);
    NumericVector bs2(replicas);
    NumericVector ratio(replicas);
```

Ejemplo con Rcpp

```
bool replace = true;

for(int i=0; i<replicas; i++) {
    bs1[i] = sum(sample(p1, n1, replace))/n1;
    bs2[i] = sum(sample(p2, n2, replace))/n2;
    ratio[i] = bs1[i]/bs2[i];
}

return ratio;
}
```

Ejemplo con Rcpp

- ▶ Se usa la interfaz de R con C++ mediante [Rcpp](#).

```
library(Rcpp)

# Se compila el programa escrito en C++
sourceCpp("boot_ratio2prop.cpp")

# Se ejecuta el programa desde R
replica = 2000

sale = boot_ratio2prop(p1, p2, replica)

lattice::histogram(sale)
```

Estudio de la media

- ▶ Por ejemplo, se puede considerar el caso de la media muestral, estudiando su precisión como estimador de la media poblacional.
- ▶ **Ejemplo:** Ratones bajo tratamiento o no para prolongar su supervivencia después de cirugía invasiva.

Tratamiento	94	197	16	38	99	141	23		
Control	52	104	146	10	51	30	40	27	46

- ▶ **Método:** comparación de las medias de tiempos de vida entre el tratamiento y el control
- ▶ Las medias de cada grupo son
 1. Tratamiento: 86.86
 2. Control: 56.22

Ejemplo ratones

- ▶ La diferencia entre ambas medias $\bar{x} - \bar{y} = 30,63$ sugiere que hay **bastante efecto** del tratamiento.
- ▶ ¿Cómo son de precisos estos estimadores?
Resulta que las muestras son muy *pequeñas*...
- ▶ En el caso de la media, el error estándar estimado basado en n valores de una *m.a.s.* x_1, x_2, \dots, x_n es

$$\sqrt{\frac{s^2}{n}}$$

donde

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Ejemplo ratones

- ▶ Si el error estándar de las estimas de las medias es pequeño, entonces la diferencia entre las medias de supervivencia sería significativa.
- ▶ Los errores estándar de x e y son $SE_x = 25,24$ y $SE_y = 14,14$
- ▶ El error estándar para la diferencia $\bar{x} - \bar{y}$ es igual a $\sqrt{25,24^2 + 14,14^2} = 28,93$
- ▶ Pero la diferencia observada es 30.63. Si se calcula el estadístico de contraste: $30,63/28,93 = 1,05$, es decir es **mucho menor** que el valor 1.96 típico de la distribución normal.
- ▶ Esto implica que el tratamiento **NO** tiene efectos significativos en la supervivencia de los ratones.

Ejemplo de ratones

```
Trata = c(94, 197, 16, 38, 99, 141, 23)
Cont = c(52, 104, 146, 10, 51, 30, 40, 27, 46)

# Grafico stem-and-leaf de grupo tratamiento
stem(Trata, scale=2)
```

The decimal point is 1 digit(s) to the right of the |

```
0 | 6
2 | 38
4 |
6 |
8 | 49
10 |
12 |
14 | 1
16 |
18 | 7
```

Ejemplo de ratones

```
# Grafico stem-and-leaf de grupo control  
stem(Cont,scale=2)
```

The decimal point is 1 digit(s) to the right of the |

```
0 | 0  
2 | 70  
4 | 0612  
6 |  
8 |  
10 | 4  
12 |  
14 | 6
```

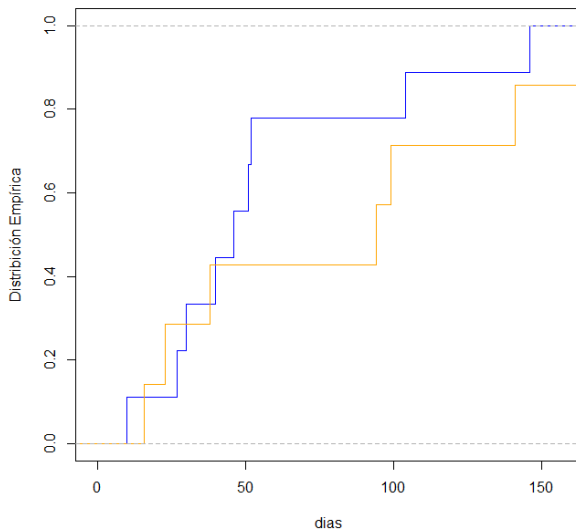

Ejemplo de ratones

```
ecdf1 = ecdf(Trata)
ecdf2 = ecdf(Cont)

X11()
plot(ecdf2, verticals=TRUE, do.points=FALSE, col='blue',
     xlab="dias",
     ylab="Distribucion Empirica",
     main="Cont(azul) / Trata(naranja)" )

plot(ecdf1, verticals=TRUE, do.points=FALSE, add=TRUE,
     col='orange')
```

Cont(azul) / Trata(naranja)



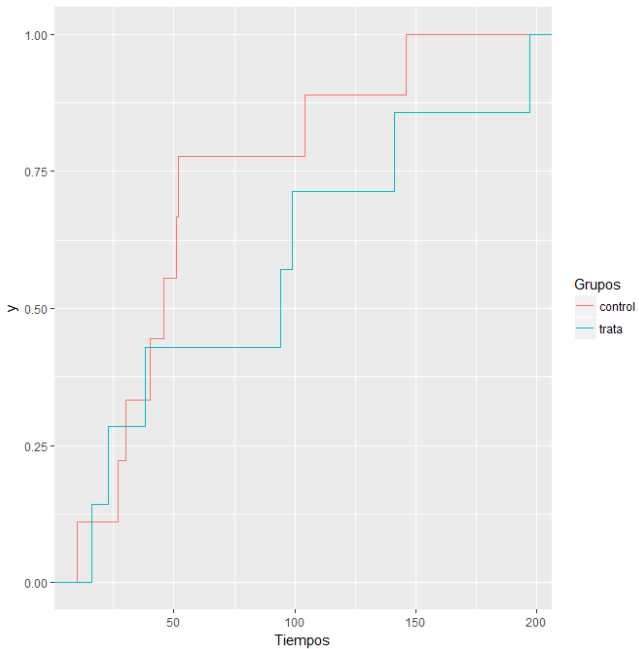
Ejemplo de ratones

Gráfico alternativo con ggplot

```
library(ggplot2)
Tiempos = c(Trata, Cont)
Grupos = as.factor(c(rep("trata", length(Trata)),
rep("control", length(Cont))))

Losdatos = data.frame(Tiempos, Grupos)

ggplot(Losdatos, aes(Tiempos, colour = Grupos)) +
  stat_ecdf()
```



Ejemplo de ratones

```
# Siguiendo a Efron y Tibshirani:  
mean(Trata)  
mean(Cont)
```

```
[1] 86.85714  
[1] 56.22222
```

```
(sdDiff =  
sqrt(var(Trata)/length(Trata)+var(Cont)/length(Cont)))
```

```
[1] 28.93607
```

```
# Aplicas el TCL  
(t = (mean(Trata) - mean(Cont)) / sdDiff)
```

```
[1] 1.058711
```

```
2*(1-pnorm(t))
```

```
[1] 0.2897316
```

Ejemplo de ratones

```
# Pones los valores en un solo vector
# Defines otro vector de 1's y 2's segun su grupo
x = matrix(c(Trata, Cont, rep(1,length(Trata)),
rep(2, length(Cont))), ncol=2)

# t-test de Student
t.test(x[,1] ~ x[,2])
```

Welch Two Sample t-test

```
data:  x[, 1] by x[, 2]
t = 1.0587, df = 9.6545, p-value = 0.3155
alternative hypothesis: true difference in means is not equal to 0

95 percent confidence interval:
 -34.15279  95.42263
sample estimates:
mean in group 1 mean in group 2
    86.85714      56.22222
```

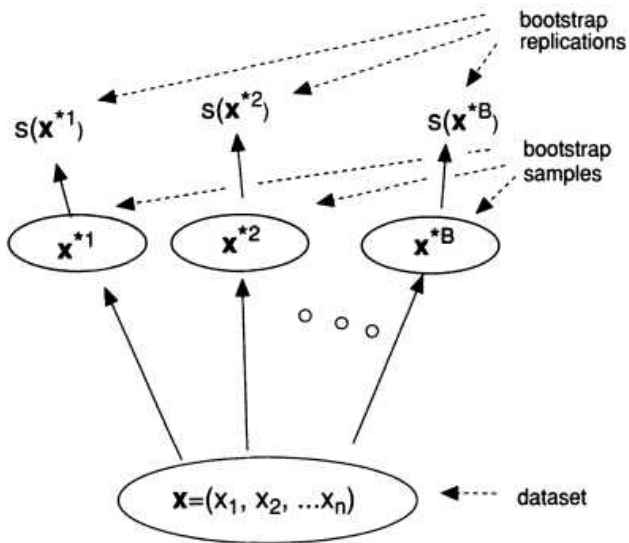
Cálculo de medianas

- ▶ El cálculo de los errores estándar es un método para estudiar la precisión de los estimadores. Pero, salvo en el caso de la media, no hay en general fórmulas concretas para estimarlos.
- ▶ Por ejemplo, si se calculan las medianas en el ejemplo de los ratones, se tiene que la mediana para el tratamiento es 94 y para el control es 46.
- ▶ La diferencia entre medianas es 48, mucho **mayor** que la diferencia entre medias.
- ▶ Pero, ¿cómo es de precisa esta estimación?
La única manera de responder es usando bootstrap en este caso.

Esquema del Bootstrap

- ▶ Supongamos que se observa una muestra $\mathbf{x} = x_1, x_2, \dots, x_n$ sobre la que se calcula un cierto estadístico $s(\mathbf{x})$.
- ▶ Por ejemplo, \mathbf{x} es el grupo control de observaciones y $s(\mathbf{x})$ es la media muestral.
- ▶ En el caso del bootstrap, se define una *muestra bootstrap* $\mathbf{x}^* = (x_1^*, x_2^*, \dots, x_n^*)$ que se obtiene muestreando n veces **con reemplazamiento** a partir de los datos originales (x_1, x_2, \dots, x_n)
- ▶ Por ejemplo si $n = 7$ una posible muestra bootstrap podría ser

$$\mathbf{x}^* = (x_5, x_7, x_5, x_4, x_7, x_3, x_1)$$



Algoritmo del Bootstrap

- ▶ Se genera un número B elevado de muestras bootstrap $\mathbf{x}^{*1}, \mathbf{x}^{*2}, \dots, \mathbf{x}^{*B}$ cada una de tamaño n . Los tamaños típicos de B para errores estándar suelen estar entre 500 y 5000.
- ▶ Para cada muestra bootstrap $b = 1, \dots, B$ se calcula el estadístico $s(\mathbf{x}^{*b})$. Por ejemplo, la *mediana*.
- ▶ El estimador bootstrap del error estándar es la desviación estándar de las B muestras bootstrap

$$\widehat{se}_{Boot} = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (s(\mathbf{x}^{*b}) - s(\cdot))^2}$$

donde

$$s(\cdot) = \frac{1}{B} \sum_{b=1}^B s(\mathbf{x}^{*b})$$

Ejemplo con la mediana

- ▶ Por ejemplo, Efron y Tibshirani obtienen en el grupo control un error estándar igual a 11.54 y en el grupo tratamiento, 36.35 basándose en $B = 100$ réplicas.
- ▶ Así, usando ese número de réplicas la diferencia de medianas observada igual a 48 se obtiene un error estándar estimado igual a $\sqrt{36,352 + 11,542} = 38,14$
- ▶ Así el estadístico es $48/38,14 = 1,26$ que tampoco es significativamente mayor que 0, aunque es mayor que en el caso de la media.
- ▶ Para la mayoría de los estadísticos no existen fórmulas explícitas que sirvan para calcular el error estándar, por ello se puede usar un procedimiento bootstrap.

Ejemplo de ratones

Se calcula el error estándar del estadístico *media muestral*.

En el caso de la *mediana muestral*, sería semejante sustituyendo el comando `mean` por `median` en el programa.

```
# Tamaños muestrales
n1 = length(Trata)
n2 = length(Cont)

n.bs = 1000    # numero de muestras bootstrap

bs1 = rep(0,n.bs)
bs2 = rep(0,n.bs)

for (i in 1:n.bs) {
  bs1[i] = mean(sample(Trata, n1, replace=TRUE))
  bs2[i] = mean(sample(Cont, n2, replace=TRUE))
}
```

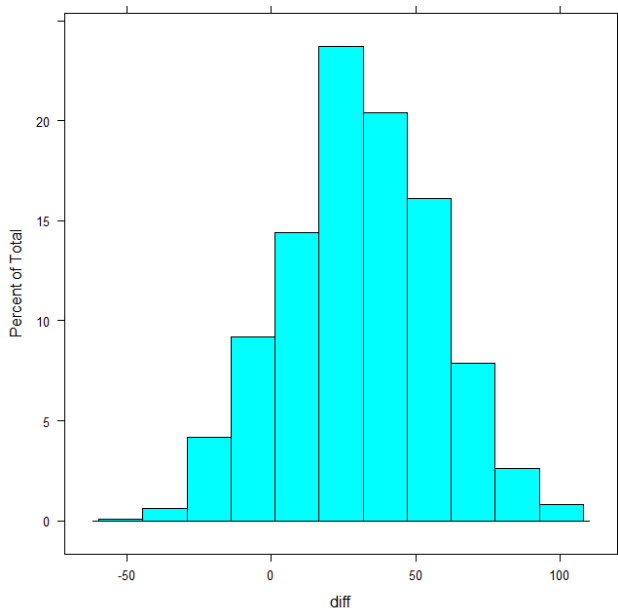
Ejemplo de ratones

```
# replicas bootstrap de estimadores de las diferencias  
diff = bs1-bs2
```

```
sd(diff) # estima error estandar
```

```
[1] 26.1389
```

```
# histograma de estimas de las diferencias  
lattice::histogram(diff)
```



Ejemplo de ratones con Rcpp

- Se escribe un programa en C++ y se graba e.j. en el fichero denominado `boot_MediaRatones.cpp`:

```
#include <Rcpp.h>
using namespace Rcpp ;
// [[Rcpp::export]]

List boot_MediaRatones(NumericVector trata,
NumericVector control, int replicas=1000) {

    int n1 = trata.size();
    int n2 = control.size();
    double sale;

    NumericVector bs1(replicas);
    NumericVector bs2(replicas);
    NumericVector diff(replicas);

    bool replace = true;
```

Ejemplo de ratones con Rcpp

```
for(int i=0; i<replicas; i++) {  
  
    bs1[i] =  
    mean(sample(trata, n1, replace));  
  
    bs2[i] =  
    mean(sample(control, n2, replace));  
  
    diff[i] = bs1[i]-bs2[i];  
}  
  
sale = sd(diff);  
  
List saletodo;  
saletodo["sd"] = sale;  
saletodo["vector"] = diff;  
return saletodo;  
}
```


Ejemplo de ratones con Rcpp

- ▶ Se usa la interfaz de R con C++ mediante **Rcpp** y se puede obtener un histograma semejante al anterior

```
library(Rcpp)

# Se compila el programa escrito en C++
sourceCpp("boot_MediaRatones.cpp")

# Se ejecuta el programa desde R
replica = 2000

sale = boot_MediaRatones(Trata, Cont, replica)

sale["sd"]

lattice::histogram(unlist(sale["vector"]),
main="Distribución error estándar",xlab="")
```

Ejemplo de ratones: Programas alternativos

```
B = 2000

# Caso de la media
mean(Trata) - mean(Cont)
sd(replicate(B, mean(sample(Trata, replace=TRUE)) -
mean(sample(Cont, replace=TRUE))))

# Caso de la mediana
median(Trata) - median(Cont)
sd(replicate(B, median(sample(Trata, replace=TRUE)) -
median(sample(Cont, replace=TRUE))))
```

En el caso de usar [Rcpp](#), el programa para la mediana sería semejante al de la media, escribiendo `median` en lugar de `mean` en el programa.

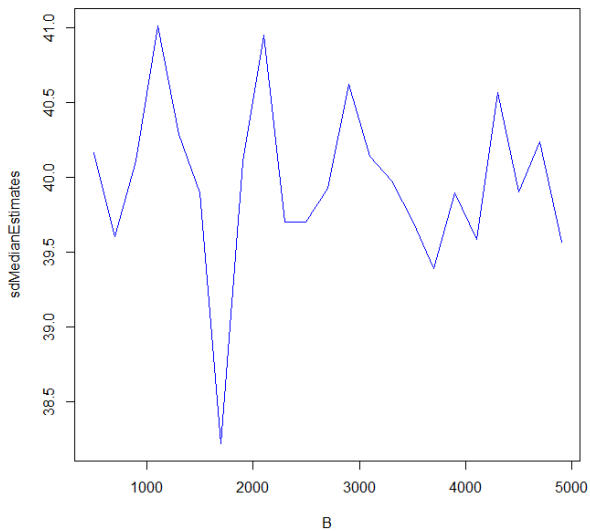
Ejemplo de ratones

```
# Programa para las medianas con distintos tamaños
# de muestras bootstrap

sdMedian = function(B){
  sd(replicate(B,
    median(sample(Trata, replace=TRUE)) -
    median(sample(Cont, replace=TRUE))))}

# Se prueban diferentes numeros de replicas
B = seq(500, 5000, 200)
sdMedianEstimates = sapply(B, sdMedian)

X11()
plot(sdMedianEstimates ~ B, type="l", col="blue")
```



Programas para el ejemplo de los ratones

- Se puede programar fácilmente el ejemplo de los ratones con las librerías `bootstrap` y `boot`.

```
library(bootstrap)

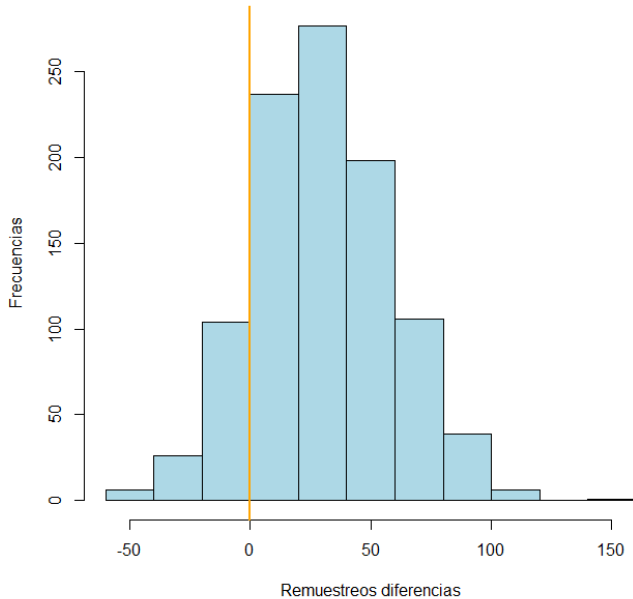
mouse.boot.c = bootstrap(mouse.c, 1000, mean)
mouse.boot.t = bootstrap(mouse.t, 1000, mean)

mouse.boot.diff = mouse.boot.t$thetastar -
mouse.boot.c$thetastar
hist(mouse.boot.diff, main='Histograma remuestreos',
col='lightblue', xlab='Remuestreos diferencias',
ylab='Frecuencias')
abline(v=0, col="orange", lwd=2)

sd(mouse.boot.diff)
```

```
[1] 28.04734
```

Histograma remuestreos



Programas para el ejemplo de los ratones

- ▶ Como introducción al manejo de la librería boot, se pueden consultar las siguientes páginas Web:

```
www.mayin.org/ajayshah/KB/R/documents/boot.html
```

```
cran.r-project.org/doc/Rnews/Rnews\_2002-3.pdf
```

```
ww2.coastal.edu/kingw/statistics/R-tutorials/resample.html
```

Programas para el ejemplo de los ratones

```
library(boot)
trat = c(94, 197, 16, 38, 99, 141, 23)
control = c(52, 104, 146, 10, 51, 30, 40, 27, 46)

bichos = data.frame(sobrevive=c(control, trat),
grupo=c(rep(1, length(control)), rep(2, length(trat))))

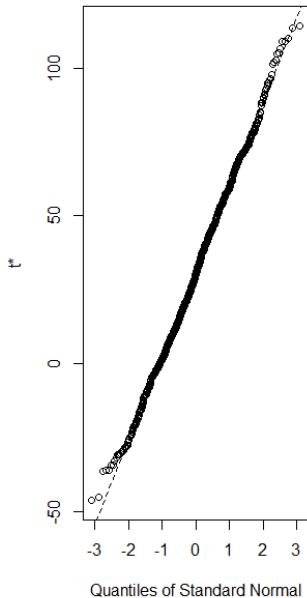
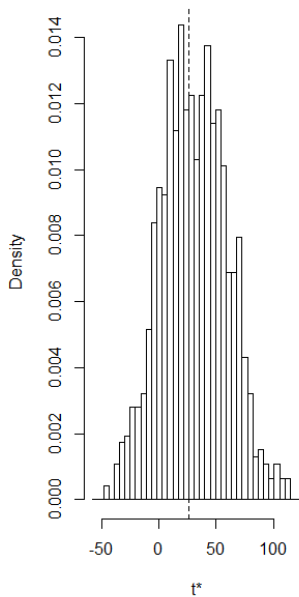
lafuncion = function(x, i) {
booty = tapply(x$sobrevive, x$grupo,
FUN=function(x) sample(x, length(x), TRUE))
diff(sapply(booty, mean)) }

bbichos = boot(data=bichos, lafuncion, R=1000)
plot(bbichos)
```

```
sd(bbichos$t)
```

```
[1] 28.60951
```


Histogram of t



Jackknife

- ▶ El **jackknife** es un método propuesto por Quenouille sobre 1950 originalmente para estimar errores estándar y sesgos.
- ▶ Dado un conjunto de datos $\mathbf{x} = (x_1, x_2, \dots, x_n)$ la muestra i -ésima $\mathbf{x}_{(i)}$ se define como el conjunto de datos original menos la observación i -ésima

$$\mathbf{x}_{(i)} = (x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$$

para $i = 1, 2, \dots, n$.

- ▶ La réplica i -ésima $\hat{\theta}_{(i)}$ del estadístico $\theta = s(\mathbf{x})$ es el valor de $s(\cdot)$ evaluado en $\mathbf{x}_{(i)}$

$$\hat{\theta}_{(i)} = s(\mathbf{x}_{(i)})$$

para $i = 1, 2, \dots, n$.