

Introducció Bootstrap. Sessió 2, 3 (NO Paramètric versus Paramètric)

Jordi Ocaña, Sergi Civit,

20 d'abril de 2018

1 Introducció: Situació de laboratori

Aquesta mostra fa el paper d'unes "dades reals" però en realitat sabem que procedeix d'una $N(15, 3)$

ATENCIÓ: AIXÒ ÉS AIXÍ PER QUE ESTEM EN UNA SITUACIÓ "DE LABORATORI", AMB DADES REALS LÒGICAMENT DESCONEIXERÍEM COMPLETAMENT ELS PARÀMETRES REALS COM A MOLT PODRÍEM TENIR ALGUNA IDEA DE LA SEVA FORMA per exemple que SEMBLA O S'ARPOXIMA A UNA NORMAL

2 Mostra

```
> # Lectura de les dades PROCEDENTS D'UNA NORMAL
> x <- c(15.54, 21.06, 16.52, 13.62, 16.14, 10.98, 13.53, 16.02, 16.79, 15.90)
> n<-length(x)
> mu <- 15
> sigma <- 3
```

3 Estudi de la distribució de l'estadístic t

A partir de la mostra anterior anem a dur a terme l'estudi bootstrap de la distribució QUE ENS PROPORCIONA LA METODOLOGIA PARAMÈTRICA
Estadístic t

```
> tStud <- function(x, mitjana)
+ {
+   (mean(x) - mitjana) / sqrt(var(x)/length(x))
+ }
```

3.1 Estudi de la distribució de l'estadístic t: Càlculs i gràfic

Obtenim:

- Mitjana mostral,
- Desviació típica mostral,
- Error estàndard de la mitjana mostral,
- Estimació de l'error estàndard de la mitjana mostral

```
> # Mitjana mostral, estimació de la veritable mitjana (que recordem que és 15)
> xBarra <- mean(x)
> xBarra

[1] 15.61

> # Desviació típica mostral, estimació de la veritable sigma (que és 3)
> s.x <- sqrt(var(x))
> s.x

[1] 2.62857

> # Veritable valor de l'error estàndard de la mitjana mostral:
> sigma.xBarra <- sigma/sqrt(n)
> sigma.xBarra

[1] 0.9486833

> # Estimació de l'error estàndard de la mitjana mostral:
> s.xBarra <- s.x/sqrt(n)
> s.xBarra

[1] 0.8312267

> # Calculem l'estadístic t i el grafiquem:
> t.x <- tStud(x, mu)
> t.x

[1] 0.7338552

> rang.t <- seq(from=-4, to=+4, by=0.1)
> dens.veritat <- dt(rang.t, df=n-1)
> windows(21,21)
> plot(rang.t, dens.veritat, type="l", col="green", ylim=c(0,0.4))
```

3.2 Estudi de la distribució de l'estadístic t: Valors crítics i Interval de confiança

Molts cops, es pot interessar coneixer certs "valors crítics" d'aquesta distribució mostral.

Aui tenim els **valors crítics** segons la "veritable" distribució mostral de t

```
> tCritStud <- qt(c(0.975, 0.025), df = n - 1)
> tCritStud
```

```
[1] 2.262157 -2.262157
```

Amb aquests valors podriem calcular un **interval de confiança paramètric** a partir d'una $t(n - 1)$ gl

```
> xBarra - tCritStud * s.xBarra
```

```
[1] 13.72963 17.49037
```

Suposem ara la **Normal com a model aproximat** i el grafiquem el tabulem i calculem el interval de confiança

```
> dens.normAprox <- dnorm(rang.t)
> lines(rang.t, dens.normAprox, type="l", col="blue")
```

Obtenim **valors crítics** segons aproximació normal:

```
> tCritNorm = qnorm(c(0.975, 0.025))
> tCritNorm
```

```
[1] 1.959964 -1.959964
```

Realitzem **Interval de Confiança** segons aproximació normal:

```
> xBarra - tCritNorm * s.xBarra
```

```
[1] 13.98083 17.23917
```

4 Bootstrap NO PARAMÈTRIC de la distribució de l'estadístic t

```
> B <- 10000
> sample(x, replace = TRUE)

[1] 10.98 13.62 16.02 16.14 16.14 15.54 16.02 15.90 16.14 16.79

> set.seed(127)
> mostres.bootstrap <- matrix(sample(x, replace=T, size=B*n), ncol=B)
> mostres.bootstrap[,1:10]
```

```

      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
[1,] 16.52 16.79 16.14 16.79 16.52 21.06 15.90 16.14 16.52 10.98
[2,] 15.54 13.62 16.14 15.90 15.54 13.62 13.53 21.06 16.52 16.02
[3,] 16.52 16.02 16.02 16.14 15.90 15.54 16.79 13.62 10.98 10.98
[4,] 21.06 21.06 15.90 16.52 21.06 16.02 16.14 15.54 13.62 15.90
[5,] 13.62 16.14 10.98 16.52 16.79 21.06 16.52 13.62 16.79 16.52
[6,] 16.02 13.62 13.62 15.90 16.02 13.62 13.62 16.52 16.79 15.54
[7,] 10.98 16.02 13.62 21.06 13.53 13.53 15.54 15.90 13.62 15.54
[8,] 16.52 21.06 16.14 15.90 10.98 16.02 21.06 16.14 16.79 21.06
[9,] 16.02 10.98 16.14 13.62 15.54 15.90 16.14 21.06 13.62 16.52
[10,] 16.79 15.54 13.62 16.02 13.53 15.90 13.62 16.52 13.53 10.98

> t.bootstrap <- apply(mostres.bootstrap, 2, tStud, mitjana=xBarra)
> t.bootstrap[1:10]

[1] 0.43437145 0.47895372 -1.37915517 1.41891456 -0.08298137 0.70798559
[7] 0.39447074 1.23138657 -1.12476800 -0.59853481

> dens.bootstrap <- density(t.bootstrap, from=-4, to=+4)
> lines(dens.bootstrap, type="l", col="red")

```

4.1 Bootstrap NO PARAMÈTRIC de la distribució de l'estadístic t: Valors crítics i Interval de Confiança

Obtenim valors crítics

```
> tCritBoot = quantile(t.bootstrap, probs = c(0.975, 0.025))
```

Interval de confiança bootstrap no paramètric

```
> xBarra - tCritBoot * s.xBarra
```

```

      97.5%      2.5%
13.84873 17.46900

```

5 Bootstrap PARAMÈTRIC: Gràfic, Valors crítics, Interval de Confiança

Suposem que: EN REALITAT CREIEM de que la forma de la distribució es NORMAL

```

> # Si en realidad nos podemos fiar de que la forma de la distribución es normal
> # Bootstrap paramètric normal:
> # Una sola remostra:
> rnorm(n, mean=xBarra, sd=s.x)

[1] 15.08031 16.34837 15.04077 14.52995 19.77730 17.58484 16.35833 13.93648
[9] 15.36730 17.20246

```

```

> B <- 10000
> set.seed(127)
> mostres.bootstrap <- matrix( rnorm(B*n, mean=xBarra, sd=s.x), ncol=B)
> t.bootstrap.param <- apply(mostres.bootstrap ,2, tStud, mitjana=xBarra)
> dens.bootstrap.param <- density(t.bootstrap.param,from=-4, to=+4)
> lines(dens.bootstrap.param, type="l", col="brown")

```

Obtenim **valors crítics**

```

> # valors crítics segons bootstrap paramètric
> tCritBoot.param = quantile(t.bootstrap.param, probs = c(0.975, 0.025))

```

Interval de confiança bootstrap paramètric

```

> xBarra - tCritBoot.param * s.xBarra

```

```

      97.5%      2.5%
13.72557 17.44400

```

6 Càlcul d'una probabilitat $P[t > \text{valor}]$

Suposem que per decidir quin és el valor que a la seva dreta i deixa 0.025 de probabilitat, hem decidit utilitzar les taules de la normal (és a dir, hem considerat que **t té distribució $N(0,1)$**)

```

> z0.05 <- qnorm(0.025, lower.tail = FALSE) # 1.959964

```

En teoria tindriem:

```

> pnorm(z0.05, lower.tail = FALSE)

```

```

[1] 0.025

```

6.1 Càlcul d'una probabilitat $P[t > \text{valor}]$ t-student n-1

```

> # Però la veritable probabilitat és:
> pt(z0.05, df = n - 1, lower.tail = FALSE)

```

```

[1] 0.04082457

```

6.2 Càlcul d'una probabilitat $P[t > \text{valor}]$ Bootstrap No paramètric

Si utilitzem l'aproximació a la veritable distribució de t que ens ha proporcionat el bootstrap no paramètric, aquesta probabilitat és:

```

> length(t.bootstrap[t.bootstrap > z0.05])/B

```

```

[1] 0.0321

```

6.3 Càlcul d'una probabilitat $P[t > \text{valor}]$ Bootstrap Paramètric

Si utilitzem l'aproximació a la veritable distribució de t que ens ha proporcionat el bootstrap paramètric, aquesta probabilitat és:

```
> length(t.bootstrap.param[t.bootstrap.param > z0.05])/B  
[1] 0.0406
```

7 Sessió 3

8 Mostra Exponencial. Estadístic t

```
> sigma <- mu  
> x <- c(8.51, 8.71, 69.19, 10.05, 23.64, 8.67, 1.51, 20.36, 1.23, 5.27)  
> n = length(x)  
> sigma.xBarra <- sigma/sqrt(n)  
> xBarra <- mean(x)  
> s.x <- sqrt(var(x))  
> s.xBarra <- s.x/sqrt(n)  
> t.x <- tStud(x, mu)  
> rang.t <- seq(from=-4, to=+4, by=0.1)
```

Ara és difícil determinar **analíticament quina és la veritable distribució mostral de l'estadístic t** quan les dades són exponencials.

8.1 Aproximació mitjançant SIMULACIÓ a la VERITABLE distribució mostral

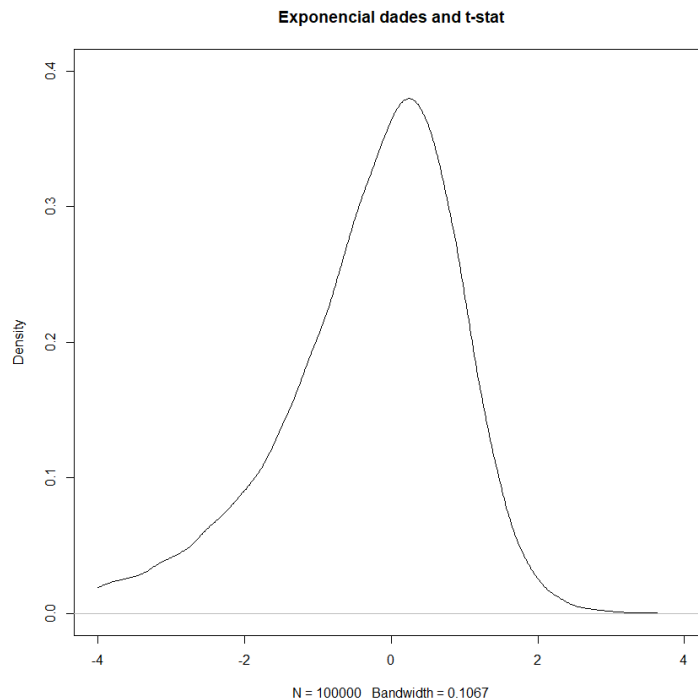
Fem una simulació utilitzant els VERITABLES VALORS DELS PARÀMETRES POBLACIONALS

PREGUNTA en què es diferencia del que fem amb bootstrap?)

```
> m <- 100000  
> mostra.t <- apply(matrix(rexp(m*n, rate=1/mu), ncol=m), 2, tStud, mitjana=mu)  
> dens.veritat <- density(mostra.t, from=-4, to=+4)  
> # Obrim una altra finestra gràfica independent:  
> windows(21,21)  
> plot(dens.veritat, type="l", col="black", ylim=c(0,0.4), main="Exponencial dades, t-stat")
```

¿Es pivotal? . Provarem amb diferents valors de " μ ". ¿És manté estable la distribució muestral?. Grafiquem (... semblà que sí)

```
> color = 0  
> for (otherMu in c(1,2,5,10,20,50)) {
```

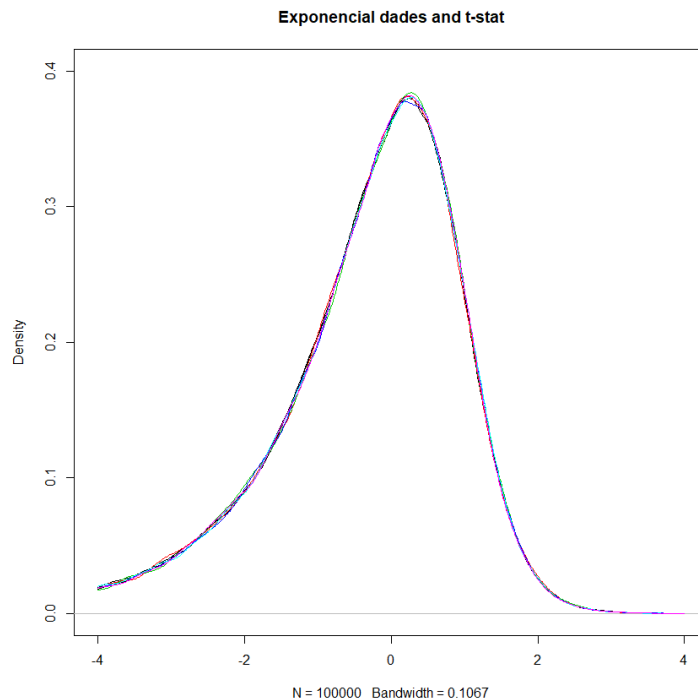


```
+ color = color + 1
+ lines(
+   density(
+     apply(matrix(rexp(m*n, rate=1/otherMu), ncol=m), 2, tStud, mitjana=otherMu),
+     from=-4, to=+4
+   ),
+   type="l", col=color, ylim=c(0,0.4)
+ )
+ }
```

8.2 Aproximació mitjançant DISTTRIBUCIÓ NORMAL I t-STUDENT a la veritable distribució mostral

```
> par(mfrow=c(2,2))
> plot(dens.veritat, type="l", col="black", ylim=c(0,0.4),main="dens.veritat versus")
> dens.normAprox <- dnorm(rang.t)
> lines(rang.t, dens.normAprox, type="l", col="blue")
> legend("topright", legend=c("Normal"))
```

L'aproximació normal sembla que no seria gaire adequada!!!



```
> plot(dens.veritat, type="l", col="black", ylim=c(0,0.4),main="dens.veritat versus")
> lines(rang.t, dt(rang.t, df = n - 1), type="l", col="green", ylim=c(0,0.4))
> legend("topright", legend=c("t-Stud"))
```

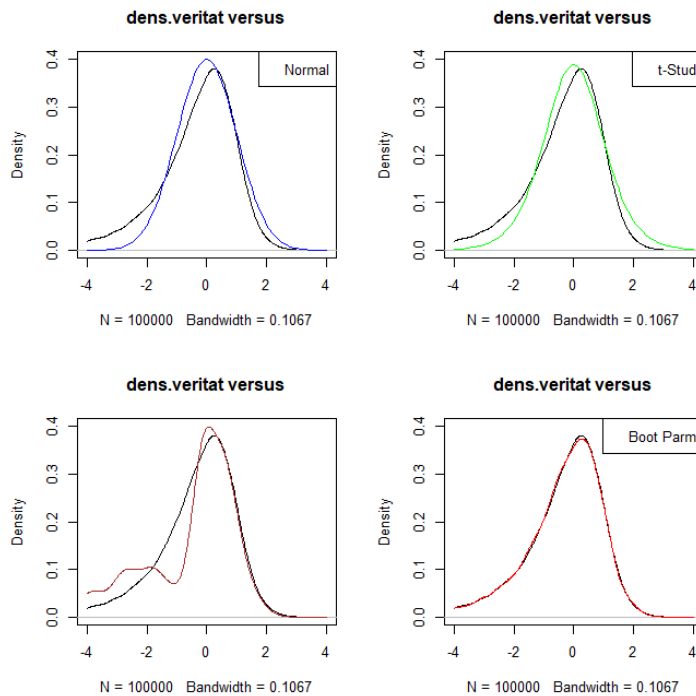
Ni tampoc la t de Student!!!

8.3 Aproximació mitjançant BOOTSTRAP NO PARA-MÈTRIC a la veritable distribució mostral

```
> # Bootstrap no paramètric:
> B <- 10000
> mostres.bootstrap <- matrix(sample(x, replace=T, size=B*n), ncol=B)
> t.bootstrap <- apply(mostres.bootstrap ,2, tStud, mitjana=xBarra)
> dens.bootstrap <- density(t.bootstrap,from=-4, to=+4)
> plot(dens.veritat, type="l", col="black", ylim=c(0,0.4),main="dens.veritat versus")
> lines(dens.bootstrap, type="l", col="brown", ylim=c(0,0.4))
>
```


8.4 Aproximació mitjançant BOOTSTRAP PARAMÈTRIC EXPONENCIAL a la veritable distribució mostral

```
> # Bootstrap paramètric exponencial:
> mostres.bootstrap <- matrix(rexp(n=B*n, rate=1/xBarra), ncol=B)
> t.bootstrap.param <- apply(mostres.bootstrap, 2, tStud, mitjana=xBarra)
> dens.bootstrap.param <- density(t.bootstrap.param, from=-4, to=+4)
> plot(dens.veritat, type="l", col="black", ylim=c(0,0.4), main="dens.veritat versus")
> lines(dens.bootstrap.param, type="l", col="red", ylim=c(0,0.4))
> legend("topright", legend=c("Boot Parm"))
```



8.5 Càlcul d'una probabilitat a la cua dreta: $P[t > 1.96]$

```
> # Càlcul d'una probabilitat a la cua dreta:  $P[t > 1.96]$ 
> z0.05 <- qnorm(0.025, lower.tail = FALSE) # 1.959964
> 1 - pnorm(z0.05) # ... aproximació normal

[1] 0.025

> 1 - pt(z0.05, df=n-1) # ... aproximació t
```

```

[1] 0.04082457
> length(mostra.t[mostra.t > z0.05])/m # ... veritable probabilitat (simulació 100000)
[1] 0.00923
> length(t.bootstrap[t.bootstrap > z0.05])/B # ... bootstrap no paramètric
[1] 0.0053
> length(t.bootstrap.param[t.bootstrap.param > z0.05])/B # ... bootstrap paramètric
[1] 0.0093

```

8.6 Càlcul d'una probabilitat a la cua esquerra: $P[t < -1.96]$

```

> pnorm(-z0.05) # ... aproximació normal
[1] 0.025
> pt(-z0.05, df=n-1) # ... aproximació t
[1] 0.04082457
> length(mostra.t[mostra.t < -z0.05])/m # ... "veritable" probabilitat
[1] 0.12156
> length(t.bootstrap[t.bootstrap < -z0.05])/B # ... bootstrap no paramètric
[1] 0.2431
> length(t.bootstrap.param[t.bootstrap.param < -z0.05])/B # ... bootstrap paramètric
[1] 0.1215

```

8.7 Càlcul d'una probabilitat bilateral: $P[|t| > 1.96]$

```

> # Càlcul d'una probabilitat bilateral:  $P[|t| > 1.96]$ 
> pnorm(-z0.05) + (1 - pnorm(z0.05))
[1] 0.05
> pt(-z0.05, df=n-1) + (1 - pt(z0.05, df=n-1))
[1] 0.08164913
> length(mostra.t[abs(mostra.t) > z0.05])/m
[1] 0.13079
> length(t.bootstrap[abs(t.bootstrap) > z0.05])/B
[1] 0.2484
> length(t.bootstrap.param[abs(t.bootstrap.param) > z0.05])/B
[1] 0.1308

```

9 Exercici: Mostra exponencial, estadístic t, n=40.

Càlcul de probabilitats: Cua drete, cua esquerra, 2 cues:

```
> # Mostra exponencial, estadístic t, n=40:
> set.seed(127)
> n <- 40
> sigma <- mu
> x <- rexp(n, rate=1/mu) #SIMULEM DADES MODEL EXPONENCIAL
> sigma.xBarra <- sigma/sqrt(n)
> xBarra <- mean(x)
> s.x <- sqrt(var(x))
> s.xBarra <- s.x/sqrt(n)
> t.x <- tStud(x, mu)
> rang.t <- seq(from=-4, to=+4, by=0.1)
> # Aproximació mitjançant simulació a la veritable distribució mostral:
> m <- 10000
> mostra.t <- apply(matrix(rexp(m*n, rate=1/mu), ncol=m) ,2, tStud, mitjana=mu)
> par(mfrow=c(2,2))
> dens.veritat <- density(mostra.t,from=-4, to=+4)
> plot(dens.veritat, type="l", col="green", ylim=c(0,0.4), main="Simulem dades Exponencials")
```

9.1 Exercici: Mostra exponencial, estadístic t, n=40. Aprox Normal.

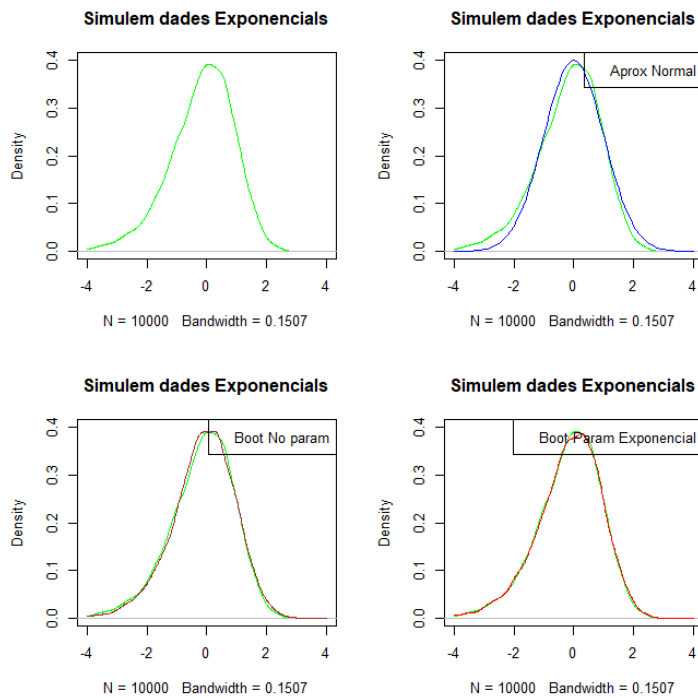
```
> # Aproximació mitjançant simulació a la veritable distribució mostral:
> dens.normAprox <- dnorm(rang.t)
> plot(dens.veritat, type="l", col="green", ylim=c(0,0.4), main="Simulem dades Exponencials")
> lines(rang.t, dens.normAprox, type="l", col="blue")
> legend("topright", legend=c("Aprox Normal"))
```

9.2 Exercici: Mostra exponencial, estadístic t, n=40. Boot No paramètric.

```
> # Bootstrap no paramètric:
> B <- 10000
> mostres.bootstrap <- matrix(sample(x, replace=T, size=B*n), ncol=B)
> t.bootstrap <- apply(mostres.bootstrap ,2, tStud, mitjana=xBarra)
> dens.bootstrap <- density(t.bootstrap,from=-4, to=+4)
> plot(dens.veritat, type="l", col="green", ylim=c(0,0.4), main="Simulem dades Exponencials")
> lines(dens.bootstrap, type="l", col="brown")
> legend("topright", legend=c("Boot No param"))
```

9.3 Exercici: Mostra exponencial, estadístic t, n=40. Bot Paramètric Exponencial.

```
> # Bootstrap paramètric exponencial:
> mostres.bootstrap <- matrix(rexp(n=B*n, rate=1/xBarra), ncol=B)
> t.bootstrap.param <- apply(mostres.bootstrap ,2, tStud, mitjana=xBarra)
> dens.bootstrap.param <- density(t.bootstrap.param,from=-4, to=+4)
> plot(dens.veritat, type="l", col="green", ylim=c(0,0.4), main="Simulem dades Exponencials")
> lines(dens.bootstrap.param, type="l", col="red")
> legend("topright", legend=c("Boot Param Exponencial"))
```



9.4 Càlcul de probabilitats. $P[t > 1.96]$

```
> # Càlcul d'una probabilitat a la cua dreta:  $P[t > 1.96]$ 
> z0.05 <- 1.959964
> 1 - pnorm(z0.05) # ... aproximació normal

[1] 0.025

> 1 - pt(z0.05, df=n-1) # ... aproximació t
```

```

[1] 0.02858665
> length(mostra.t[mostra.t > z0.05])/m # ... veritable probabilitat
[1] 0.0086
> length(t.bootstrap[t.bootstrap > z0.05])/B # ... bootstrap no paramètric
[1] 0.0134
> length(t.bootstrap.param[t.bootstrap.param > z0.05])/B # ... bootstrap paramètric
[1] 0.0102

```

9.5 Càlcul de probabilitats. $P[t < -1.96]$

```

> # Càlcul d'una probabilitat a la cua esquerra:  $P[t < -1.96]$ 
> pnorm(-z0.05) # ... aproximació normal
[1] 0.025
> pt(-z0.05, df=n-1) # ... aproximació t
[1] 0.02858665
> length(mostra.t[mostra.t < -z0.05])/m # ... "veritable" probabilitat
[1] 0.0625
> length(t.bootstrap[t.bootstrap < -z0.05])/B # ... bootstrap no paramètric
[1] 0.0547
> length(t.bootstrap.param[t.bootstrap.param < -z0.05])/B # ... bootstrap paramètric
[1] 0.0647

```

9.6 Càlcul de probabilitats. $P[|t| > 1.96]$

```

> # Càlcul d'una probabilitat bilateral:  $P[|t| > 1.96]$ 
> pnorm(-z0.05) + (1 - pnorm(z0.05))
[1] 0.05
> pt(-z0.05, df=n-1) + (1 - pt(z0.05, df=n-1))
[1] 0.0571733
> length(mostra.t[abs(mostra.t) > z0.05])/m
[1] 0.0711

```

```
> length(t.bootstrap[abs(t.bootstrap) > z0.05])/B
[1] 0.0681
> length(t.bootstrap.param[abs(t.bootstrap.param) > z0.05])/B
[1] 0.0749
>
>
```