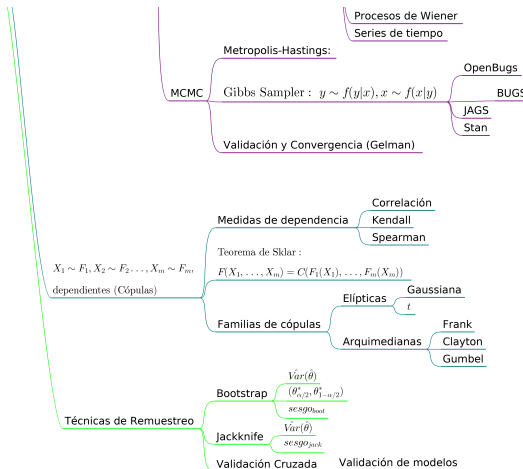


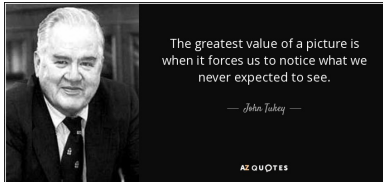




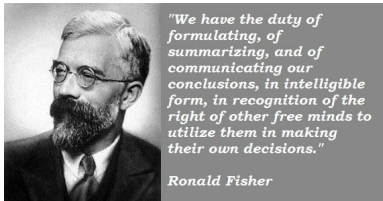
# ¡Ya acabamos!



# Introducción



**Figura:** John Tukey, 1915-2000.



**Figura:** Ronald Avery Fisher, 1890-1962.

- El Remuestreo se refiere a un conjunto de técnicas estadísticas, computacionalmente intensivas y no paramétricas, que estiman la distribución de una población basadas en muestreo aleatorio *con reemplazo*.
- Se considera que una muestra aleatoria  $\{X_1, X_2, \dots, X_n\}$  como si fuera una *población* finita y se generan muestras aleatorias de la misma muestra para estimar características poblacionales y hacer inferencia de la población muestreada.
- Las técnicas de remuestreo permiten asignar *medidas de ajuste* (en términos de sesgo, varianza, intervalos de confianza, errores de predicción o de algunas otras medidas) a los estimados basados en muestras.
- Estas técnicas son usualmente no paramétricas, y varias son tan antiguas como la estadística misma. Por ejemplo, las técnicas de permutación son de Fisher (1935) y Pitmann (1937); la validación cruzadas fue propuesta por Kurtz en 1948, y el Jackknife fue propuesto por Maurice Quenouille en 1949 y John Tukey en 1958 fue quien le dio nombre a la técnica.













# Bootstrap



## Simulación

## Ejemplo de las diferencias I

Seguimos con el ejemplo de las diferencias, podemos encontrar la distribución de las muestras bootstrap para la mediana  $\theta = q_{.5}$ :

```
d <- c(3,5,-3,6)
quantile(d,0.5)

50%
4

Boot <- NULL
B <- 300
Boot <- matrix(0, nrow=B, ncol=4)
for(i in 1:B) Boot[i,] <- sample(d,replace = T)
head(Boot)
```

	[,1]	[,2]	[,3]	[,4]
[1,]	-3	3	5	3
[2,]	-3	3	-3	3
[3,]	-3	-3	5	5
[4,]	5	3	6	5
[5,]	5	5	3	3
[6,]	6	5	5	-3

## Ejemplo de las diferencias II

```
medianaboot <- apply(Boot, 1, quantile, 0.5)
mean(medianaboot)

[1] 3.418333

sd(medianaboot) #un estimador de la dispersión de la mediana

[1] 2.201967

hist(medianaboot, prob=T, main="Distribución empírica de la muestra bootstrap para la
mediana", breaks=30)
abline(v=quantile(d, 0.5), lwd=3, col="red")
```











## Ejemplo I

## Correlación (Efron & Tibshirani)

De una población de  $N = 82$  escuelas americanas de leyes, se toma una muestra de tamaño  $n = 15$ . Se miden los resultados promedios de dos scores: LSAT (promedio de evaluaciones de la prueba de admisión a Leyes) y GPA (promedio de licenciatura).

Los datos están en la variable `law` que se obtiene al cargar el paquete `bootstrap`. En este ejercicio el parámetro de interés es el coeficiente de correlación  $\rho$ . A partir de la muestra obtenemos el coeficiente de correlación muestral,  $\hat{\rho} = 0.776$ . Noten también que la correlación muestral es un estimador *plug-in*.

```
library(bootstrap)
data("law")
plot(law, main="Relación entre LSAT y GPA en 15 escuelas")
```



## Correlación (Efron & Tibshirani)

```
B <- 400
n <- nrow(law) #tamaño de la muestra, 15
R <- numeric(B)

for (b in 1:B){
  i <- sample(1:n,size=n, replace=T)
  LSAT <- law$LSAT[i]
  GPA <- law$GPA[i]
  R[b] <- cor(LSAT,GPA)
}

sd(R)

[1] 0.1409769

hist(R,prob=T,breaks=30)
lines(density(R),col="red",lwd=2)
```



## Correlación (Efron & Tibshirani)

El estimador bootstrap del error estándar de  $\hat{\rho}$  es 0.1409769 y claramente su distribución no es normal. La teoría dice que el estimador basado en normalidad del error estándar de  $\hat{\rho}$  es

$$\hat{\sigma}(\hat{\rho}) = \frac{1 - \hat{\rho}^2}{\sqrt{n - 3}} = 0.115$$







# Estimación bootstrap del sesgo II

## Estimador bootstrap de sesgo

$$\widehat{sesgo}_{boot}(\hat{\theta}) = \frac{1}{B} \sum_{i=1}^B \hat{\theta}^{(b)} - \hat{\theta}$$

donde  $\hat{\theta}$  es el estimador del parámetro con la muestra original.



# Ejemplo I

## precios y distribución de la media recortada

Los datos siguientes corresponden a una muestra de precios de venta de propiedades en Seattle en 2002. Los datos no distinguen entre propiedades residenciales que son la mayoría, pero hay algunas comerciales en la muestra, lo cual puede incrementar el precio promedio de venta muestral. Una estadística más resistente a valores extremos es la media recortada.

```
precios <- c(142, 175, 197.5, 149.4, 705, 232, 50, 146.5, 155, 1850,
            132.5, 215, 116.7, 244.9, 290, 200, 260, 449.9, 66.407, 164.95,
            362, 307, 266, 166, 375, 244.95, 210.95, 265, 296, 335,
            335, 1370, 256, 148.5, 987.5, 324.5, 215.5, 684.5, 270, 330,
            222, 179.8, 257, 252.95, 149.95, 225, 217, 570, 507, 190)

mean(precios)

[1] 329.2571

mean(precios,trim=0.25)

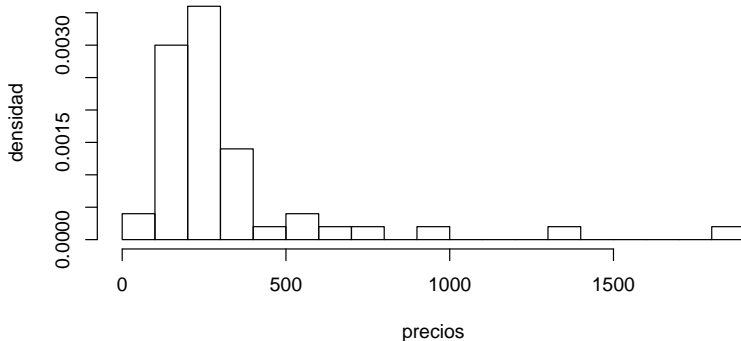
[1] 244.0019

hist(precios, probability=T,main="histograma de los precios", ylab="densidad",breaks=20)
```

## Ejemplo II

precios y distribución de la media recortada

histograma de los precios



# Ejemplo I

## precios y distribución de la media recortada

- ¿Qué podemos decir de la distribución muestral de  $\bar{x}_{25}\%$ ? No mucho. Pero podemos estimar lo que necesitamos con bootstrap.

```
media.recortada <- function(x,i){mean(x[i],trim=0.25)}
boot1 <- boot(data = precios, statistic = media.recortada, R=1000)
boot1
```

ORDINARY NONPARAMETRIC BOOTSTRAP

Call:

```
boot(data = precios, statistic = media.recortada, R = 1000)
```

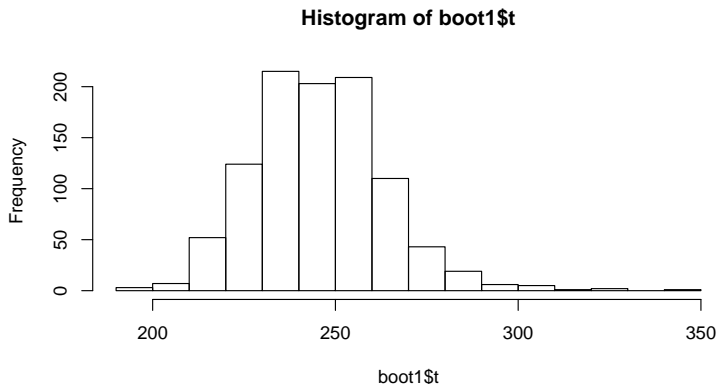
Bootstrap Statistics :

	original	bias	std. error
t1*	244.0019	1.752694	18.02073

```
hist(boot1$t) #gráfica de las replicas bootstrap
```

## Ejemplo II

precios y distribución de la media recortada



- En el resultado, se puede concluir que la forma de la distribución de la media recortada es parecida a la normal.

## Ejemplo III

precios y distribución de la media recortada

- La estimación del sesgo es 1.7526942, que es un sesgo pequeño relativo al tamaño de la estadística. La estadística (la media recortada de la muestra) tiene un sesgo pequeño como un estimado del parámetro (que es la media recortada de la población).
- El estimador bootstrap del error estándar es 18.0207255.





# Jackknife

- El *jackknife* es un método de remuestreo propuesto por Quenouille (1949) para estimar sesgo, y por Tukey (1958) para estimar error estándar.
- *Jackknife* es un tipo de *validación cruzada* en donde se deja un dato fuera a la vez: si  $\mathbf{x} = (x_1, \dots, x_n)$  es una muestra aleatoria, hay  $n$  muestras jackknife  $\mathbf{x}_{-i}$  que es la muestra original sin la  $i$ -ésima observación.
- Si  $\hat{\theta} = T_n(\mathbf{x})$ , entonces la réplica  $i$ -ésima jackknife es  $\hat{\theta}_{-i} = T_{n-1}(\mathbf{x}_{-i})$ .
- Se puede generalizar a considerar  $g$  grupos de tamaño  $h$  tales que  $gh = n$  y se calcula la  $i$ -ésima réplica jackknife dejando uno de los  $g$  grupos afuera.

# Estimador jackknife del sesgo I

## Estimador Jackknife del sesgo

El estimador jackknife del sesgo es

$$sesgo_{jack} = (n - 1) \left( \frac{\sum_{i=1}^n \hat{\theta}_{-i}}{n} - \hat{\theta} \right)$$

y el estimador corregido por sesgo es:

$$\hat{\theta}_{jack} = \hat{\theta} - sesgo_{jack} = n\hat{\theta} - \frac{n-1}{n} \sum_{i=1}^n \hat{\theta}_{-i}$$

## Ejemplo jackknife

En nuestro ejemplo inicial,  $\hat{\theta} = \bar{d} = 2.75$ , y el sesgo está dado por:

```
d <- c(6, -3, 5, 3)
dadj <- NULL #guarda los valores ajustados
for(i in 1:4) dadj[i] <- mean(d[-i])
sesgo <- 3*(mean(dadj)-mean(d))
sesgo

[1] 0

dadj

[1] 1.666667 4.666667 2.000000 2.666667
```

## Ejemplo sesgo Jackknife I

Consideren una muestra aleatoria  $X_1, \dots, X_n \sim \mathbf{Bernoulli}(\theta)$ . Queremos estimar  $\theta^2$ .

- El estimador máximo verosímil de  $\theta^2$  es  $\bar{x}^2$ . Sin embargo, este estimador de  $\theta^2$  tiene sesgo, el cuál es fácil de calcular considerando que  $Y = \sum_{i=1}^n X_i \sim \mathbf{Bin}(n, \theta)$  y  $E(Y^2) = n\theta(1 - \theta) + n^2\theta^2$ :

$$E(\hat{\theta}^2) = \theta^2 + \frac{\theta(1 - \theta)}{n}$$

- De acuerdo a las definiciones previas, el estimador jackknife está dado por:

$$\begin{aligned} \hat{\theta}_{jack}^2 = \bar{x}^2 - sesgo_{jack} &= \bar{x}^2 - (n-1) \left( \frac{1}{n} \sum_{i=1}^n \bar{x}_{-i}^2 - \bar{x}^2 \right) \\ &= n\bar{x}^2 - \frac{n-1}{n} \sum_{i=1}^n \bar{x}_{-i}^2 \end{aligned}$$

# Ejemplo sesgo Jackknife II

## ● numéricamente:

```
set.seed(1)
n <- 100 #tamaño de muestra
x <- rbinom(n, size=1, p=0.3)
x

[1] 0 0 0 1 0 1 1 0 0 0 0 0 0 0 1 0 1 1 0 1 1 0 0 0 0 0 0 0 1 0 0 0 0 0 1
[36] 0 1 0 1 0 1 0 1 0 1 0 0 1 0 0 1 0 0 1 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 1 0 1
[71] 0 1 0 0 0 1 1 0 1 1 0 1 0 0 1 0 1 0 0 0 0 0 0 0 1 1 1 0 0 1 0

thetahat2 <- mean(x)^2 #estimador máximo verosímil
thetahat2

[1] 0.1024

xj <- NULL #calcula los estimados jack
for(i in 1:n) xj[i] <- mean(x[-i])^2
sesgo <- (n-1)*(mean(xj)-thetahat2)
sesgo

[1] 0.00219798

#Estimador jack corregido
thetahat.jack <- n*thetahat2-(n-1)*sum(xj)/n
thetahat.jack

[1] 0.100202
```

## Notas al jackknife I

- Noten que si definimos  $\tilde{\theta}_i = n\hat{\theta} - (n-1)\hat{\theta}_{-i}$  como *pseudovalores*, entonces  $\hat{\theta}_{jack}$  es el promedio de estos pseudovalores:

$$\hat{\theta}_{jack} = \frac{1}{n} \sum_{i=1}^n \tilde{\theta}_i$$

En el ejemplo considerado, los pseudovalores son iguales a los valores originales.

- El estimador jackknife tiene un factor  $n-1$ . Para ver su origen, consideren el caso donde  $\theta = \sigma^2$  la varianza de la población. El estimador plug-in de la varianza es:

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

Este estimador es sesgado para  $\sigma^2$  con sesgo:

$$\text{sesgo}(\hat{\theta}) = E(\hat{\theta} - \sigma^2) = \frac{n-1}{n} \sigma^2 - \sigma^2 = -\frac{\sigma^2}{n}.$$

## Notas al jackknife II

Cada replicación jackknife calcula el estimado  $\hat{\theta}_{(i)}$  sobre una muestra de tamaño  $n - 1$ , así que el sesgo en la replicación jackknife es  $-\sigma^2/(n - 1)$  Entonces

$$\begin{aligned}
 E(\hat{\theta}_{(i)} - \hat{\theta}) &= E(\hat{\theta}_{-i} - \theta) - E(\hat{\theta} - \theta) \\
 &= \text{sesgo}(\hat{\theta}_{-i}) - \text{sesgo}(\hat{\theta}) \\
 &= -\frac{\sigma^2}{n-1} - \left(-\frac{\sigma^2}{n}\right) \\
 &= -\frac{\sigma^2}{n(n-1)} = \frac{\text{sesgo}(\hat{\theta})}{n-1}
 \end{aligned}$$

Entonces el estimador jackknife con factor  $n - 1$  da un estimado correcto de sesgo en el estimador plug-in de la varianza.

- El estimador jackknife puede fallar cuando la estadística  $\hat{\theta}$  no es una función “suave” (en el sentido usual del cálculo). Por ejemplo, la mediana  $q_{.5}$  no es una función suave de los datos.



# Estimador jackknife de error estándar

## Estimador Jackknife de varianza

El estimador jackknife del error estándar del estimador es:

$$\hat{se}_{jack} = \sqrt{\frac{n-1}{n} \sum_{i=1}^n \left( \hat{\theta}_{-i} - \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{-i} \right)^2}$$

## Ejemplo sesgo y error estándar jackknife I

Las siguientes corresponden a mediciones de cierta hormona en el torrente sanguíneo de ocho sujetos después de ponerse un parche médico (Efron y Tibshirani). El parámetro de interés es:

$$\theta = \frac{E(new) - E(old)}{E(old) - E(placebo)}$$

Si  $|\theta| \leq 0.20$ , esto indica bioequivalencia entre los parches viejo y nuevo. La estadística es  $\bar{Y}/\bar{Z}$ .

```
data(patch, package="bootstrap")
patch
```

	subject	placebo	oldpatch	newpatch	z	y
1	1	9243	17649	16449	8406	-1200
2	2	9671	12013	14614	2342	2601
3	3	11792	19979	17274	8187	-2705
4	4	13357	21816	23798	8459	1982
5	5	9055	13850	12560	4795	-1290
6	6	6290	9806	10157	3516	351
7	7	12412	17208	16570	4796	-638
8	8	18806	29044	26325	10238	-2719

A continuación calculamos el sesgo y el error estándar jackknife.





## Ejemplo I

Para los datos de los parches, estimamos el error estándar del error estándar para  $\hat{\theta}$ .

```
data(patch, package="bootstrap")
n <- nrow(patch)
y <- patch$y
z <- patch$z
B <- 2000
theta.b <- numeric(B)

#almacenamiento de los indices muestrados
indices <- matrix(0,nrow=B, ncol=n)

#jackknife-after-bootstrap paso 1: corre el bootstrap
for (b in 1:B) {
  i <- sample(1:n, size = n, replace = T)
  y <- patch$y[i]
  z <- patch$z[i]
  theta.b[b] <- mean(y)/mean(z)
  #se guardan los indices para el jackknife
  indices[b, ] <- i
}
```



## Intervalos de confianza bootstrap I

Hay diferentes maneras de calcular intervalos de confianza bootstrap. Varian en su facilidad de cómputo y en su exactitud. Estos se basan en el error estándar bootstrap  $\hat{s}e_{boot} = \sqrt{v_{boot}}$ . Los métodos de intervalos de confianza incluyen:

- bootstrap normal estándar
- bootstrap básico
- bootstrap percentil
- bootstrap  $t$

A continuación revisaremos los supuestos asociados con cada uno de estos casos:

## Bootstrap normal estándar

Usualmente este intervalo no es adecuado a menos que  $\hat{\theta}$  tenga una distribución parecida a la normal.

## Intervalo de confianza bootstrap normal estándar

Si  $\hat{\theta}$  es un estimador del parámetro  $\theta$  y se supone que el error estándar del estimador es  $se(\hat{\theta})$ , entonces si  $\hat{\theta}$  es una media muestral y el tamaño de muestra es grande, el CLT implica que  $Z = \frac{\hat{\theta} - E(\hat{\theta})}{se(\hat{\theta})}$  es aproximadamente normal. Un intervalo aproximado de  $100(1 - \alpha) \%$  de confianza para  $\theta$  es

$$\hat{\theta} \pm z_{\alpha/2} \cdot \hat{se}(\hat{\theta})$$



## Intervalo de confianza bootstrap básico I

Este intervalo transforma la distribución de las replicaciones sustrayendo la estadística observada. Los cuantiles de la muestra transformada se utilizan para determinar los intervalos de confianza:

## Confianza bootstrap básico

El intervalo de  $100(1 - \alpha) \%$  de confianza está dado por:

$$(2\hat{\theta} - \hat{\theta}_{1-\alpha/2}, 2\hat{\theta} - \hat{\theta}_{\alpha/2})$$

En el caso paramétrico, se tiene que si  $T$  es un estimador de  $\theta$ ,

$$P(T - \theta > q_\alpha) = 1 - \alpha \implies P(T - q_\alpha > \theta) = 1 - \alpha$$

Entonces, un intervalo de confianza con  $100(1 - 2\alpha)\%$  con errores  $\alpha$  en las colas inferior y superior iguales está dado por  $(t - q_{1-\alpha}, t - q_\alpha)$ . En el bootstrap la distribución de  $T$  es generalmente desconocida, pero los cuantiles se pueden estimar por un método aproximado:

- 1 Calcular los cuantiles muestrales  $\hat{\theta}_\alpha$  de la función de distribución empírica de las replicaciones  $\hat{\theta}^*$ .

## Intervalo de confianza bootstrap básico II

- ② Denotar los cuantiles  $\alpha$  de  $\hat{\theta}^* - \hat{\theta}$  por  $b_\alpha$ . Entonces  $\hat{b}_\alpha = \hat{\theta}_\alpha - \hat{\theta}$  es un estimador de  $b_\alpha$
- ③ Un límite de confianza superior para el intervalo de  $100(1 - \alpha)\%$  está dado por

$$\hat{\theta} - \hat{b}_{\alpha/2} = \hat{\theta} - (\hat{\theta}_{\alpha/2} - \hat{\theta}) = 2\hat{\theta} - \hat{\theta}_{\alpha/2}.$$

- ④ Similarmente, un límite de confianza inferior para el intervalo de  $100(1 - \alpha)\%$  está dado por

$$\hat{\theta} - \hat{b}_{1-\alpha/2} = \hat{\theta} - (\hat{\theta}_{1-\alpha/2} - \hat{\theta}) = 2\hat{\theta} - \hat{\theta}_{1-\alpha/2}.$$

# Intervalo percentil bootstrap

Los intervalos percentil se basan en la distribución empírica de las réplicas bootstrap como la distribución de referencia .

## Intervalo percentil

A partir de la distribución bootstrap de  $\hat{\theta}^*$  se puede calcular un intervalo de confianza no paramétrico. El  $(1 - \alpha) \times 100\%$  intervalo basado en percentiles es

$$(\hat{\theta}_{(B \cdot \alpha/2)}^*, \hat{\theta}_{(B \cdot (1 - \alpha/2))}^*)$$

basado en los cuantiles  $B \cdot \alpha/2$  y  $B \cdot (1 - \alpha/2)$  de la muestra bootstrap.

- Efron y Tibshirani muestran que los intervalos percentil tienen ventajas teóricas sobre los intervalos estándar normales y un mejor desempeño de cobertura.

## Ejemplo de intervalos de confianza I

Para los datos del experimento del parche, obtendremos intervalos de confianza normal, básico, y percentil usando las funciones de `R boot` y `boot.ci` en el paquete `boot`

```
library(boot)
data(patch, package="bootstrap")
theta.boot <- function(datos, indice) {
  #función para calcular la estadística
  y <- datos[indice,1]
  z <- datos[indice,2]
  mean(y)/mean(z)
}
y <- patch$y #lee los datos
z <- patch$z
datos <- cbind(y,z)
boot.obj <- boot(datos, statistic=theta.boot, R=2000)
```

Los resultados de los intervalos de confianza se dan a continuación:

---

### Calculations and Intervals on Original Scale



## Intervalos $BC_n$ I

## Intervalos $BC_n$

Para calcular un intervalo de confianza de  $100(1 - \alpha)\%$  para  $\theta$ , realizamos los siguientes ajustes:

- 1 Calcular  $z_0 = \Phi^{-1}\left(\frac{\#\{\hat{\theta}_i^* < \hat{\theta}\}}{B}\right)$ . Es el efecto del sesgo, mide el sesgo de la mediana de las replicaciones bootstrap.

- ## 2 Calcular

$$a = \frac{\sum_{b=1}^B (\hat{\theta}_{(i)} - \bar{\theta}_{(\cdot)})^3}{6 \left[ \sum_{b=1}^B (\hat{\theta}_{(i)} - \bar{\theta}_{(\cdot)})^2 \right]^{3/2}}$$

Este cálculo corresponde al cociente del estimador de aceleración.

- 3 Define  $\alpha_1 = \Phi \left[ z_0 + \frac{z_0 - z_{\alpha/2}}{1 - a(z_0 - z_{\alpha/2})} \right]$  y  $\alpha_2 = \Phi \left[ z_0 + \frac{z_0 + z_{\alpha/2}}{1 - a(z_0 + z_{\alpha/2})} \right]$ . El intervalo está dado por  $(\hat{\theta}_{(l^*)}^*, \hat{\theta}_{(u^*)}^*)$ , donde  $l^* = B\alpha_1$  y  $u^* = B\alpha_2$ .

## Ejemplo intervalo $BC_a$

Para los datos del parche,

```
boot.ci(boot.obj, type=c("bca"))
```

```
BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
```

```
Based on 2000 bootstrap replicates
```

```
CALL :
```

```
boot.ci(boot.out = boot.obj, type = c("bca"))
```

```
Intervals :
```

```
Level      BCa
```

```
95%      (-0.2213,  0.1923 )
```

```
Calculations and Intervals on Original Scale
```





## Intervalo bootstrap $t$ . II

- se calcula o estima el error estándar  $\hat{se}(\hat{\theta}^{(b)})$  (un estimado separado para cada muestra bootstrap: esto es, se obtienen bootstraps de la muestra actual  $x^{(b)}$ , no de  $x$ ).
- se calcula la estadística  $t^{(b)} = \frac{\hat{\theta}^{(b)} - \hat{\theta}}{\hat{se}(\hat{\theta}^{(b)})}$
- ③ La muestra de replicadas  $t^{(1)}, \dots, t^{(B)}$  es la distribución de referencia para bootstrap  $t$ . Encontrar los cuantiles  $t_{\alpha/2}^*$  y  $t_{1-\alpha/2}^*$ .
- ④ Calcular  $\hat{se}(\hat{\theta})$ , la desviación estándar muestral de las replicadas  $\hat{\theta}^{(b)}$ .
- ⑤ Calcular  $(\hat{\theta} - t_{1-\alpha/2}^* \hat{se}(\hat{\theta}), \hat{\theta} - t_{\alpha/2}^* \hat{se}(\hat{\theta}))$
- La principal desventaja de este intervalo es que se tiene bootstrap anidado, lo que puede demandar tiempo de cómputo alto.





### Histograma de tiempos de reparación ILEC



Los datos no son normales. Sin embargo, sabemos que la media tiende a distribuirse como una normal.

## Ejemplos adicionales. Media poblacional IV

```
n <- length(ilec)
summary(ilec)

  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.000   0.730   3.590   8.412   7.080  191.600

#Intervalo para la media:
mean(ilec) + c(-1,1)*qt(.975,df = n-1 ,lower.tail = T)*sd(ilec)/sqrt(n)

[1] 7.705276 9.117945
```

- Con una muestra bootstrap, se puede reproducir la forma y dispersión de la distribución. En el caso de la media es obvio, pero esto no es evidente para otras estadísticas. Ese es justamente la fortaleza del bootstrap, que nos permite obtener una estimación de la distribución de la estadística, no sólo para la media, sino para estadísticas mucho más complejas.

## Ejemplos adicionales. Media poblacional $V$

- La media de la distribución bootstrap está centrada cerca de la media de la muestra obtenida, no la media poblacional. Así que la media de muestra bootstrap tiene sesgo como estimador de la media poblacional. Sin embargo, el tamaño del sesgo de la estimación es parecido al sesgo que puede tener la media muestral respecto a la media poblacional.

```
B <- 1000
n <- length(ilec)
z <- numeric(B)
for(b in 1:B){
  i <- sample(1:n, replace=T)
  R[b] <- mean(ilec[i])
}
hist(R, prob=T)
```







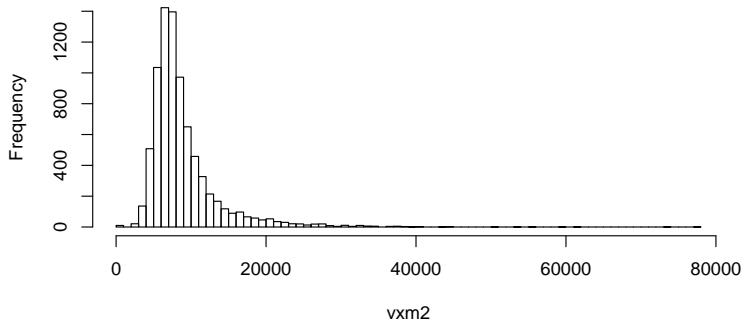
## Ejemplo 2: Intervalos de confianza I

- Consideren los costos por metro cuadrado de vivienda en México. Los datos corresponden al Promedio de valor del mercado por metro cuadrado de construcción, que se puede obtener de la Sociedad Hipotecaria Federal (avaluos) a partir de avalúos. Los datos corresponden a 8,076 registros de todos los municipios de México registrados en 2015, y los datos consideran todo tipo de vivienda, desde aquella de interés social como la residencial.
- La variable de interés es `vxm2`, el valor de la vivienda por metro cuadrado de construcción, en pesos. Los datos están en el archivo `SHF-PromedioValorMercadoxMetroCuadradoConstruccion.csv` y se encuentran en el blog.

```
vxm2 <-
read.csv("~/Dropbox/Public/SimS18-I/data/SHF-PromedioValorMercadoxMetroCuadradoConstruccion.csv") $vxm2
hist(vxm2, breaks=100)
```

## Ejemplo 2: Intervalos de confianza II

### Histogram of vxm2



summary (vxm2)

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0	6242	7616	8861	9848	77161

## Ejemplo 2: Intervalos de confianza III

- La distribución de los precios no es normal, y está 'contaminada' por diferentes tipos de propiedades. Para intentar corregir este efecto, consideremos un estimador robusto del parámetro de localización, que puede ser la media recortada del 25% que es la media del 50% de los datos que están en el centro de la distribución. Con bootstrap estimemos su distribución.

```
summary(vxm2)

  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
    0     6242     7616    8861    9848   77161

n <- length(vxm2)
mean(vxm2,trim = .25)

[1] 7747.663

z <- boot::boot(data = vxm2, statistic = function(x,i)mean(x[i],trim=.25), R = 1000)
par(mfrow=c(1,2))
hist(z$t)
qqnorm(z$t)
```



## Ejemplo 2: Intervalos de confianza V

```
z

ORDINARY NONPARAMETRIC BOOTSTRAP

Call:
boot::boot(data = vxm2, statistic = function(x, i) mean(x[i],
  trim = 0.25), R = 1000)

Bootstrap Statistics :
  original      bias    std. error
t1* 7747.663 -0.2707323    32.13932
```

- La distribución bootstrap es aproximadamente normal, tiene sesgo pequeño relativo al valor de la media, y tiene un error estándar de 32.1393157. Entonces podemos calcular un intervalo de confianza.

```
boot::boot.ci(z, type=c("norm", "basic", "perc"))

BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
Based on 1000 bootstrap replicates

CALL :
boot::boot.ci(boot.out = z, type = c("norm", "basic", "perc"))

Intervals :
Level      Normal              Basic              Percentile
95%   (7685, 7811 )   (7686, 7815 )   (7681, 7809 )
Calculations and Intervals on Original Scale
```











## Ejemplo Pruebas de permutación II

- Varianzas diferentes y desconocidas:  $\sigma_1 \neq \sigma_2$ . En este caso, la estadística de prueba es de la forma (aproximada):

$$t_0 = \frac{\bar{x} - \bar{y}}{\sqrt{S_1^2/n_1 + S_2^2/n_2}} \sim t_{(\nu)}$$

y los grados de libertad se obtienen usando la aproximación de Satterthwaite:

$$\nu = \frac{(S_1^2/n_1 + S_2^2/n_2)^2}{\frac{(S_1^2/n_1)^2}{n_1+1} + \frac{(S_2^2/n_2)^2}{n_2+1}} - 2$$

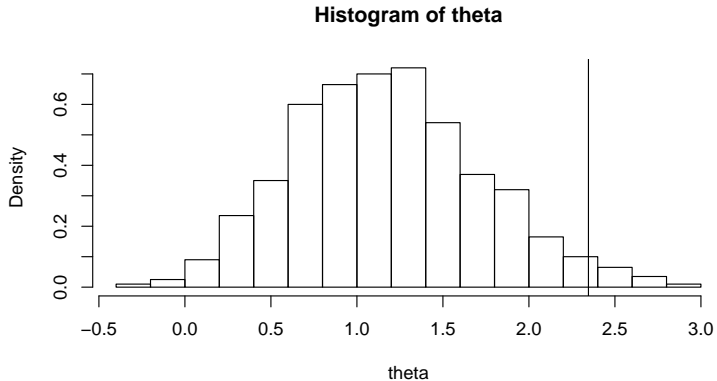
- En los dos casos anteriores, se tienen aproximaciones que pueden fallar si los datos no siguen una distribución normal.
- Las pruebas de permutación simplifican mucho el problema y se obtienen resultados más exactos.

## Ejemplo de pruebas de permutación I

- Bajo la hipótesis nula, esperaríamos que las diferencias de medias muestrales  $\hat{\theta} = \bar{x} - \bar{y}$  fueran cercanas a 0. Si  $H_0$  no es cierta, entonces tenderíamos a observar valores “grandes” de  $\hat{\theta}$ . Habiendo observado  $\hat{\theta}$ , definimos el  $p$ -value como  $pv = P(\hat{\theta}^* \geq \hat{\theta} | H_0)$  (también se denota como ASL, *achieved significance level*).
- Si  $H_0$  es cierta, cualquiera de las observaciones pudo haber venido de cualquiera de las poblaciones. Podemos entonces combinar *todas* las  $n_1 + n_2$  observaciones de las dos muestras<sup>3</sup>, y extraer *sin reemplazo*  $n_1$  de ellas para crear el primer grupo, y las  $n_2$  observaciones restantes forman el otro grupo. Con estas nuevas muestras calculamos la estadística de prueba (cualquiera que sea para probar esta hipótesis) y repetimos el proceso  $B$  veces.

## Simulación

## Ejemplo de pruebas de permutación III



<sup>3</sup>En este sentido son indistinguibles e intercambiables, como se comentó anteriormente

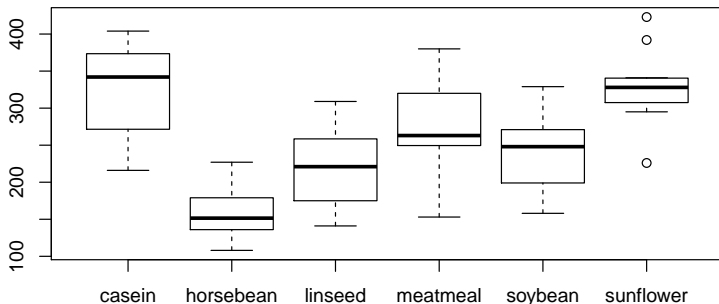


## Simulación





0 1



Los pesos ordenados para esas dos muestras son los siguientes:

```
x <- with(chickwts, sort(weight[feed=="soybean"]))
y <- with(chickwts, sort(weight[feed=="linseed"]))
x

[1] 158 171 193 199 230 243 248 248 250 267 271 316 327 329

y

[1] 141 148 169 181 203 213 229 244 257 260 271 309
```

Ambos grupos se pueden comparar utilizando alguna estadística de prueba, basada en la media, mediana, media recortada, o cualquier otra medida de tendencia central. Por ejemplo, considerando las medias, y usando la estadística  $t$  bajo el supuesto de normalidad

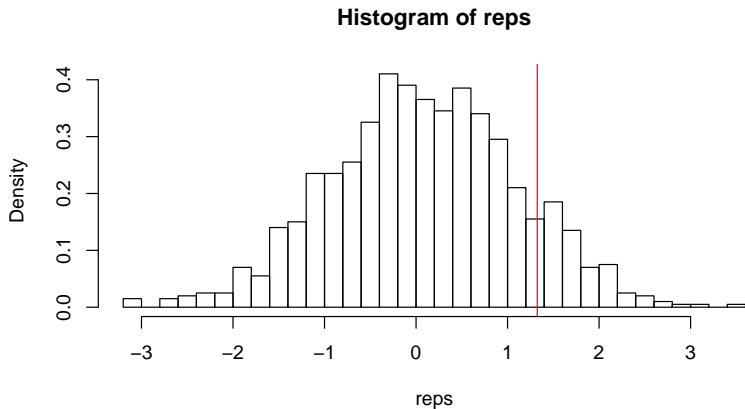
```
t.test(x=x,y=y)
```

Welch Two Sample t-test

```
data: x and y
t = 1.3246, df = 23.63, p-value = 0.198
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -15.48547  70.84262
sample estimates:
mean of x mean of y
246.4286  218.7500
```



## Ejemplos V



Para una prueba de dos colas, el  $ASL$  es  $2\hat{p}$  si  $\hat{p} \leq 0.5$  (y es  $2(1 - \hat{p})$  cuando  $\hat{p} > 0.5$ ). En este ejemplo,  $ASL = 0.238$ . Por lo tanto, no se tiene evidencia para rechazar la hipótesis nula de igualdad de medias.











```
par(mfrow=c(2,2))
par(oma=c(1,1,1,1),mar=c(4,4,1,1))
a <- seq(10,40,0.1) #sucesión para graficar los ajustes

L1 <- lm(y ~ x)
plot(x,y,main="Lineal",pch=16)
yhat1 <- L1$coef[1] + L1$coef[2]*a
lines(a,yhat1,lwd=2)

L2 <- lm(y ~ x + I(x^2))
plot(x,y,main="Cuadrático",pch=16)
yhat2 <- L2$coef[1] + L2$coef[2]*a + L2$coef[3]*a^2
lines(a,yhat2,lwd=2)

L3 <- lm(log(y) ~ x)
plot(x,y,main="Exponencial",pch=16)
logyhat3 <- L3$coef[1] + L3$coef[2]*a
yhat3 <- exp(logyhat3)
lines(a,yhat3,lwd=2)

L4 <- lm(log(y) ~ log(x))
plot(log(x),log(y),main="Log-Log",pch=16)
logyhat4 <- L4$coef[1] + L4$coef[2]*log(a)
lines(log(a),logyhat4,lwd=2)
```



- Una vez que el modelo es ajustado, se evalúa el ajuste.

## Procedimiento para estimar el error de predicción usando validación cruzada (uno afuera)

- 1 Para  $k = 1, \dots, n$  dejar la observación  $(x_k, y_k)$  para ser el punto de prueba y usar las observaciones restantes para ajustar el modelo.
  - a. Ajusta el modelo usando sólo  $n - 1$  observaciones en el conjunto de entrenamiento.
  - b. Calcular la respuesta predictiva  $\hat{y}_k = \hat{\beta}_0 + \hat{\beta}_1 x_k$  para el punto de prueba
  - c. Calcula el error de predicción  $e_k = y_k - \hat{y}_k$ .
- 2 Estima la media de los errores de predicción al cuadrado  $\hat{\sigma}_\epsilon^2 = \frac{1}{n} \sum_{k=1}^n e_k^2$ .

9 1