

Tema 3 Estimació puntual

L'objectiu principal de la **inferència estadística** és obtenir informació sobre tota una població a partir de només una mostra, com quan volem saber si un brou és fat o salat tastant-ne només una cullerada. El primer tipus d'informació que ens sol interessar és què val qualche paràmetre d'alguna variable aleatòria poblacional (una proporció, una mitjana...), per exemple per poder escriure un titular com el següent:



Figura 3.1: <https://www.efesalud.com/miopía-estudio-universitarios>

Aquest 60% no s'ha obtingut fent passar a tots els universitaris espanyols un test de miopia, ni tan sols demanant-los a tots si són miops o no, sinó que simplement s'ha pres una mostra d'universitaris, s'hi ha observat un 60% de miops i s'ha extrapolat aquesta proporció a tot el col·lectiu d'universitaris espanyols.

El procés d'intentar endevinar el valor d'un paràmetre d'una població a partir d'una mostra se'n diu **estimació puntual**, i és el que tractarem en aquest tema.

En aquest curs, sempre suposarem que empram mostres aleatòries i gairebé sempre que aquestes mostres aleatòries són a més simples. Per tant, si no diem el contrari, d'ara endavant **quan parlem de mostres sempre suposarem que són mostres aleatòries simples**, encara que no ho diguem explícitament

per no carregar massa el text.

3.1 Definicions bàsiques

Per estimar el valor d'un paràmetre d'una variable aleatòria poblacional, en prenem una mostra (aleatòria simple) i calculam qualche cosa amb els valors que la formen. Què calculam? Doncs un **estimador**: alguna funció adequada aplicada als valors de la mostra, i que dependrà del que volguem estimar.

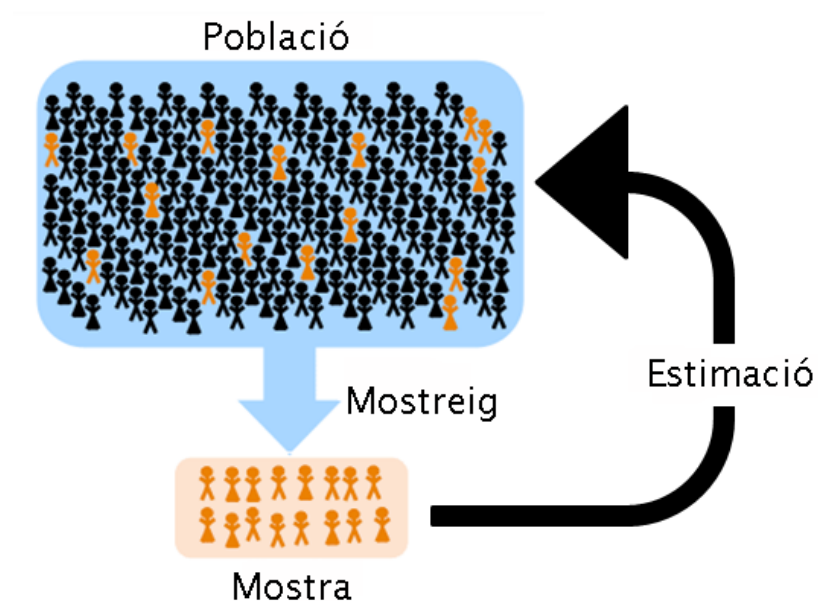


Figura 3.2: Població *versus* mostra

Per exemple:

- Si volem estimar l'alçada mitjana dels estudiants de la UIB, prendrem una mostra d'estudiants de la UIB, els amidarem i calcularem la **mitjana aritmètica** de les seves alçades.
- Si volem estimar la proporció d'estudiants de la UIB que han passat la COVID-19, prendrem una mostra d'estudiants de la UIB, els farem un test d'anticossos i calcularem la **proporció mostral** de positius en la mostra.

Formalment:

- Tenim una **variable aleatòria poblacional** X , definida sobre una **població**.
- Una **mostra aleatòria simple** de mida n de X és un vector (X_1, \dots, X_n) format per n còpies *independents* de X .

Cada variable X_i és una còpia de “Prenem un subjecte de la població i hi mesuram X ”.

- Una **realització** de la mostra aleatòria simple (X_1, \dots, X_n) és un vector $(x_1, \dots, x_n) \in \mathbb{R}^n$ de valors presos per aquestes variables aleatòries.

És a dir, amb (X_1, \dots, X_n) repetim n vegades (independents les unes de les altres) el procés de prendre un subjecte de la població i mesurar-hi X . Cada vegada que ho fem, obtenim un vector de números, al que diem una **realització** de la mostra.

A la lliçó anterior a aquestes realitzacions les déiem directament “mostres aleatòries simples de valors de X ”; no passeu ànsia, en sortir d’aquest “formalment” els ho tornarem a dir.

- Un **estimador** és una variable aleatòria $f(X_1, \dots, X_n)$ obtinguda aplicant una funció f a una mostra aleatòria simple (X_1, \dots, X_n) .

Aquest estimador s’aplica a les realitzacions de la mostra i dóna nombres reals.

Un estimador és una **variable aleatòria**, definida sobre la població formada per les mostres aleatòries simples de la població de partida. Per tant, té funció de densitat, funció de distribució (que genèricament anomenarem **distribució mostral**, per indicar que refereix a la probabilitat que li passi alguna cosa al valor del estimador sobre una mostra), esperança, desviació típica, etc.

Com ja us hem dit, i com que no hi ha necessitat de filar tan prim, d’ara endavant cometrem l’abús de llenguatge de dir **mostra (aleatòria simple)** tant al vector de variables aleatòries (X_1, \dots, X_n) com a una realització (x_1, \dots, x_n) i hi ometrem els parèntesis.

Com ja hem comentat a la Secció 1.5, si la mida N de la població és MOLT més gran que la mida n de la mostra (per fixar idees, hem dit que si $N \geq 1000n$), els resultats per a mostres aleatòries simples valen (aproximadament) per a mostres aleatòries sense reposició, perquè les variables aleatòries que formen la mostra sense reposició són gairebé idèntiques i independents i les repeticions són improbables.

Els estimadors tenen sempre sentit per a mostres en general, però gairebé tots els teoremes que estableixen les seves propietats són vertaders només sota determinades restriccions (mostra aleatòria simple, condicions extra sobre X , ...), per la qual cosa les seves conseqüències tan sols són segures sota aquestes restriccions.

3.2 Mitjana mostral

Quan volem estimar el valor mitjà d'una variable sobre una població, en prenem una mostra de valors i calculam la seva mitjana aritmètica, no és ver? Doncs això és la mitjana mostral.

Donada una variable aleatòria X , la seva **mitjana mostral** (de mostres aleatòries simples) **de mida** n és la variable aleatòria \bar{X} “Prenem una mostra aleatòria simple de mida n de X i calculam la mitjana aritmètica dels seus valors”. És a dir, formalment, la **mitjana mostral** de mida n de X és la variable aleatòria obtinguda prenent n còpies independents X_1, \dots, X_n de la variable aleatòria X i calculant

$$\bar{X} = \frac{X_1 + \dots + X_n}{n}$$

Fixau-vos que definim la mitjana mostral només per a mostres aleatòries simples. Naturalment, té sentit definir-la per a mostres qualssevol, però llavors la seva distribució mostral deixa de complir les propietats que donarem en aquesta secció. El mateix advertiment val per als estimadors que definim en les pròximes seccions.

Com a conseqüència del comportament d'esperances i variàncies de combinacions lineals, tenim el resultat següent:

Teorema 3.1 *Siguin X una variable aleatòria d'esperança μ_X i desviació típica σ_X , i \bar{X} la seva mitjana mostral (de mostres aleatòries simples) de mida n . Aleshores*

- a. *El valor esperat de \bar{X} és $E(\bar{X}) = \mu_X$.*
- b. *La desviació típica de \bar{X} és $\sigma(\bar{X}) = \sigma_X/\sqrt{n}$.*

En efecte, com que

$$\bar{X} = \frac{1}{n}X_1 + \cdots + \frac{1}{n}X_n$$

i les variables X_1, \dots, X_n són còpies de X , i per tant tenen totes esperança μ_X i variància σ_X^2 , tenim que

$$\mu_{\bar{X}} = \overbrace{\frac{1}{n}\mu_X + \cdots + \frac{1}{n}\mu_X}^n = \mu_X$$

i, si X_1, \dots, X_n són independents,

$$\sigma_{\bar{X}} = \sqrt{\overbrace{\frac{1}{n^2}\sigma_X^2 + \cdots + \frac{1}{n^2}\sigma_X^2}^n} = \sqrt{\frac{n}{n^2}\sigma_X^2} = \frac{\sigma_X}{\sqrt{n}}$$

Si us hi fixau, per demostrar que $E(\bar{X}) = \mu_X$ no hem fet servir que X_1, \dots, X_n siguin independents. De fet, **la igualtat $E(\bar{X}) = \mu_X$ també és vertadera per a mitjanes mostrals de mostres aleatòries sense reposició.** En canvi, la igualtat $\sigma_{\bar{X}} = \sigma_X/\sqrt{n}$ sí que requereix que les mostres aleatòries siguin simples.

Per tant:

- \bar{X} és un estimador puntual de μ_X .

- $E(\bar{X}) = \mu_X$, la qual cosa significa que:
 - La mitjana de les mitjanes mostrals de totes les mostres aleatòries de mida n de X torna a ser la mitjana μ_X de X .
 - *Esperam que la mitjana mostral doni μ_X* : si repetíssim moltes vegades el procés de prendre una mostra aleatòria de mida n i calcular-ne la mitjana mostral, molt probablement el valor mitjà d'aquestes mitjanes s'acostaria molt a μ_X .
- $\sigma(\bar{X}) = \sigma_X / \sqrt{n}$ indica que la dispersió de les mitjanes mostrals creix amb la dispersió de X i decreix amb la mida n de la mostra, tendint a 0 quan $n \rightarrow \infty$.

L'efecte de la mida de les mostres sobre la variabilitat de \bar{X} és raonable. Quan prenem mostres aleatòries d'una variable i en calculam la mitjana, el més normal és que dins cada mostra els valors més petits se compensin amb els més grans, i que com a conseqüència les mitjanes siguin més homogènies que els valors de la variable.

Vaja: si triam una persona a l'atzar, no és molt improbable que faci, jo què sé, 2.10 m. Però si prenem una mostra aleatòria de 50 persones, és molt més difícil que la mitjana de les seves alçades sigui 2.10 m. El que hi esperaríem és que les alçades dels més alts s'hi compensin amb les alçades dels més baixos i tot plegat doni una mitjana més "típica".

A la desviació típica de \bar{X} li diem l'**error estàndard**, o **típic**, de \bar{X} . Per tant, l'error estàndard de la mitjana mostral de mida n de X és σ_X / \sqrt{n} .

Exemple 3.1 El fitxer **tests.txt** que trobareu a l'url

<https://raw.githubusercontent.com/AprendeR-UIB/MatesII/master/Dades/tests.txt> conté les notes (sobre 100) de tests dels estudiants de Matemàtiques I de fa uns cursos. El guardam en un vector anomenat `tests` :

```
tests=scan("https://raw.githubusercontent.com/AprendeR-UIB/MatesII/master/Da
```

Considerarem la població dels estudiants de Matemàtiques I d'aquell curs i com a variable aleatòria d'interès X la seva nota de tests sobre 100. Per tant, aquest vector `tests` conté els valors de la variable aleatòria d'interès sobre tots els individus de la població. La seva mida és

```
N=length(tests)
```

```
N
```

```
## [1] 185
```

La seva mitjana, que és la mitjana poblacional μ_X , és

```
mu=mean(tests)
```

```
mu
```

```
## [1] 55.43243
```

Si en prenem una mostra aleatòria simple, per exemple de mida $n = 40$, la seva mitjana mostral no té per què coincidir amb la mitjana poblacional:

```
n=40
```

```
MAS=sample(tests,n,replace=TRUE) # Una mostra aleatòria simple
```

```
x.barra=mean(MAS) # La seva mitjana mostral
```

```
x.barra
```

```
## [1] 53.5
```

Però si prenem *moltes* mostres aleatòries simples, la mitjana de les seves mitjanes és molt probable que sí que s'acosti a la mitjana poblacional. Vegem si tenim sort amb cent mil mostres:

```
mitjanes=replicate(10^5,mean(sample(tests,n,replace=TRUE)))  
mean(mitjanes)
```

```
## [1] 55.4187
```

Vegem ara que la desviació típica d'aquesta mostra de mitjanes s'acosta a l'error típic de la mitjana mostral, no a la desviació típica de la població:

- La desviació típica poblacional:

```
sigma=sd(tests)*sqrt((N-1)/N)  
sigma
```

```
## [1] 21.38241
```

- La desviació típica de la mostra de mitjanes:

```
sd(mitjanes)
```

```
## [1] 3.384683
```

- L'error típic de la mitjana mostral:

```
sigma/sqrt(n)
```

```
## [1] 3.380856
```


Veiem que les mitjanes mostrals presenten una dispersió molt més petita que la variable poblacional original. Gràficament, als histogrames de les Figures 3.3 i 3.4 podeu veure com les mitjanes estan més concentrades al voltant de 55 que les notes originals.

Recordau del Teorema 2.6 que una combinació lineal de variables aleatòries normals independents torna a ser normal. Com que la mitjana mostral de mostres aleatòries simples és una combinació lineal de variables aleatòries independents, obtenim el resultat següent:

Teorema 3.2 Si X és una variable aleatòria normal $N(\mu_X, \sigma_X)$, la seva mitjana mostral \bar{X} de mostres aleatòries simples de mida n és normal

$$N\left(\mu_X, \frac{\sigma_X}{\sqrt{n}}\right).$$

El teorema següent diu que la conclusió del teorema anterior és aproximadament vertadera si la mida n de les mostres aleatòries simples és gran:

Teorema 3.3 (Teorema Central del Límit) Siguin X una variable aleatòria qualsevol d'esperança μ_X i desviació típica σ_X . Quan $n \rightarrow \infty$, la funció de distribució de la seva mitjana mostral \bar{X} de mostres aleatòries simples de mida n tendeix a la d'una variable normal

$$N\left(\mu_X, \frac{\sigma_X}{\sqrt{n}}\right).$$

Com us podeu imaginar, quan un resultat l'anomenen **Teorema Central** de qualche cosa és perquè és molt important.

Normalment aplicarem el Teorema Central del Límit de la manera següent:

Siguin X una variable aleatòria qualsevol d'esperança μ_X i desviació típica σ_X . Si la mida n de les mostres (aleatòries simples) és gran, la mitjana mostral \bar{X} és aproximadament normal $N(\mu_X, \sigma_X/\sqrt{n})$.

Per fixar una fita, en aquest curs entendrem que n és prou gran com per poder aplicar aquest “resultat” quan és més gran o igual que 40, potser menys com més se sembli \bar{X} a una normal i potser més si la \bar{X} és molt diferent d’una normal.

A partir d’ara, sovint cometrem l’abús de llenguatge d’ometre l’adverbi “aproximadament” de l’expressió anterior, i direm simplement que si n és gran, \bar{X} és normal. Però heu de tenir present que aquest “és normal” en realitat vol dir “la seva distribució és aproximadament la d’una variable normal”.

Exemple 3.2 Suposem que tenim una variable aleatòria X de mitjana poblacional $\mu_X = 3$ i desviació típica poblacional $\sigma_X = 0.2$ i que en prenem mostres aleatòries simples de mida 100. Pel Teorema Central del Límit, la distribució de la mitjana mostral \bar{X} és (aproximadament)

$$N\left(3, \frac{0.2}{\sqrt{100}}\right) = N(3, 0.02)$$

Exemple 3.3 Tornem a la situació de l’Exemple 3.1. Teníem les notes guardades en un vector anomenat **tests**. Amb l’histograma següent podem veure que aquestes notes no tenen pinta de seguir una distribució normal.

```
fact.trans=hist(tests,plot=FALSE)$counts[1]/hist(tests,plot=FALSE)$density[1]
hist(tests,col="light blue",xlab="Notes dels tests",
      ylab="Freqüències",main="")
curve(fact.trans*dnorm(x,mean(tests),sd(tests)),col="red",lwd=2,add=TRUE)
```

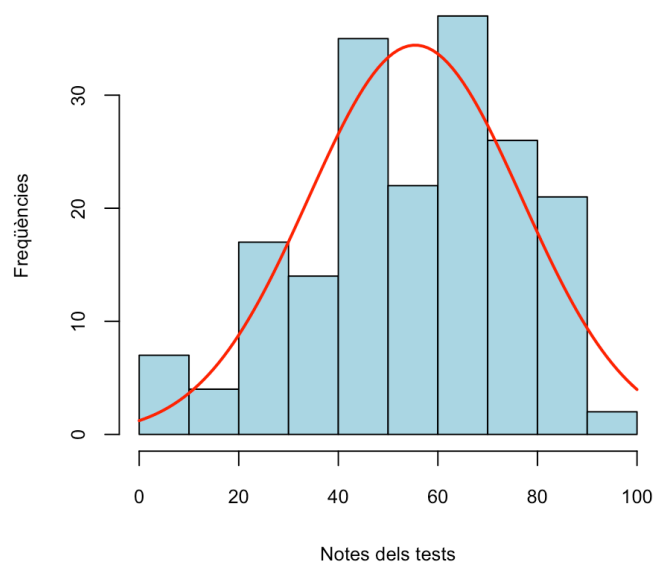


Figura 3.3: Histograma de les notes de tests

A l'Exemple 3.1 també hem construït un vector anomenat **mitjanes** format per 10^5 mitjanes mostrals de mostres aleatòries simples de notes de mida 40. Pel Teorema Central del Límit, aquestes mitjanes mostrals haurien de seguir aproximadament una distribució normal, malgrat que la “població original” (les notes dels tests) no sigui normal. Vegem-ho amb un histograma, on hem afegit la densitat de la normal $N(\mu_X, \sigma_X/\sqrt{n})$ predita pel Teorema Central del Límit.

```
fact.trans.m=hist(mitjanes,plot=FALSE)$counts[1]/hist(mitjanes,plot=FALSE)$c
hist(mitjanes,col="light blue",xlab="Mitjanes",
     ylab="Freqüències",main="")
curve(fact.trans.m*dnorm(x,mu,sigma/sqrt(n)),col="red",lwd=2,add=TRUE)
```

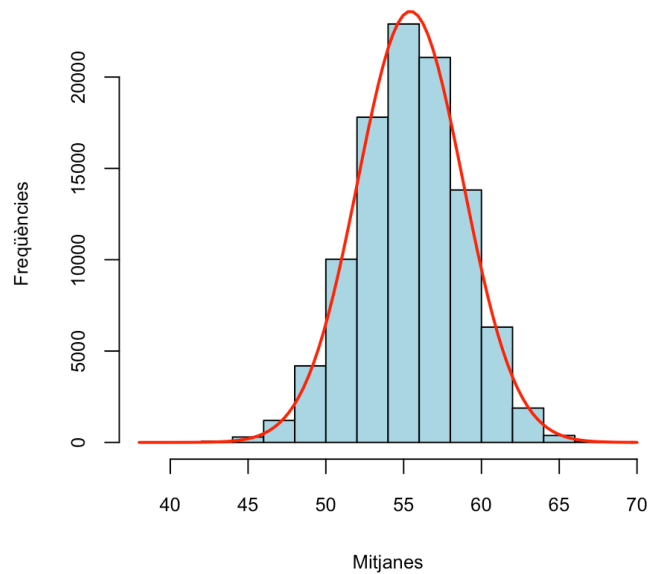


Figura 3.4: Histograma de les mitjanes de mostres de notes de tests

L'exemple següent és un tipus de pregunta que més endavant ens preocuparà molt.

Exemple 3.4 L'alçada d'una espècie de matolls té valor mitjà 115 cm, amb una desviació típica de 25 cm. Si prenem una mostra aleatòria simple de 100 matolls d'aquesta espècie, quina és la probabilitat que la mitjana mostral de les alçades sigui més petita que 110 cm?

Diguem \bar{X} a la variable aleatòria definida per les alçades d'aquests matolls. Pel Teorema Central del Límit, la mitjana mostral \bar{X} de mostres aleatòries simples de 100 alçades segueix una distribució $N(115, 25/\sqrt{100}) = N(115, 2.5)$. Llavors, la probabilitat que ens demanen és

$$P(\bar{X} < 110)$$

que podem calcular amb

```
round(pnorm(110, 115, 2.5), 4)
```

```
## [1] 0.0228
```

Un 2.28% de les mostra aleatòries simples de 100 matolls d'aquesta espècie tenen la mitjana de les alçades més petita que 110 cm.

3.3 Proporció mostral

Quan volem estimar la proporció de subjectes d'una població que tenen una determinada característica, en prenem una mostra i hi calculam la proporció de subjectes amb aquesta característica. Aquesta serà la **proporció mostral** de subjectes amb aquesta característica en la nostra mostra.

Sigui X una variable aleatòria poblacional de Bernoulli amb **probabilitat d'èxit** p_X . És a dir, X pren els valors 1 (èxit) o 0 (fracàs) i p_X és la proporció de subjectes de la població en els quals val 1. Recordau que $E(X) = p_X$ i $\sigma_X = \sqrt{p_X(1 - p_X)}$.

La **proporció mostral** (de mostres aleatòries simples) **de mida** n de X , \hat{p}_X , és la variable aleatòria que consisteix a prendre una mostra aleatòria simple de mida n de X i calcular-ne la proporció d'èxits: és a dir, comptar-hi el nombre total d'èxits i dividir el resultat per n .

Formalment, sigui X_1, \dots, X_n una mostra aleatòria simple de mida n de X . Sigui $S_n = \sum_{i=1}^n X_i$, que és la variable aleatòria que compta el nombre d'èxits en una mostra aleatòria simple de mida n . Aleshores, la **proporció mostral** de mida n de X és

$$\hat{p}_X = \frac{S_n}{n} = \frac{\sum_{i=1}^n X_i}{n}.$$

Recordau que si prenem mostres aleatòries simples, S_n **és una variable aleatòria binomial** $B(n, p_X)$. Però és un doi dir que la proporció mostral \hat{p}_X és una variable aleatòria binomial, ni que només sigui perquè les variables aleatòries binomials prenen valors nombres naturals i els valors que pot prendre \hat{p}_X són fraccions entre 0 i 1.

Fixau-vos que \hat{p}_X és un cas particular de la mitjana mostral \bar{X} , per tant per a les proporcions mostrals val tot el que hem dit per a mitjanes mostrals:

Teorema 3.4 Sigui X una variable aleatòria de Bernoulli amb probabilitat d'èxit p_X .

Aleshores, la proporció mostral de mostres aleatòries simples de mida n de X , \hat{p}_X , satisfà que:

1. $E(\hat{p}_X) = p_X$

2. $\sigma(\hat{p}_X) = \sqrt{\frac{p_X(1 - p_X)}{n}}$

3. Pel Teorema Central del Límit, si la mida n de la mostra és gran, la distribució de \hat{p}_X és aproximadament la d'una variable normal

$$N\left(p_X, \sqrt{\frac{p_X(1 - p_X)}{n}}\right)$$

i per tant

$$\frac{\hat{p}_X - p_X}{\sqrt{p_X(1 - p_X)/n}}$$

és aproximadament $N(0, 1)$.

Els valors de l'esperança i la desviació típica de \hat{p}_X es poden deduir dels de S_n sense necessitat d'invocar els de la mitjana mostral. Com que S_n és $B(n, p_X)$, sabem que $E(S_n) = np_X$ i $\sigma(S_n)^2 = np_X(1 - p_X)$ i per tant

$$E(\hat{p}_X) = E\left(\frac{S_n}{n}\right) = \frac{E(S_n)}{n} = \frac{np_X}{n} = p_X$$

$$\sigma(\hat{p}_X) = \sigma\left(\frac{S_n}{n}\right) = \frac{\sigma(S_n)}{n} = \frac{\sqrt{np_X(1 - p_X)}}{n} = \sqrt{\frac{p_X(1 - p_X)}{n}}$$

Alguns comentaris:

- $E(\hat{p}_X) = p_X$: Esperam que la proporció mostral sigui igual a la proporció poblacional d'èxits (si hi ha, diguem, un 20% d'èxits a la població, quin percentatge d'èxits “esperau” trobar a la mostra? Un 20%, no?).

És a dir, si repetíssim moltes vegades el procés de prendre una mostra aleatòria simple de mida n d'una variable aleatòria de Bernoulli X i calcular-ne la proporció mostral d'èxits, molt probablement la mitjana d'aquestes proporcions mostrals s'acostaria molt a p_X .

- En particular, \hat{p}_X serveix per estimar p_X , naturalment.
- $\sigma(\hat{p}_X) = \sqrt{p_X(1 - p_X)/n}$: la variabilitat dels resultats de \hat{p}_X decreix amb n i tendeix a 0 quan $n \rightarrow \infty$.

Pel que fa a la dependència de $\sigma(\hat{p}_X)$ respecte de p_X si la n és fixada, observau a la Figura 3.5 que $\sqrt{p_X(1 - p_X)}$ creix entre 0 i 0.5 i decreix entre 0.5 i 1, assolint el valor màxim a $p_X = 0.5$.

```
curve(sqrt(x*(1-x)), xlab="p", ylab="", lwd=2)
```

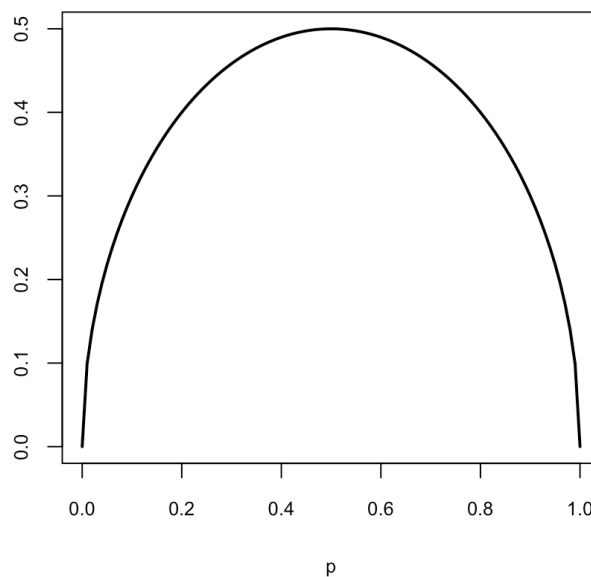


Figura 3.5: Gràfica de $\sqrt{p(1-p)}$

- $\sqrt{p_X(1 - p_X)/n}$ és l'**error estàndard**, o **típic**, de \hat{p}_X . L'estimam amb l'**error estàndard**, o **típic**, de la mostra $\sqrt{\hat{p}_X(1 - \hat{p}_X)/n}$.

- Sovint cometrem l'abús de llenguatge d'ometre l'adverbi “aproximadament” de l'apartat (3) del teorema anterior, i direm simplement que si n és gran, \hat{p}_X és normal. Però, repetim, hem de recordar que aquest “és normal” en realitat vol dir “la seva distribució és aproximadament la d'una variable normal”.

Exemple 3.5 Tornem una altra vegada a la situació dels Exemples 3.1 i 3.3. Traduïm el fitxer de notes de tests en un vector binari: 0 per suspens (haver tret menys de 50) i 1 per aprovat (haver tret 50 o més):

```
# Iniciam totes les notes a 1
aprovs=rep(1,length(tests))
# Posam 0 on la nota del test és suspesa
aprovs[which(tests<50)]=0
```

Aquest vector `aprovs` el podem entendre com els valors sobre la nostra població d'estudiants de la variable poblacional de Bernoulli Y que ens diu si un estudiant aprovà o suspengué els tests. La seva probabilitat poblacional d'èxit (aprovat) p_Y serà la proporció d'estudiants aprovats:

```
p_Y=sum(aprovs)/N
round(p_Y,4)
```

```
## [1] 0.5946
```

Ara n'extreurem 10^5 mostres aleatòries simples de mida $n = 40$, en calcularem les proporcions mostrals d'aprovats i comprovarem si es confirmen les conclusions del teorema anterior.

```
n=40
props.mostrals=replicate(10^5,mean(sample(aprovs,n,rep=TRUE)))
```


La mitjana d'aquest vector de proporcions hauria de ser propera a la proporció poblacional d'aprovat $p_Y = 0.5946$.

```
round(mean(props.mostrals),4)
```

```
## [1] 0.5954
```

Vegem ara la seva desviació típica:

```
round(sd(props.mostrals),4)
```

```
## [1] 0.0775
```

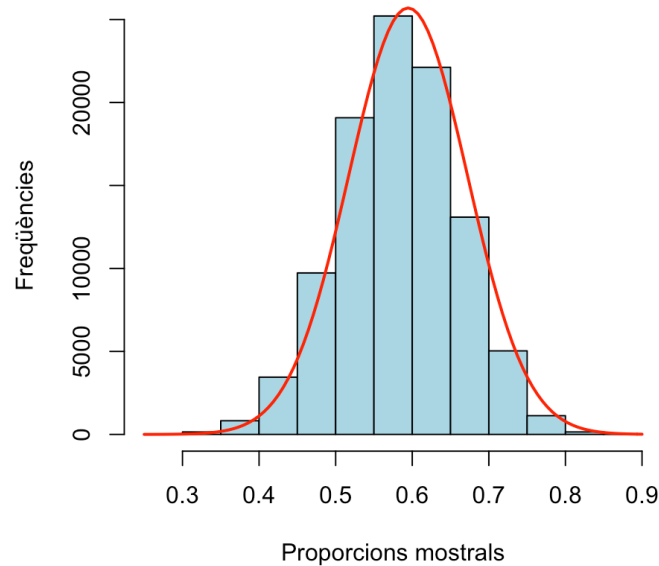
Pel Teorema 3.4, sabem que això hauria de ser proper a $\sqrt{p_Y(1 - p_Y)/n}$

```
round(sqrt(p_Y*(1-p_Y)/n),4)
```

```
## [1] 0.0776
```

I pel Teorema Central del Límit, aquestes proporcions mostrals haurien de seguir aproximadament una distribució normal $N(p_Y, \sqrt{p_Y(1 - p_Y)/n})$. Vegem-ho amb un histograma:

```
fact.trans.p=hist(props.mostrals,plot=FALSE)$counts[1]/hist(props.mostrals,
hist(props.mostrals,col="light blue",xlab="Proporcions mostrals",
ylab="Freqüències",main="Histograma de la mostra de proporcions")
curve(fact.trans.p*dnorm(x,p_Y,sqrt(p_Y*(1-p_Y)/n)),
col="red",lwd=2,add=TRUE)
```

Histograma de la mostra de proporcions

I això que la mida de les mostres, 40, no és especialment gran.

Exemple 3.6 Un 59.1% dels estudiants de la UIB són dones. Hem pres una mostra més o menys aleatòria de 60 estudiants de la UIB i hi hem trobat 40 dones, dos terços. Ens demanam si 40 de 60 és una quantitat raonable de dones en una mostra aleatòria simple d'estudiants de la UIB, o si són moltes (atès que hi esperaríem al voltant d'un 59% de dones).

Aquesta pregunta, que serà molt típica d'aquí a pocs temes, la traduïm en la pregunta següent:

Si prenem una mostra aleatòria simple de 60 estudiants, quina és la probabilitat que la proporció mostral de dones sigui més gran o igual que $2/3$?

La manera més correcta de respondre aquesta qüestió és emprar que el nombre S_{60} de dones en mostres aleatòries simples de 60 estudiants de la UIB segueix de manera exacta una distribució binomial $B(60, 0.591)$. Com que el 66.67% de la pregunta en realitat representa 40 dones, la probabilitat demanada és exactament

```
round(1-pbinom(39,60,0.591),4)
```

```
## [1] 0.1441
```

Recordau que si X és una variable aleatòria discreta que pren valors enters, com ara la binomial, $P(X \geq 40) = 1 - P(X \leq 39)$.

Això ens diu que, de mitjana, 1 de cada 7 mostres aleatòries simples de 60 estudiants de la UIB conté almenys 40 dones. Naturalment, això només és exacte si la proporció poblacional “59.1%” és exacta.

Una altra opció seria aprofitar el Teorema Central del Límit, segons el qual la proporció mostrat \hat{p}_X de dones en mostres aleatòries simples de 60 estudiants de la UIB segueix una distribució aproximadament normal amb $\mu = 0.591$ i

$$\sigma = \sqrt{\frac{0.591(1 - 0.591)}{60}} = 0.0635$$

Per tant, la probabilitat que $\hat{p}_X \geq 2/3$ és (recordau, aproximadament)

```
round(1-pnorm(2/3,0.591,sqrt(0.591*(1-0.591)/60)),4)
```

```
## [1] 0.1166
```

L’aproximació seria més bona si haguéssim efectuat la correcció de continuïtat. Diguem Y a la normal $N(0.591, 0.0635)$. Com que $\hat{p}_X = S_{60}/60$, seria millor aproximar

$$P(S_{60} \geq 40) = 1 - P(S_{60} \leq 39)$$

per

$$1 - P(Y \leq 39.5/60)$$

```
round(1-pnorm(39.5/60,0.591,sqrt(0.591*(1-0.591)/60)),4)
```

```
## [1] 0.1444
```

En el cas de la proporció mostral, de vegades considerarem que s'han pres **mostres aleatòries sense reposició**. En aquest cas, la distribució del nombre d'èxits S_n en una mostra segueix una distribució hipergeomètrica. D'aquí deduïm, exactament igual que en el cas de mostres aleatòries simples, que seguim tenint que $E(\hat{p}_X) = p_X$, però ara, si N és la mida de la població,

$$\sigma(\hat{p}_X) = \sqrt{\frac{p_X(1-p_X)}{n}} \cdot \sqrt{\frac{N-n}{N-1}}.$$

Recordau que al factor

$$\sqrt{\frac{N-n}{N-1}}$$

que transforma $\sigma(\hat{p}_X)$ per a mostres aleatòries simples en la desviació típica de \hat{p}_X per a mostres aleatòries sense reposició li diem el **factor de població finita**, i és el que transformava la desviació típica d'una variable binomial (que compta èxits en mostres aleatòries simples) en la desviació típica d'una variable hipergeomètrica (que compta èxits en mostres aleatòries sense reposició).

Però recordau que si la mida de la població N és molt gran comparat amb n , podem suposar que una mostra aleatòria sense reposició és simple.

Exemple 3.7 Tornem a la situació de l'Exemple 3.5. Què passa si prenem les mostres aleatòries de notes de tests sense reposició?

Prenguem ara 10^5 mostres aleatòries sense reposició de 40 notes de tests.

```
props.norep=replicate(10^5,mean(sample(aprovs,n)))
```

Un altre cop, la mitjana d'aquest vector de proporcions mostrals hauria de ser propera a la proporció poblacional d'aprovals $p_Y = 0.5946$.

```
round(mean(props.norep),4)
```

```
## [1] 0.5946
```

Calculem ara la desviació típica d'aquest vector:

```
round(sd(props.norep),4)
```

```
## [1] 0.069
```

Pel que acabam d'explicar, la desviació típica d'aquest vector de proporcions mostrals de mostres sense reposició hauria de ser molt propera a

$$\sqrt{\frac{p_Y(1-p_Y)}{n}} \cdot \sqrt{\frac{N-n}{N-1}}$$

on N és la mida de la població, és a dir, la longitud del vector `aprovs`, i n la mida de les mostres. Vegem si és veritat:

```
round(sqrt(p_Y*(1-p_Y)/n)*sqrt((N-n)/(N-1)),4)
```

```
## [1] 0.0689
```

Però si prenem mostres aleatòries sense reposició i la població no és molt més gran que les mostres, ja no es pot aplicar el Teorema Central del Límit: encara que les mostres siguin grans, la proporció mostral no té per què ser aproximadament normal.

3.4 Variància mostral

Donada una variable aleatòria X , direm:

- La **variància mostral** (de mostres aleatòries simples) **de mida** n , \tilde{S}_X^2 , a la variable aleatòria que consisteix a prendre una mostra aleatòria simple de mida n de X i calcular la variància mostral dels seus valors.
- **Desviació típica mostral** (de mostres aleatòries simples) **de mida** n , \tilde{S}_X , a la variable aleatòria que consisteix a prendre una mostra aleatòria simple de mida n de X i calcular la desviació típica mostral dels seus valors.

Formalment, sigui X_1, \dots, X_n una mostra aleatòria simple de mida n d'una variable aleatòria X . Aleshores

$$\tilde{S}_X^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}, \quad \tilde{S}_X = +\sqrt{\tilde{S}_X^2}$$

A més, de tant en tant també farem servir la **variància** i la **desviació típica** “a seques”:

$$S_X^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n} = \frac{(n-1)}{n} \tilde{S}_X^2$$

$$S_X = +\sqrt{S_X^2}$$

La variància (a seques) admet la expressió senzilla següent:

$$S_X^2 = \frac{\sum_{i=1}^n X_i^2}{n} - \bar{X}^2$$

En efecte:

$$\begin{aligned}
\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n} &= \frac{1}{n} \sum_{i=1}^n (X_i^2 - 2\bar{X}X_i + \bar{X}^2) \\
&= \frac{1}{n} \left(\sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + n\bar{X}^2 \right) \\
&= \frac{\sum_{i=1}^n X_i^2}{n} - 2\bar{X} \frac{\sum_{i=1}^n X_i}{n} + \frac{n\bar{X}^2}{n} \\
&= \frac{\sum_{i=1}^n X_i^2}{n} - 2\bar{X} \cdot \bar{X} + \bar{X}^2 = \frac{\sum_{i=1}^n X_i^2}{n} - \bar{X}^2
\end{aligned}$$

Tenim els dos resultats següents. El primer ens diu que **esperam** que la variància mostral d'una mostra aleatòria simple de X valgui la variància σ_X^2 de X , en el sentit usual que si prenem mostres aleatòries simples de X de mida n gran i calculam les seves variàncies mostrals, molt probablement obtenim de mitjana un valor molt proper a σ_X^2 .

Teorema 3.5 Si X és una variable aleatòria de desviació típica σ_X i \tilde{S}_X és la seva variància mostral de mida n ,

$$E(\tilde{S}_X^2) = \sigma_X^2$$

per a qualsevol n .

I per tant **no esperam** que la variància “a seques” d'una mostra aleatòria simple valgui σ_X^2 . Això ho podeu comprovar fàcilment, perquè S_X^2 s'obté a partir de \tilde{S}_X^2 canviant el denominador,

$$S_X^2 = \frac{n-1}{n} \tilde{S}_X^2$$

i per tant

$$E(S_X^2) = \frac{n-1}{n} E(\tilde{S}_X^2) = \frac{n-1}{n} \sigma_X^2$$

El segon resultat ens diu que **si la variable X és normal**, un múltiple adequat de \tilde{S}_X^2 té distribució mostral coneguda, la qual cosa ens permetrà calcular probabilitats d'esdeveniments relatius a \tilde{S}_X^2 .

Teorema 3.6 Si X es $N(\mu_X, \sigma_X)$ i \tilde{S}_X és la seva variància mostral de mida n , la variable aleatòria

$$\frac{(n-1)\tilde{S}_X^2}{\sigma_X^2}$$

té distribució coneguda: χ_{n-1}^2 (es llegeix **khi quadrat amb $n-1$ graus de llibertat**).

La lletra χ en català es llegeix *khi*; en castellà, *ji*; i en anglès, *chi*, pronunciat *xai*.



De la distribució χ_ν^2 , on ν són els **graus de llibertat**, heu de saber que:

- Per definició, és la distribució de la suma dels quadrats de ν variables aleatòries normals estàndard independents. És a dir, si Z_1, Z_2, \dots, Z_ν són variables $N(0, 1)$ independents, la variable

$$Z_1^2 + Z_2^2 + \cdots + Z_\nu^2$$

té distribució χ_ν^2 .

- Per tant, és una distribució contínua
- La ν és el paràmetre del que depèn la seva densitat
- Amb R és `chisq`
- Si X_ν és una variable aleatòria amb distribució χ_ν^2 , aleshores $\mu_{X_\nu} = \nu$ i $\sigma_{X_\nu}^2 = 2\nu$
- Per a ν petits, la distribució d'una χ_ν^2 és asimètrica amb una cua a la dreta.

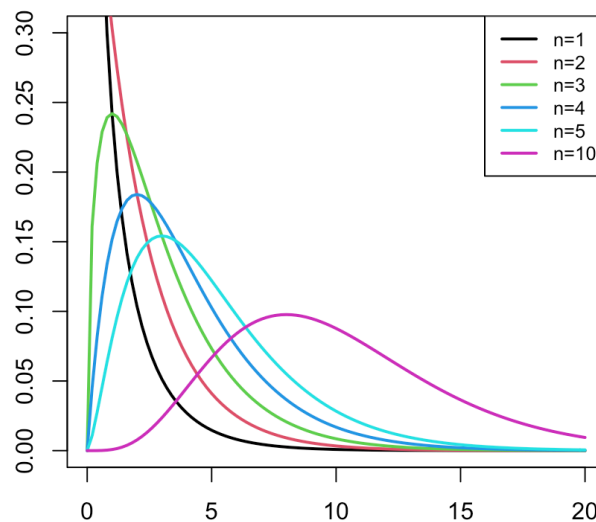
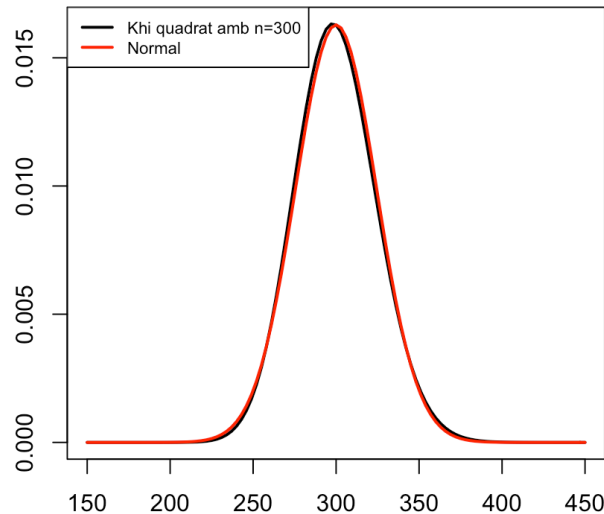


Figura 3.6: Algunes densitats de variables χ^2

- A mida que ν creix, i atès que és la distribució d'una suma de ν variables aleatòries independents, pel Teorema Central del Límit es va aproximant a una distribució normal $N(n, \sqrt{2n})$.

Figura 3.7: χ^2 vs normal

Tornem un instant a això dels *graus de llibertat*. Per què diem que la variància (mostral o a seques) té $n - 1$ graus de llibertat?

Doncs perquè si volem construir un conjunt de n nombres x_1, \dots, x_n que tenguin variància un valor donat, posem y_0 , aleshores en principi podem escollir $n - 1$ d'ells, x_1, \dots, x_{n-1} , com vulguem i aleshores el darrer, x_n , queda bastant fixat. En matemàtiques això se sol expressar dient que “tenim $n - 1$ graus de llibertat a l'hora d'escollir x_1, \dots, x_n amb variància fixada y_0 ”.

En efecte, si fixam el valor $y_0 \geq 0$ de la variància i volem trobar x_1, \dots, x_n tals que

$$y_0 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n} = \frac{\sum_{i=1}^n x_i^2}{n} - \bar{x}^2$$

vegem que per a qualssevol valors de x_1, \dots, x_{n-1} , el valor de x_n queda fixat per una equació quadràtica:

$$\begin{aligned}
ny_0 &= \sum_{i=1}^n x_i^2 - n\bar{x}^2 = \sum_{i=1}^n x_i^2 - n\left(\frac{\sum_{i=1}^n x_i}{n}\right)^2 \\
&= \sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n} \\
&= \frac{1}{n} \left(n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2 \right) \\
&= \frac{1}{n} \left(n \sum_{i=1}^{n-1} x_i^2 + n\mathbf{x}_n^2 - \left(\sum_{i=1}^{n-1} x_i \right)^2 \right. \\
&\quad \left. - 2 \left(\sum_{i=1}^{n-1} x_i \right) \mathbf{x}_n - \mathbf{x}_n^2 \right) \\
&= \frac{1}{n} \left((n-1)\mathbf{x}_n^2 - 2 \left(\sum_{i=1}^{n-1} x_i \right) \mathbf{x}_n \right. \\
&\quad \left. + n \sum_{i=1}^{n-1} x_i^2 - \left(\sum_{i=1}^{n-1} x_i \right)^2 \right)
\end{aligned}$$

d'on (multiplicant els dos costats de la igualtat per n i dividint-los per $n - 1$) obtenim, finalment, l'equació de segon grau en \mathbf{x}_n

$$\mathbf{x}_n^2 - \frac{2 \sum_{i=1}^{n-1} x_i}{n-1} \mathbf{x}_n + \frac{n \sum_{i=1}^{n-1} x_i^2 - \left(\sum_{i=1}^{n-1} x_i \right)^2 - n^2 y_0^2}{n-1} = 0$$

Per tant, fixat y_0 i un cop escollits x_1, \dots, x_{n-1} , el darrer valor x_n ha de ser per força una solució d'aquesta equació de segon grau.

Fixau-vos que aquesta equació no sempre té solució real, perquè pot tenir el discriminant negatiu. Per tant exageràvem un poc dient que podíem triar x_1, \dots, x_{n-1} "com vulguem". Per exemple, si voleu que la variància sigui 0 i preneu x_1, \dots, x_{n-1} no tots iguals, podeu estar ben segurs que no trobareu cap x_n que satisfaci aquesta equació: per tenir variància 0, x_1, \dots, x_n han de ser tots iguals. Però el que ha de quedar clar és que un cop escollits x_1, \dots, x_{n-1} , el valor de x_n ja no pot ser qualsevol, pot prendre com a màxim dos valors diferents.

Anau alerta:

- Si la variable poblacional X no és normal, la conclusió del Teorema 3.6 no és vertadera.
- Encara que X sigui normal, $E(\tilde{S}_X) \neq \sigma_X$.
- Ja ho hem comentat abans, però ho repetim: si S_X^2 és la variància “a seques” (dividint per n en comptes de per $n - 1$), $E(S_X^2) \neq \sigma_X^2$.

Exemple 3.8 Suposem que el pes en néixer dels nadons segueix una distribució normal, i s'estima que la seva desviació típica (en g) és 800. Hem anotat els pesos de tots els recent nats amb SIDA d'una ciutat durant 2 anys. El teniu al vector `pesos.SIDA` següent:

```
pesos.SIDA=c(2466,3941,2807,3118,2098,3175,3515,3317,3742,3062,
             3033,2353,2013,3515,3260,2892,1616,4423,3572,2750,
             2807,2807,3005,3374,2722,2495,3459,3374,1984,2495,
             3062,3005,2608,2353,4394,3232,2013,2551,2977,3118,
             2637,1503,2438,2722,2863,2013,3232,2863)
```

Quina és la probabilitat que una mostra (aleatòria simple) de pesos de recent nats de la mateixa mida que aquesta tengui una desviació típica mostral més petita que la d'aquesta mostra?

La variable d'interès és \bar{X} : “Prenem un recent nat i pesam el seu pes en g”. Ens diuen que és normal amb $\sigma = 800$. Pel que fa a la nostra mostra de pesos, la seva mida n és

```
n=length(pesos.SIDA)
```

```
n
```

```
## [1] 48
```

i la seva desviació típica mostral és

```
s.tilda=round(sd(pesos.SIDA),1)
s.tilda
```

```
## [1] 623.4
```

Sigui \tilde{S}_X la desviació típica mostral de mida 48 de la variable X . Ens demanen $P(\tilde{S}_X < 623.4)$. Això tal qual no ho sabem calcular, perquè no sabem la funció de distribució de \tilde{S}_X . Però sí que sabem la de

$$\frac{(n-1)\tilde{S}_X^2}{\sigma_X^2} = \frac{47\tilde{S}_X^2}{800^2}$$

Aquesta variable té distribució χ_{47}^2 . Per tant el que hem de fer és traduir la probabilitat que volem calcular en termes d'aquesta variable:

$$P(\tilde{S}_X < 623.4) = P\left(\frac{47\tilde{S}_X^2}{800^2} < \frac{47 \cdot 623.4^2}{800^2}\right) = P(\chi_{47}^2 < 28.54)$$

i això val

```
round(pchisq(round((n-1)*s.tilda^2/800^2,2),n-1),4)
```

```
## [1] 0.0153
```

Per tant, només un 1.5% de les mostres aleatòries simples de 47 recent nats tenen una desviació típica mostral més petita que la de la nostra mostra de recent nats amb SIDA.

3.5 La distribució t de Student

Recordau que si la variable poblacional X és $N(\mu_X, \sigma_X)$ i prenem mostres aleatòries simples de mida n , la variable

$$\frac{\bar{X} - \mu_X}{\sigma_X / \sqrt{n}}$$

és normal estàndard. Des del punt de vista teòric això és útil per obtenir fórmules, però normalment no ens serveix per calcular la probabilitat que a \bar{X} li passi alguna cosa, perquè gairebé mai sabrem la desviació típica poblacional σ_X . Què passa si l'estimam per mitjà de \tilde{S}_X amb la mateixa mostra amb la qual calculam \bar{X} ? Doncs que el resultat següent ens salva el dia, perquè la variable que obtenim té distribució coneguda.

Teorema 3.7 *Sigui X una variable $N(\mu_X, \sigma_X)$. Siguin \bar{X} i \tilde{S}_X les seves mitjana mostral i desviació típica mostral de mostres aleatòries simples de mida n , respectivament. Aleshores, la variable aleatòria*

$$T = \frac{\bar{X} - \mu_X}{\tilde{S}_X / \sqrt{n}}$$

*segueix una distribució coneguda, anomenada **t de Student amb $n - 1$ graus de llibertat**, t_{n-1} .*

Al denominador \tilde{S}_X / \sqrt{n} li diem l'**error estàndard**, o **típic**, de la mostra: estima l'error estàndard σ_X / \sqrt{n} de \bar{X} .

De la distribució t de Student amb ν graus de llibertat, t_ν , heu de saber que:

- És contínua
- Amb R és `t`
- El **nombre de graus de llibertat** ν és un paràmetre del que depèn la seva distribució
- Si T_ν és una variable amb distribució t_ν , aleshores $\mu_{T_\nu} = 0$ i $\sigma_{T_\nu}^2 = \nu / (\nu - 2)$ (en realitat això només és veritat si $\nu \geq 3$, però no fa falta recordar-ho).
- La funció de densitat d'una variable T_ν és simètrica al voltant de 0 (com la d'una $N(0, 1)$):

$$P(T_\nu \leq -x) = P(T_\nu \geq x) = 1 - P(T_\nu \leq x)$$

Per tant 0 també és la seva mediana.

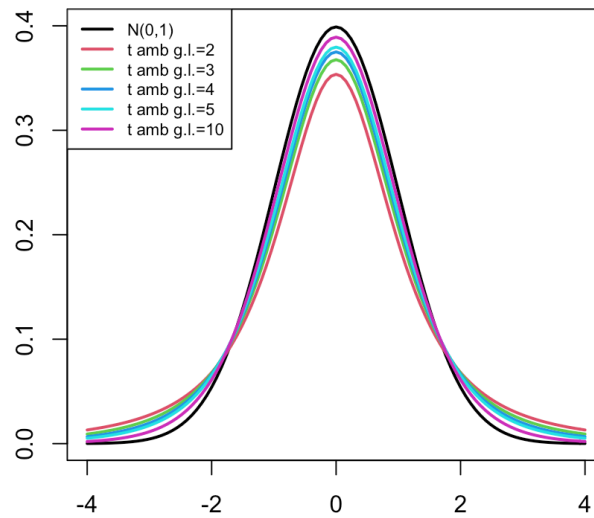


Figura 3.8: Densitats d'algunes t de Student

- Si ν és gran, la distribució d'una variable T_ν és aproximadament la d'una $N(0, 1)$ però amb més variància (perquè $\nu/(\nu - 2) > 1$) i per tant un poc més aplatada.

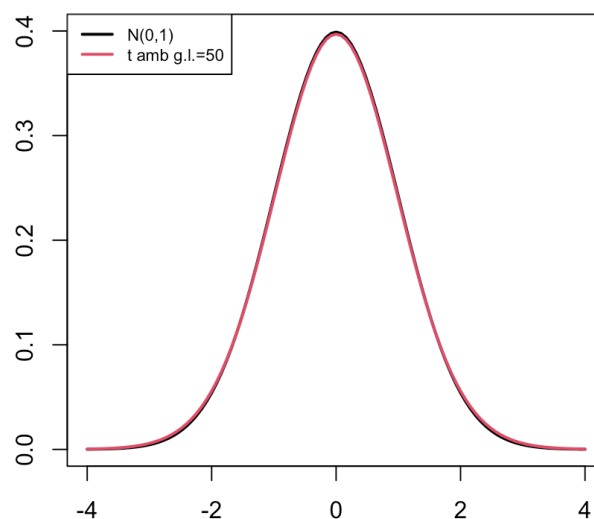


Figura 3.9: t amb molts graus de llibertat vs Normal estàndard

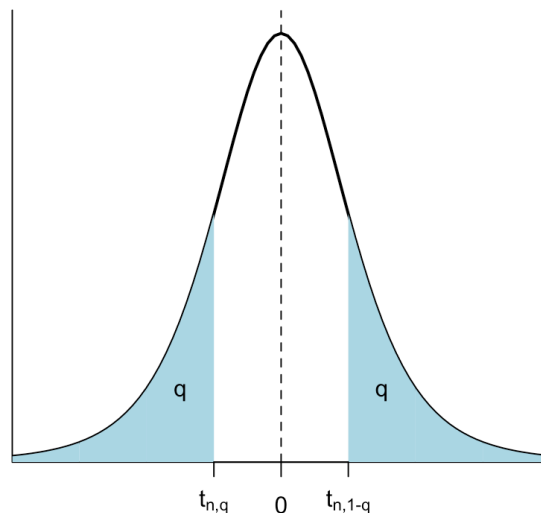
El fet que una t de Student sigui més aplatada que una normal estàndard Z implica que les cues de la t tenen major probabilitat que les de Z (fixau-vos que als gràfics anteriors els extrems de les densitats de les t estan per damunt dels de la de Z), la qual cosa es tradueix en el fet que és més probable obtenir valors lluny del 0 amb una t de Student que amb una $N(0, 1)$.

Indicarem amb $t_{\nu,q}$ el q -quantil d'una variable aleatòria T_{ν} que segueix una distribució t_{ν} . És a dir, $t_{\nu,q}$ és el valor tal que

$$P(T_{\nu} \leq t_{\nu,q}) = q$$

Per la simetria de la distribució t_{ν} ,

$$t_{\nu,q} = -t_{\nu,1-q}.$$



Hi ha algunes propietats dels quantils de la t de Student que heu de saber, per poder aplicar-les quan no tingueu a l'abast R o una apli per calcular quantils:

- $t_{\nu,q} \approx z_q$ si ν és molt gran, posem $\nu \geq 200$. Per exemple

```
qt(0.975,200) # t_{200,0.975}
```



```
## [1] 1.971896
```

```
qnorm(0.975) # z_0.975
```

```
## [1] 1.959964
```

- $t_{\nu,0.95}$ (per a $10 \leq \nu \leq 200$) està entre 1.65 i 1.8; ho podeu aproximar $t_{n,0.95} \approx 1.7$
- $t_{n,0.975}$ (per a $10 \leq n \leq 200$) està entre 1.97 i 2.2; ho podeu aproximar $t_{n,0.975} \approx 2$

Comprovau amb R les afirmacions sobr els quantils de la t de Student dels darrers dos punts.

Abans de tancar aquesta secció, recordau que, donada una variable aleatòria X , no heu de confondre:

- **Desviació típica (o estàndard) de la variable aleatòria, σ_X :** El paràmetre poblacional, normalment desconegut
- **Desviació típica (o estàndard) (sigui mostral, \tilde{S}_X , o a seques, S_X) d'una mostra:** L'estadístic que calculam sobre la mostra i que quantifica la dispersió de la mostra
- **Error típic (o estàndard) d'un estimador:** La desviació típica de la variable aleatòria que defineix l'estimador, normalment desconeguda
- **Error típic (o estàndard) d'una mostra:** Estimació de l'error típic de la mitjana mostral (o de la proporció mostral) a partir d'una mostra; servirà per calcular intervals de confiança. És \tilde{S}_X/\sqrt{n} .

3.6 “Bons” estimadors

3.6.1 Estimadors no esbiaixats

Un estimador puntual $\hat{\theta}$ d'un paràmetre poblacional θ és **no esbiaixat** (**insesgado**, en castellà) quan el seu valor esperat és precisament el valor poblacional del paràmetre, és a dir, quan

$$E(\hat{\theta}) = \theta$$

El **biaix** d'un estimador $\hat{\theta}$ d'un paràmetre θ és la diferència $E(\hat{\theta}) - \theta$

Exemples: Ja hem vist a les seccions anteriors que

- $E(\bar{X}) = \mu_X$. Per tant, \bar{X} és un estimador no esbiaixat de μ_X .
- $E(\hat{p}_X) = p_X$. Per tant, \hat{p}_X és un estimador no esbiaixat de p_X .
- $E(\tilde{S}_X^2) = \sigma_X^2$. Per tant, \tilde{S}_X^2 és un estimador no esbiaixat de σ_X^2 .
- Com que $S_X^2 = \frac{n-1}{n} \tilde{S}_X^2$, tenim que $E(S_X^2) = \frac{n-1}{n} \sigma_X^2$. Per tant, en aquest cas, S_X^2 és un **estimador esbiaixat** de σ_X^2 , amb biaix

$$\mu_{S_X^2} - \sigma_X^2 = \frac{n-1}{n} \sigma_X^2 - \sigma_X^2 = -\frac{\sigma_X^2}{n} \xrightarrow{n \rightarrow \infty} 0$$

Diem en aquest cas que **el biaix tendeix a 0**.

- $E(\tilde{S}_X), E(S_X) \neq \sigma_X$ ni tan sols quan X és normal. Per tant, \tilde{S}_X i S_X són estimadors **esbiaixats** de σ_X . Quan X és normal, el seu biaix tendeix a 0.

Recordau que, en general, $E(X^2) \neq E(X)^2$. Per tant, no hauríem d'esperar que $E(\tilde{S}_X) = \sqrt{E(\tilde{S}_X^2)}$.

Per si qualche dia us cal saber-ho, si X és normal

$$E(\tilde{S}_n) = \begin{cases} \sigma_X \sqrt{\frac{2}{(n-1)\pi}} \cdot \frac{2^{n-2}(\frac{n}{2}-1)!^2}{(n-2)!} & \text{si } n \text{ és parell} \\ \sigma_X \sqrt{\frac{2\pi}{n-1}} \cdot \frac{(n-2)!}{2^{n-2}(\frac{n-1}{2}-1)!^2} & \text{si } n \text{ és imparell} \end{cases}$$

Per tant, per obtenir un estimador no esbiaixat per a σ_X només heu de dividir \tilde{S}_X pel factor que acompanyi la σ_X a la dreta d'aquesta fórmula.

3.6.2 Estimadors eficients

Donats dos estimadors $\hat{\theta}_1, \hat{\theta}_2$ del mateix paràmetre θ , direm que $\hat{\theta}_1$ és **més eficient**, o **més precís**, que $\hat{\theta}_2$ quan l'error típic de $\hat{\theta}_1$ és més petit que el de $\hat{\theta}_2$:

$$\sigma(\hat{\theta}_1) < \sigma(\hat{\theta}_2).$$

Normalment, només comparam l'eficiència de dos estimadors quan són no esbiaixats (o, com a molt, quan el seu biaix tendeix a 0). En aquest cas, que $\hat{\theta}_1$ sigui més eficient que $\hat{\theta}_2$ significa que la seva variabilitat és més petita i que per tant les estimacions són **més precises** amb $\hat{\theta}_1$: es concentren més al voltant del paràmetre θ que volem estimar.

Exemples:

- Si X és normal, \bar{X} és l'estimador no esbiaixat més eficient de la mitjana poblacional μ_X .
- Si X és Bernoulli, \hat{p}_X és l'estimador no esbiaixat més eficient de la proporció poblacional p_X .
- Si X és normal, \tilde{S}_X^2 és l'estimador no esbiaixat més eficient de la variància poblacional σ_X^2 .

Exemple 3.9 Sigui X una variable aleatòria normal $N(\mu_X, \sigma_X)$. Considerem la mediana $Me = Q_{0.5}$ d'una mostra aleatòria simple de X com a estimador puntual de μ_X , que coincideix amb la mediana de X per la simetria de les variables normals.

Resulta que $E(Me) = \mu_X$ però

$$\sigma^2(Me) \approx \frac{\pi}{2} \cdot \frac{\sigma_X^2}{n} \approx 1.57 \cdot \frac{\sigma_X^2}{n} = 1.57\sigma_{\bar{X}}^2$$

Per tant, si X és normal, la mediana Me també és un estimador no esbiaixat de μ_X , però és menys eficient que \bar{X} . Per això preferim emprar la mitjana mostral per estimar μ_X .

Hem dit que si la població és normal, \tilde{S}_X^2 és l'estimador no esbiaixat més eficient de la variància poblacional σ_X^2 . La variància a seques

$$S_X^2 = \frac{(n-1)}{n} \tilde{S}_X^2$$

és més eficient, perquè

$$\sigma(S_X^2) = \sqrt{\frac{(n-1)}{n}} \sigma(\tilde{S}_X^2) < \sigma(\tilde{S}_X^2),$$

però és un estimador esbiaixat de σ_X^2 , amb biaix que tendeix a 0.

Si n és petit, és millor fer servir la variància mostral \tilde{S}_X^2 per estimar la variància, ja que el biaix pot desplaçar substancialment l'estimació, però si n és gran, el biaix de S_X^2 ja és petit i es pot fer servir S_X^2 : de fet, si n és molt gran, dividir per n o per $n-1$ no varia gaire el resultat i per tant \tilde{S}_X^2 i S_X^2 donen valors molt semblants.

3.6.3 Estimadors màxim versemblants

Un estimador d'un paràmetre és **màxim versemblant** quan, aplicat a una mostra aleatòria simple, dóna el valor del paràmetre que fa màxima la probabilitat d'obtenir aquesta mostra.

En realitat, l'estimació màxim versemblant d'un paràmetre el que fa màxim és el producte dels valors de la funció densitat de la variable aleatòria poblacional aplicada als elements de la mostra. Quan la variable aleatòria és discreta, això coincideix amb el que hem dit, perquè la probabilitat d'obtenir un valor concret és la funció densitat aplicada a aquest valor. Però quan la variable aleatòria poblacional és contínua, la probabilitat d'obtenir una mostra concreta és sempre 0 i no té sentit parlar de maximitzar aquest 0. Per això es pren la funció densitat.

En aquest curs no ens complicarem la vida i entendrem que el que maximitzam és la probabilitat d'obtenir la mostra.

Exemple 3.10 Suposem que tenim una variable aleatòria Bernoulli X de probabilitat d'èxit p_X desconeguda. Donada una mostra aleatòria simple x_1, \dots, x_n de X , siguin \hat{p}_x la seva proporció mostral i

$$P(x_1, \dots, x_n \mid p_X = p)$$

la probabilitat d'obtenir la mostra quan la probabilitat poblacional p_X és igual p . Un estimador per a p_X és màxim versemblant quan, aplicat a cada mostra aleatòria simple x_1, \dots, x_n de X , ens dóna el valor de p que fa que

$$P(x_1, \dots, x_n \mid p_X = p)$$

sigui el màxim possible.

Quin creieu que és l'estimador màxim versemblant de p_X ? Doncs sí, la proporció mostral \hat{p}_X .

Teorema 3.8 El valor de p per al qual $P(x_1, \dots, x_n \mid p_X = p)$ és màxim és \hat{p}_x .

La demostració és senzilla. Suposau que dins x_1, \dots, x_n hi ha m 1s i $n - m$ 0s, de manera que $\hat{p}_X = m/n$. Aleshores, la probabilitat d'obtenir x_1, \dots, x_n és

$$P(x_1, \dots, x_n \mid p_X = p) = p^m (1 - p)^{n-m}$$

Per trobar el valor de p que fa aquest probabilitat màxima, derivau respecte de p i estudiau el signe de la derivada, i concloureu que el màxim es dóna efectivament a $p = m/n$.

La proporció \hat{p}_x és el valor que fa màxima la probabilitat d'obtenir la nostra mostra, no a l'enrevés: **No és el valor més probable de p_X condicionat a la nostra mostra.** Vaja, no confongueu

$$P(x_1, \dots, x_n \mid p_X = p) \text{ amb } P(p_X = p \mid x_1, \dots, x_n).$$

D'això darrer no en sabem trobar el màxim sense alguna hipòtesi sobre la distribució de probabilitat dels valors possibles de p_X .

Alguns altres estimadors màxim versemblants:

- \bar{X} és l'estimador màxim versemblant del paràmetre λ d'una variable aleatòria Poisson
- \bar{X} és l'estimador màxim versemblant de la mitjana μ_X d'una variable aleatòria normal
- S_X^2 i S_X (la variància i desviació típica a seques, no les mostrals!) són els estimadors màxim versemblants de la variància σ_X^2 i la desviació típica σ_X d'una variable aleatòria normal

3.7 Estimació de poblacions

3.7.1 Estimació de poblacions numerades

Exemple 3.11 Un dia vaig voler estimar quants taxis hi havia a Palma. Per fer-ho, assegut en un bar del Passeig Marítim vaig apuntar les llicències dels 40 primers taxis que passaren. Els entraré directament en un vector de R.

```
taxis=c(1217,600,883,1026,150,715,297,137,508,134,38,961,538,1154,
        314,1121,823,158,940,99,977,286,1006,1207,264,1183,1120,
        498,606,566,1239,860,114,701,381,836,561,494,858,187)
sort(taxis)
```

```
## [1] 38 99 114 134 137 150 158 187 264 286 297 314 381 494
## [16] 508 538 561 566 600 606 701 715 823 836 858 860 883 940
## [31] 977 1006 1026 1120 1121 1154 1183 1207 1217 1239
```

Puc estimar quants taxis hi ha a Palma a partir d'aquesta mostra? Us pot semblar una beneitura de pregunta, però aquest és un problema de rellevància històrica, com podeu consultar en [aquest article](#).

La solució d'aquest problema és donada pel resultat següent:

Teorema 3.9 *Sigui X una variable aleatòria **uniforme** sobre $\{1, 2, \dots, N\}$ (és a dir, X pot prendre tots els valors entre 1 i N , tots amb la mateixa probabilitat $1/N$), i sigui x_1, \dots, x_n una mostra aleatòria sense reposició de X . Sigui $m = \max(x_1, \dots, x_n)$. Aleshores, l'estimador no esbiaixat més eficient de N és*

$$\widehat{N} = m + \frac{m - n}{n}$$

Suposau que teniu x_1, \dots, x_n ordenats en ordre creixent, $x_1 < \dots < x_n$, de manera que $x_n = m$. Calculem la longitud mitjana dels “forats” a l'esquerra de cada valor x_i .

- A l'esquerra de x_1 hi “falten” els nombres $1, 2, \dots, x_1 - 1$, per tant hi ha un forat de $x_1 - 1$ nombres.
- Entre x_1 i x_2 hi “falten” els nombres $x_1 + 1, \dots, x_2 - 1$, per tant hi ha un forat de $x_2 - x_1 - 1$ nombres.
- En general, entre cada x_{i-1} i x_i hi ha un forat de $x_i - x_{i-1} - 1$ nombres (hi “falten” els nombres $x_{i-1} + 1, \dots, x_i - 1$)

Per tant, la mitjana de les longituds d'aquests “forats” és

$$\frac{(x_1 - 1) + (x_2 - x_1 - 1) + \cdots + (x_n - x_{n-1} - 1)}{n} \\ = \frac{x_n - n}{n} = \frac{m - n}{n}$$

El que fa l'estimador \widehat{N} és sumar al màxim de la mostra, m , aquesta longitud mitjana dels forats entre membres de la mostra. És a dir, estimam que la mida de la població és tal que a la dreta del màxim de la nostra mostra hi ha un “forat” de mida la mitjana dels forats de la mostra.

Exemple 3.12 Continuem amb l'Exemple 3.11. Emprant la fórmula anterior, obtenim

```
max(taxis)+(max(taxis)-length(taxis))/length(taxis)
```

```
## [1] 1268.975
```

la qual cosa em permeté estimar que hi havia 1269 taxis a Palma. En realitat, consultant la web de l'Ajuntament, després vaig saber que en aquell moment n'hi havia 1246.

Exemple 3.13 Fem un experiment. Generarem a l'atzar una mida N d'una població grandeta, i suposarem que els individus de la població estan numerats d'1 a N . A continuació, prendrem 100 mostres aleatòries sense reposició de la nostra població i amb cada una d'aquestes mostres estimarem la N emprant la fórmula que hem donat. Al final, calcularem la mitjana d'aquestes estimacions i la compararem amb el valor real de N , que no descobrirem fins el final.

Perquè l'experiment sigui reproduïble, fixarem la llavor d'aleatorietat, però perquè no cregueu que fem trampes amb aquesta llavor, el que farem serà generar a l'atzar la llavor d'aleatorietat amb la funció `sample`.

```
Llavor=sample(1000,1)
Llavor
```



```
## [1] 580
```

```
set.seed(Llavor)
```

Ara generam la mida N de la població com un nombre a l'atzar entre 5000 i 10000.

```
N=sample(5000:10000,1)
```

Suposarem per tant que hi ha N individus a la nostra població, numerats de l'1 a l'N. Ara generarem 100 mostres aleatòries sense reposició d'aquesta població, i ens quedarem amb la mida i el valor màxim de cada una d'elles, que és l'únic que necessitam saber. Les mides les generarem a l'atzar entre, posem, 25 i 75:

```
Mostra=function(a,b,P){  
  # a i b: mides màxima i mínima de la mostra; P: mida de la població  
  n=sample(a:b,1) # Mida de la mostra  
  X=sample(P,n,rep=FALSE) # Mostra aleatòria  
  c(n,max(X)) # Parell (mida, màxim)  
}  
Mostres=replicate(100,Mostra(25,75,N))  
Mostres
```

```

##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10] [,11] [,12] [,13]
## [1,]   43   53   53   40   67   44   40   73   67   67   61   26   37
## [2,] 6682 6684 6684 6652 6673 6623 6607 6540 6670 6657 6657 6455 6528
##      [,15] [,16] [,17] [,18] [,19] [,20] [,21] [,22] [,23] [,24] [,25] [,
## [1,]    75    29    56    29    33    74    51    64    42    61    39
## [2,] 6686 6651 6632 6617 6671 6470 6379 6565 6637 6684 6423 6
##      [,27] [,28] [,29] [,30] [,31] [,32] [,33] [,34] [,35] [,36] [,37] [,
## [1,]    48    28    57    28    63    46    32    31    44    69    31
## [2,] 6105 6522 6293 6327 6691 6498 6464 6522 6652 6606 6298 6
##      [,39] [,40] [,41] [,42] [,43] [,44] [,45] [,46] [,47] [,48] [,49] [,
## [1,]    27    65    70    30    42    56    53    40    33    25    61
## [2,] 6562 6676 6509 5655 6370 6641 6653 6674 6341 6612 6641 6
##      [,51] [,52] [,53] [,54] [,55] [,56] [,57] [,58] [,59] [,60] [,61] [,
## [1,]    59    42    60    52    64    38    26    71    75    48    39
## [2,] 6657 6596 6545 6681 6528 6583 6656 6692 6588 6603 6569 6
##      [,63] [,64] [,65] [,66] [,67] [,68] [,69] [,70] [,71] [,72] [,73] [,
## [1,]    41    37    57    70    51    29    69    46    52    64    56
## [2,] 6480 6535 6693 6694 6422 6665 6630 6654 6679 6414 6609 6
##      [,75] [,76] [,77] [,78] [,79] [,80] [,81] [,82] [,83] [,84] [,85] [,
## [1,]    59    57    56    69    39    70    58    59    48    50    37
## [2,] 6675 6624 6223 6618 6480 6664 6681 6536 6525 6671 6669 6
##      [,87] [,88] [,89] [,90] [,91] [,92] [,93] [,94] [,95] [,96] [,97] [,
## [1,]    70    44    39    47    56    65    43    32    60    57    42
## [2,] 6697 6490 6677 6619 6492 6609 6587 6653 6564 6674 6152 6
##      [,99] [,100]
## [1,]    75    36
## [2,] 6398 6623

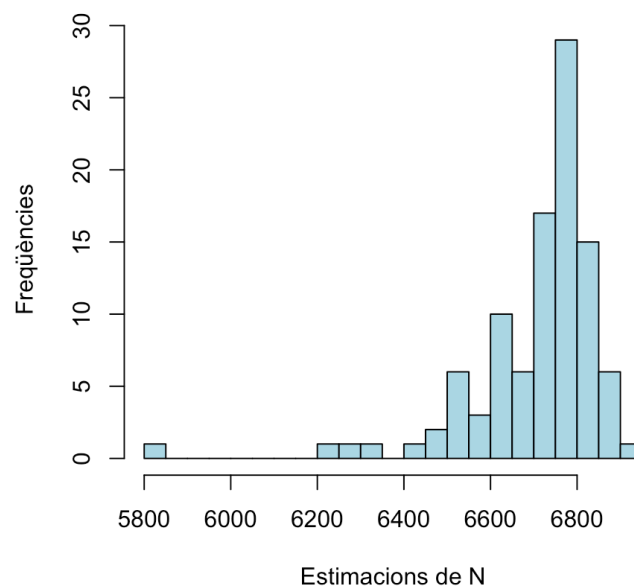
```

En aquesta matriu `Mostres`, cada columna correspon a una mostra aleatòria: la primera filera és la seva mida n i la segona filera el màxim m . Ara, amb cada una d'aquestes mostres, podem estimar la mida N de la població per mitjà de la fórmula $m + (m - n)/n$. Donarem aquestes estimacions ordenades de menor a major.

```
Estimacions=Mostres[2,]+(Mostres[2,]-Mostres[1,])/Mostres[1,]
round(sort(Estimacions),1)
```

```
## [1] 5842.5 6231.2 6297.5 6333.1 6402.4 6463.0 6482.3 6500.2 6503.1 6513
## [11] 6520.7 6532.2 6546.9 6552.0 6556.4 6586.7 6601.0 6606.9 6616.9 6628
## [21] 6629.0 6636.5 6637.0 6638.3 6645.2 6645.8 6653.1 6659.9 6665.0 6666
## [31] 6672.4 6674.8 6700.7 6702.3 6703.4 6703.7 6709.7 6710.6 6712.9 6722
## [41] 6725.1 6726.0 6731.4 6736.4 6739.2 6739.2 6739.6 6748.9 6749.4 6752
## [51] 6753.9 6755.2 6755.4 6758.2 6758.6 6758.8 6765.1 6768.6 6768.8 6771
## [61] 6771.6 6772.5 6774.1 6777.5 6777.7 6783.3 6785.3 6787.1 6788.6 6790
## [71] 6791.7 6792.6 6794.0 6794.6 6795.2 6796.2 6797.7 6798.2 6802.2 6803
## [81] 6804.0 6806.0 6806.4 6808.5 6809.1 6809.1 6809.4 6817.3 6836.4 6839
## [91] 6844.2 6847.2 6848.2 6859.9 6872.2 6875.5 6879.3 6888.2 6893.8 6911
```

```
hist(Estimacions, breaks=20,col="light blue",
     xlab="Estimacions de N",ylab="Freqüències",main="")
```



Com veieu, obtenim estimacions que van de 5842.5 a 6911. La mitjana d'aquestes estimacions és

```
round(mean(Estimacions),1)
```

```
## [1] 6704.5
```

És hora de descobrir el valor de N , per veure si hi hem fet a prop:

```
N
```

```
## [1] 6701
```

No hem fet molt enfora, com veieu.

3.7.2 Marca-recaptura

Suposem que en una població hi ha N individus, en capturam K (tots diferents), els marcam i els tornam a amollar. Al cap de poc temps, en capturam n , dels quals resulta que k estan marcats. A partir d'aquestes dades, volem estimar el valor de N .

Si suposam que N i K no han canviat de la primera a la segona captura (cap individu no ha abandonat la població ni se n'hi ha incorporat cap de nou), aleshores la variable aleatòria X definida per “Capturam un individu i miram si està marcat” és Bernoulli $Be(p)$ amb $p = K/N$, on coneixem la K i volem estimar la N .

Sigui ara x_1, \dots, x_n la mostra capturada en segon lloc. La seva proporció mostral d'individus marcats és $\hat{p}_X = k/n$. Com que \hat{p}_X és l'estimador màxim versemblant de p , estimam que

$$\frac{K}{N} = \frac{k}{n}$$

d'on, aïllant la N , estimam que

$$N = \frac{n \cdot K}{k}.$$

En resum, l'estimador

$$\widehat{N} = \frac{n \cdot K}{k}$$

maximitza la probabilitat d'obtenir k individus marcats en una mostra aleatòria de n individus. És l'**estimador màxim versemblant** de N a partir de K , k i n ; també se li diu **estimador de Lincoln-Petersen**. Fixau-vos que aquest estimador no fa res més que traduir la proporció "Si he trobat k individus marcats en un conjunt de n individus, què ha de valer el nombre total N de individus perquè hi hagi en total K individus marcats?"

Exemple 3.14 Suposem que hem marcat 15 peixos d'un llac, i que en una captura posterior de 10 peixos, n'hi ha 4 de marcats. Quants peixos conté el llac?

Ho estimarem amb l'estimador de Lincoln-Petersen:

$$\widehat{N} = \frac{15 \cdot 10}{4} = 37.5$$

Per tant, estimam que hi haurà entre 37 i 38 peixos al llac.

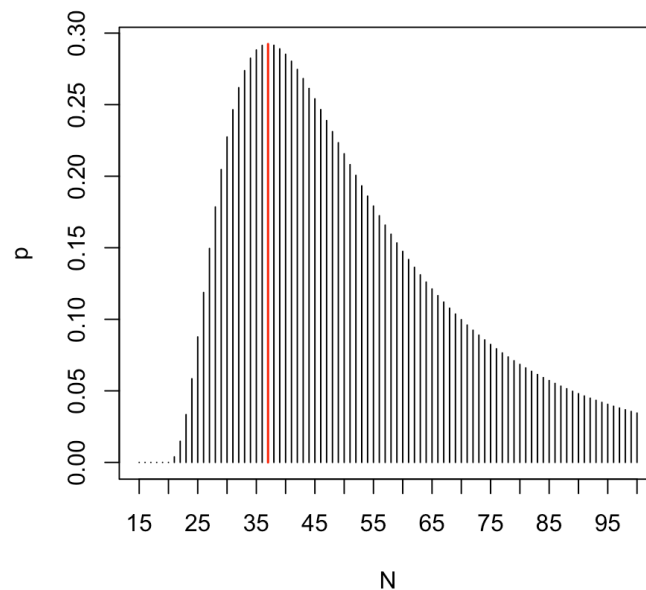
En aquest cas podem comprovar la màxima versemblança d'aquesta estimació, calculant la probabilitat d'obtenir 4 individus marcats en una mostra aleatòria de 10 individus d'una població de N individus on n'hi ha 15 de marcats i trobant la N que maximitza aquesta probabilitat. Per fer-ho, recordem que si una població està formada per K subjectes marcats i $N - K$ subjectes no marcats, el nombre de subjectes marcats en mostres aleatòries sense reposició de mida n segueix una distribució hipergeomètrica $H(K, N - K, n)$. Per tant, per a cada possible N , la probabilitat que en una mostra de 10 peixos del nostre llac n'hi hagi 4 de marcats serà $\text{dhyper}(4, 15, N - 15, 10)$.

```
N=15:1000 # Rang de possibles valors de N
p=dhyper(4, 15, N-15, 10) # Probabilitats de 4 marcats en 10
Nmax=N[which(p==max(p))] # N que maximitza la probabilitat
Nmax
```

```
## [1] 37
```

Aquest N_{\max} és la N que fa màxima la probabilitat que en una mostra de 10 peixos del nostre llac n'hi hagi 4 de marcats. Vegem-ho amb un gràfic:

```
plot(N[1:86],p[1:86],type="h",xaxp=c(15,100,17),xlab="N",ylab="p")
points(Nmax,dhyper(4,15,Nmax-15,10),type="h",col="red",lwd=1.5)
```



I què hagués passat si no haguéssim trobat cap peix marcat a la mostra?

Quan la mida de la mostra és petita, és més convenient emprar l'**estimador de Chapman**:

$$\widehat{N} = \frac{(n+1) \cdot (K+1)}{k+1} - 1$$

La idea és que afegim a la població un individu extra i marcat, que suposam que també capturam a la segona mostra. Llavors, aplicam l'estimador anterior i finalment restam 1, per descomptar l'individu marcat extra que realment no pertany a la població que volem estimar. D'aquesta manera ja no tenim el problema de dividir per 0 si k ens dóna 0.

En la situació de l'Exemple 3.14, aquest estimador dóna

$$\widehat{N} = \frac{16 \cdot 11}{5} - 1 = 34.2$$

i ens fa estimar una població total d'uns 34 peixos. Abans hem obtingut entre 37 i 38 peixos.

Quina de les dues estimacions s'acosta més a la realitat? Ni idea, no ho podem saber. Amb una altra recaptura segurament haguéssim obtingut resultats diferents.

L'estimador de Lincoln-Petersen

$$\widehat{N} = \frac{n \cdot K}{k}$$

és esbiaixat, amb biaix que tendeix a 0. L'estimador de Chapman és menys esbiaixat però no és màxim versemblant.

Exemple 3.15 Fem un experiment similar al de l'Exemple 3.13. Generarem a l'atzar una mida N d'una població grandeta i en marcarem una certa quantitat K . A continuació, prendrem 50 mostres aleatòries sense reposició de la nostra població i amb cada una d'aquestes mostres estimarem la N emprant els dos estimadors que hem explicat en aquesta subsecció. Al final, calcularem les mitjanes d'aquestes estimacions i les compararem amb el valor real de N , que no descobrirem fins el final. Com a l'Exemple 3.13, fixarem la llavor d'aleatorietat a l'atzar.

```
Llavor2=sample(1000,1)
Llavor2
```

```
## [1] 206
```

```
set.seed(Llavor2)
```

Ara generam la mida N de la població com un nombre a l'atzar entre 5000 i 10000.

```
N=sample(5000:10000,1)
```

Ara en capturam i marcam K; per fixar idees, prendrem $K=200$.

$K=200$

Per simplificar, suposarem que els N individus de la nostra població estan numerats de l'1 a l'N i que els marcats són els K primers. Ara generarem 100 mostres aleatòries sense reposició d'aquesta població, i ens quedarem amb la mida i el nombre d'individus marcats (és a dir, el nombre de valors menor o iguals a $K=200$ en la mostra). Les mides les generarem a l'atzar entre, posem, 100 i 150:

```
Mostra=function(a,b,P,M){  
  # a i b: mides màxima i mínima de la mostra; P: mida de la població;  
  # M: nombre de marcats  
  n=sample(a:b,1) # Mida de la mostra  
  X=sample(P,n,rep=FALSE) # Mostra aleatòria  
  c(n,length(which(X<=M))) # Parell (mida, nombre de marcats)  
}  
Mostres=replicate(100, Mostra(100,150,N,K))  
Mostres
```



```

##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10] [,11] [,12] [,13]
## [1,] 119 101 101 109 130 121 129 107 136 108 120 146 124
## [2,] 5 2 5 1 5 4 6 3 4 3 2 5 3
##      [,15] [,16] [,17] [,18] [,19] [,20] [,21] [,22] [,23] [,24] [,25] [,
## [1,] 106 131 120 149 113 143 117 125 132 101 109
## [2,] 1 3 4 5 2 3 4 5 1 4 3
##      [,27] [,28] [,29] [,30] [,31] [,32] [,33] [,34] [,35] [,36] [,37] [,
## [1,] 138 143 109 105 141 110 125 150 109 103 141
## [2,] 4 3 8 2 3 1 5 3 2 6 8
##      [,39] [,40] [,41] [,42] [,43] [,44] [,45] [,46] [,47] [,48] [,49] [,
## [1,] 107 132 135 118 131 117 110 115 135 138 127
## [2,] 2 4 3 4 3 5 1 3 6 2 6
##      [,51] [,52] [,53] [,54] [,55] [,56] [,57] [,58] [,59] [,60] [,61] [,
## [1,] 124 144 134 115 129 102 111 130 102 107 111
## [2,] 2 4 4 1 5 6 5 5 2 1 4
##      [,63] [,64] [,65] [,66] [,67] [,68] [,69] [,70] [,71] [,72] [,73] [,
## [1,] 119 131 115 150 145 125 119 149 129 112 146
## [2,] 4 3 2 1 7 1 3 1 5 4 7
##      [,75] [,76] [,77] [,78] [,79] [,80] [,81] [,82] [,83] [,84] [,85] [,
## [1,] 119 128 124 147 148 135 139 150 123 142 108
## [2,] 2 5 4 4 5 6 4 6 5 8 2
##      [,87] [,88] [,89] [,90] [,91] [,92] [,93] [,94] [,95] [,96] [,97] [,
## [1,] 112 141 131 144 111 112 147 103 112 108 107
## [2,] 2 7 5 4 1 6 3 4 4 1 2
##      [,99] [,100]
## [1,] 139 149
## [2,] 1 5

```

En aquesta matriu `Mostres`, cada columna correspon a una mostra aleatòria: la primera filera és la seva mida n i la segona filera el nombre d'individus marcats a la mostra. Ara, amb cada una d'aquestes mostres, podem estimar la mida N de la població per mitjà de l'estimador de Lincoln-Petersen.

```

EstimacionsLP=Mostres[1,]*K/Mostres[2,]
round(sort(EstimacionsLP),1)

```

```

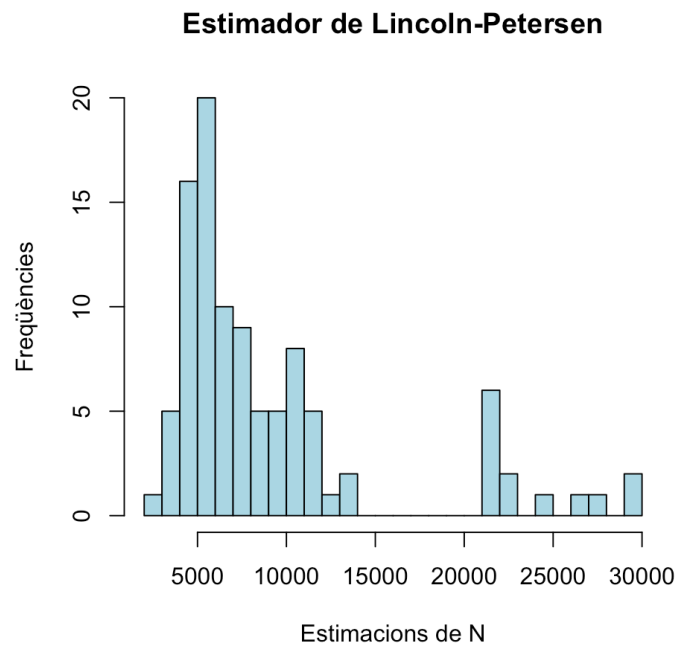
## [1] 2725.0 3400.0 3433.3 3525.0 3550.0 3733.3 4028.6 4040.0 41
## [10] 4171.4 4233.3 4300.0 4440.0 4500.0 4500.0 4680.0 4760.0 48
## [19] 4920.0 5000.0 5000.0 5000.0 5050.0 5120.0 5150.0 5160.0 51
## [28] 5200.0 5200.0 5240.0 5320.0 5550.0 5600.0 5600.0 5840.0 58
## [37] 5900.0 5920.0 5950.0 5960.0 5960.0 6000.0 6050.0 6150.0 62
## [46] 6600.0 6700.0 6700.0 6750.0 6800.0 6900.0 6950.0 7133.3 72
## [55] 7200.0 7200.0 7266.7 7350.0 7666.7 7666.7 7933.3 8266.7 87
## [64] 8733.3 8733.3 9000.0 9400.0 9533.3 9533.3 9800.0 10000.0 101
## [73] 10200.0 10500.0 10700.0 10700.0 10800.0 10800.0 10900.0 11200.0 113
## [82] 11500.0 11900.0 12000.0 12400.0 13800.0 13800.0 21200.0 21400.0 216
## [91] 21800.0 22000.0 22000.0 22200.0 23000.0 25000.0 26400.0 27800.0 298
## [100] 30000.0

```

```

hist(EstimacionsLP, breaks=20,col="light blue",xlab="Estimacions de N",ylab=

```



Com veieu, obtenim estimacions que van de 2725 a 30000. La mitjana de les estimacions és

```
round(mean(EstimacionsLP),1)
```

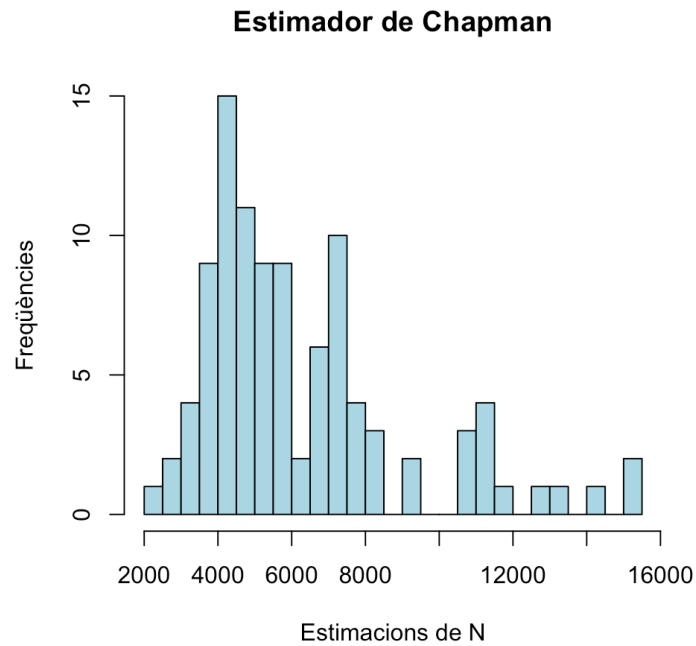
```
## [1] 9246
```

També podem emprar l'estimador de Chapman:

```
EstimacionsCh=(Mostres[,1]+1)*(K+1)/(Mostres[,2]+1)-1
round(sort(EstimacionsCh),1)
```

```
## [1] 2455.7 2956.6 2985.3 3170.3 3192.7 3243.7 3416.0 3566.8 36
## [10] 3674.4 3692.4 3731.9 3751.0 3904.1 3904.1 3952.0 4019.0 40
## [19] 4099.4 4153.0 4179.8 4220.0 4220.0 4320.5 4334.9 4354.0 43
## [28] 4387.5 4387.5 4421.0 4488.0 4501.4 4541.6 4541.6 4742.6 47
## [37] 4823.0 4863.2 4903.4 4923.5 4983.8 4990.5 5024.0 5024.0 50
## [46] 5345.6 5426.0 5426.0 5426.0 5466.2 5476.2 5506.4 5526.5 55
## [55] 5627.0 5828.0 5828.0 5828.0 5828.0 5948.6 6029.0 6280.2 66
## [64] 6632.0 6632.0 6833.0 6833.0 6900.0 7101.0 7134.5 7235.0 72
## [73] 7235.0 7235.0 7302.0 7302.0 7369.0 7436.0 7570.0 7586.8 76
## [82] 7771.0 8039.0 8106.0 8374.0 9312.0 9312.0 10752.5 10853.0 109
## [91] 11054.0 11154.5 11154.5 11255.0 11657.0 12662.0 13365.5 14069.0 150
## [100] 15174.5
```

```
hist(EstimacionsCh, breaks=20,col="light blue",xlab="Estimacions de N",ylab=
```



Com veieu, obtenim estimacions que van de 2455.7 a 15174. La mitjana d'aquestes estimacions és

```
round(mean(EstimacionsCh),1)
```

```
## [1] 6292.7
```

És hora de descobrir el valor de N , per veure si hi hem fet a prop:

```
N
```

```
## [1] 7225
```

Com veieu, amb l'estimador de Chapman ens hem fet més a prop del valor real de N que amb el màxim versemblant.

3.8 Test de la lliçó 3

(1) Tenim una variable aleatòria X normal de mitjana μ i desviació típica σ . Prenem mostres aleatòries simples de mida n , i indicam amb \tilde{S}_X la seva desviació típica mostral. Quina o quines de les afirmacions següents són vertaderes?

1. $E(\tilde{S}_X^2) = \sigma^2$.
2. $E(\tilde{S}_X) = \sigma$.
3. \tilde{S}_X^2 segueix una distribució χ^2 amb $n - 1$ graus de llibertat.
4. $(n - 1)\tilde{S}_X^2/\sigma^2$ segueix una distribució χ^2 amb $n - 1$ graus de llibertat.
5. Totes les altres respostes són falses.

(2) Quina o quines de les afirmacions següents sobre la mitjana mostral són vertaderes?

1. Si la distribució poblacional és normal, sempre coincideix amb la mediana de la mostra.
2. Sempre serveix per estimar la mitjana poblacional.
3. Sempre serveix per estimar la mediana poblacional.
4. Si la distribució poblacional és normal, serveix per estimar la mediana poblacional.
5. Cap de les altres respostes és correcta.

(3) L'error estàndard de la mitjana mostral de mida $n \geq 2$ (marca totes les continuacions correctes):

1. Mesura la variabilitat de les observacions que formen la mostra.
2. És l'exactitud amb què es mesura cada observació de la mostra.
3. Mesura la variabilitat de les mitjanes mostrals de mostres aleatòries simples.
4. Mesura la precisió amb què les mitjanes mostrals de mostres aleatòries simples estimen la mitjana poblacional.
5. És proporcional a la mitjana mostral.
6. És més gran que la desviació típica de la població.
7. Sobre cada mostra val la desviació típica de la mostra.
8. Cap de les altres respostes és correcta.

(4) La proporció d'afectats per una determinada malaltia en una població és del 10%. Si estimam aquesta proporció poblacional repetidament a partir de mostres de mida 1000, aquestes estimacions segueixen una distribució que (marca totes les afirmacions correctes):

1. És una distribució mostral.
2. És aproximadament normal.
3. Té mitjana 0.1.
4. Té variància 90.
5. És binomial.
6. Cap de les altres respostes és correcta.

(5) Què hem de fer per disminuir a la meitat l'error estàndard d'una proporció?

1. Hem d'augmentar en un 50% la mida de la mostra
2. Hem de doblar la mida de la mostra.
3. Hem de quadruplicar la mida de la mostra.
4. Hem de dividir per 2 la mida de la mostra.
5. Hem de dividir per 4 la mida de la mostra.
6. Cap de les altres respostes és correcta.

(6) La probabilitat que els individus d'una determinada població tenguin una determinada característica C és p . Prenem mostres aleatòries simples de mida n d'aquesta població, i indicam amb \hat{p}_X la seva proporció mostral. Quina o quines de les afirmacions següents són vertaderes?

1. \hat{p}_X té sempre distribució binomial $B(n, p)$.
2. \hat{p}_X té sempre distribució normal.
3. Si n és gran, \hat{p}_X té distribució aproximadament binomial $B(n, p)$.
4. Si n és gran, \hat{p}_X té distribució aproximadament normal.
5. L'error estàndard de \hat{p}_X és $\sqrt{p(1-p)/n}$.

(7) Sigui X una variable aleatòria que no és constant. Si en prenem mostres aleatòries simples més grans (marca totes les continuacions correctes):

1. La mitjana mostral sempre disminueix.
2. L'error estàndard de la mitjana sempre disminueix.
3. L'error estàndard de la mitjana sempre augmenta.
4. La variància mostral sempre augmenta.
5. El nombre de graus de llibertat de l'estimador χ^2 associat a la variància mostral sempre augmenta.

6. Cap de les altres respostes és correcta.

(8) La longitud d'una determinada espècie d'animals té un valor mitjà de μ cm. Si prenem mostres aleatòries simples de 20 exemplars, calculam la seva mitjana mostral \bar{X} i la seva desviació típica mostral \tilde{S}_X (marca la continuació més correcta):

1. L'estadístic $\frac{\bar{X}-\mu}{\tilde{S}_X/\sqrt{20}}$ segueix sempre una llei normal.
2. L'estadístic $\frac{\bar{X}-\mu}{\tilde{S}_X/\sqrt{20}}$ segueix sempre una llei t de Student.
3. L'estadístic $\frac{\bar{X}-\mu}{\tilde{S}_X/\sqrt{20}}$ segueix una llei normal si la longitud segueix una llei normal.
4. L'estadístic $\frac{\bar{X}-\mu}{\tilde{S}_X/\sqrt{20}}$ segueix una llei t de Student si la longitud segueix una llei normal.
5. L'estadístic $\frac{\bar{X}-\mu}{\tilde{S}_X/\sqrt{20}}$ no segueix mai ni una llei normal ni una llei t de Student, perquè les mostres no són prou grans.

(9) En una mostra de 100 dones es va obtenir una concentració mitjana d'hemoglobina en sang de 10 amb una desviació típica de 2. Quin és l'error típic de la mostra?

1. 0.02
2. 0.04
3. 0.2
4. 0.4
5. 1
6. Cap dels valors anteriors

(10) Què significa que un estimador d'un paràmetre d'una variable aleatòria sigui no esbiaixat?

1. Que la distribució mostral de l'estimador és normal.
2. Que aplicat a una mostra aleatòria simple sempre dona el valor poblacional del paràmetre.
3. Que el seu valor esperat és igual al valor poblacional del paràmetre.
4. Que aplicat a una mostra aleatòria simple sempre dona el valor esperat del paràmetre.
5. Que el seu error típic és petit.

(11) La concentració en sang a les persones d'un determinat metabolit (en mg/ml) té una distribució $N(23, 3)$. Quina de les afirmacions següents és vertadera?

1. Aproximadament un 90% de les mostres aleatòries de 100 individus tendran la seva mitjana entre 22.4 i 23.6 mg/ml.
2. Aproximadament un 95% de les mostres aleatòries de 100 individus tendran la seva mitjana entre 22.4 i 23.6 mg/ml.
3. Aproximadament un 99% de les mostres aleatòries de 100 individus tendran la seva mitjana entre 22.4 i 23.6 mg/ml.
4. Més d'un 99% de les mostres aleatòries de 100 individus tendran mitjana igual a 23.
5. Cap de les afirmacions anteriors és vertadera.

(12) Sigui X una variable aleatòria $N(\mu_X, 2)$ i sigui \bar{X} la mitjana mostral de mida 10 de X . Quina de les afirmacions següents és vertadera?

1. La desviació típica de \bar{X} és igual a 2.
2. La desviació típica de \bar{X} és menor que 2.
3. La desviació típica de \bar{X} és major que 2.
4. Que la desviació típica de \bar{X} sigui major, menor o igual que 2 depèn de μ_X .
5. Cap de les afirmacions anteriors és vertadera.