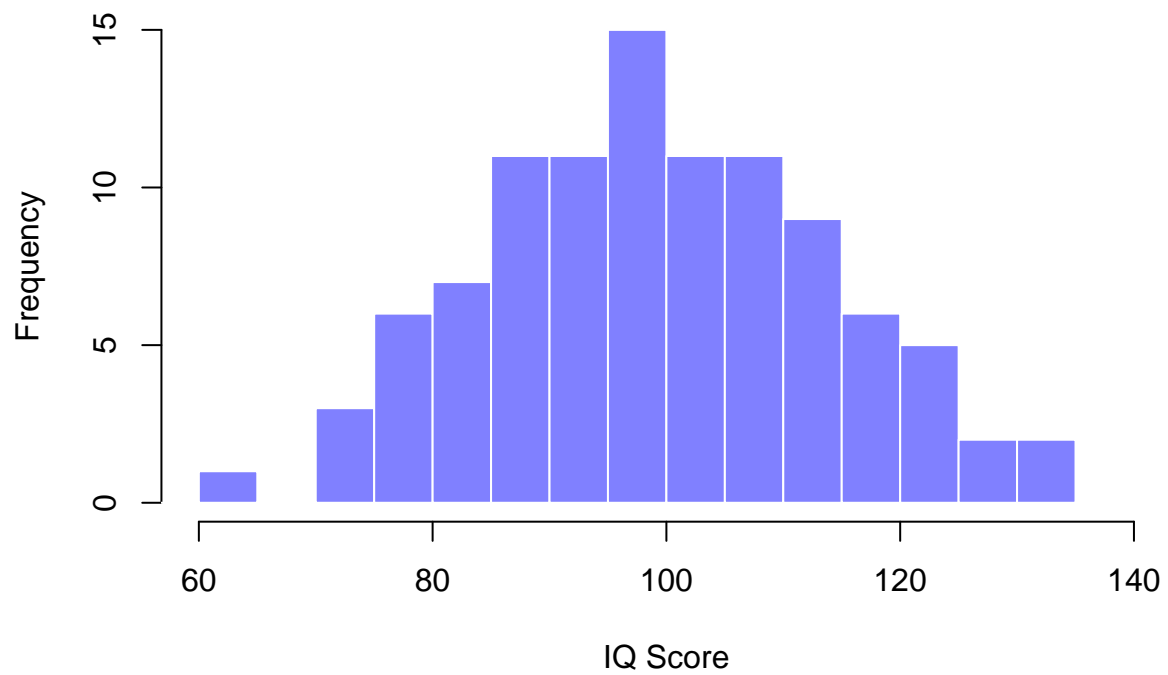
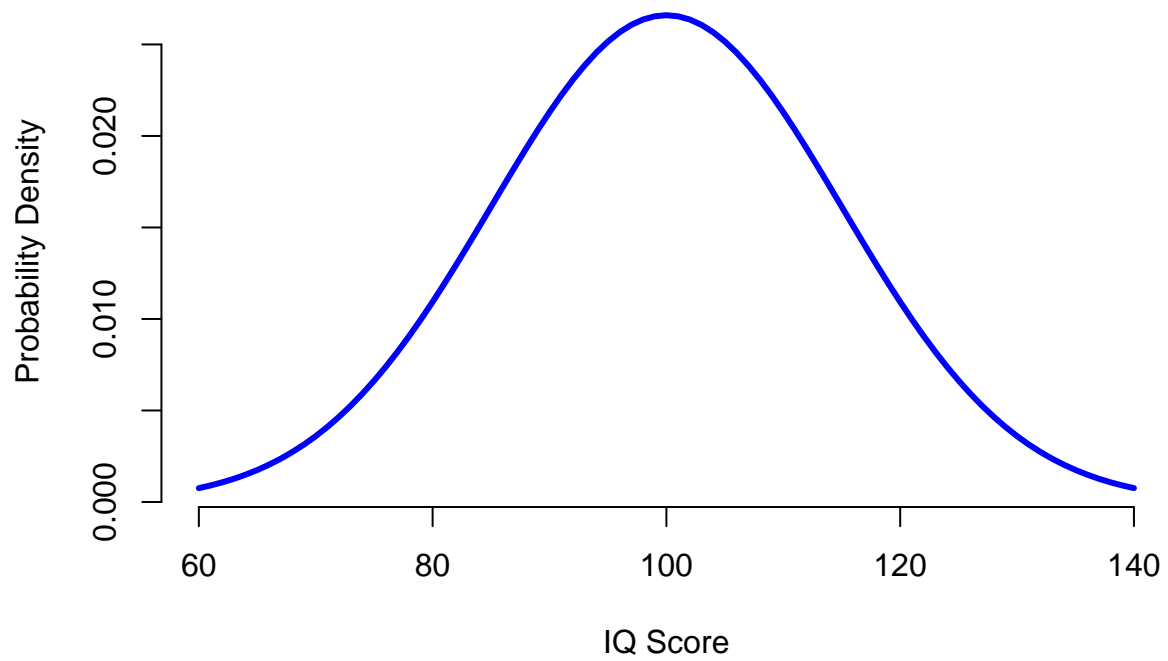


Prepa3

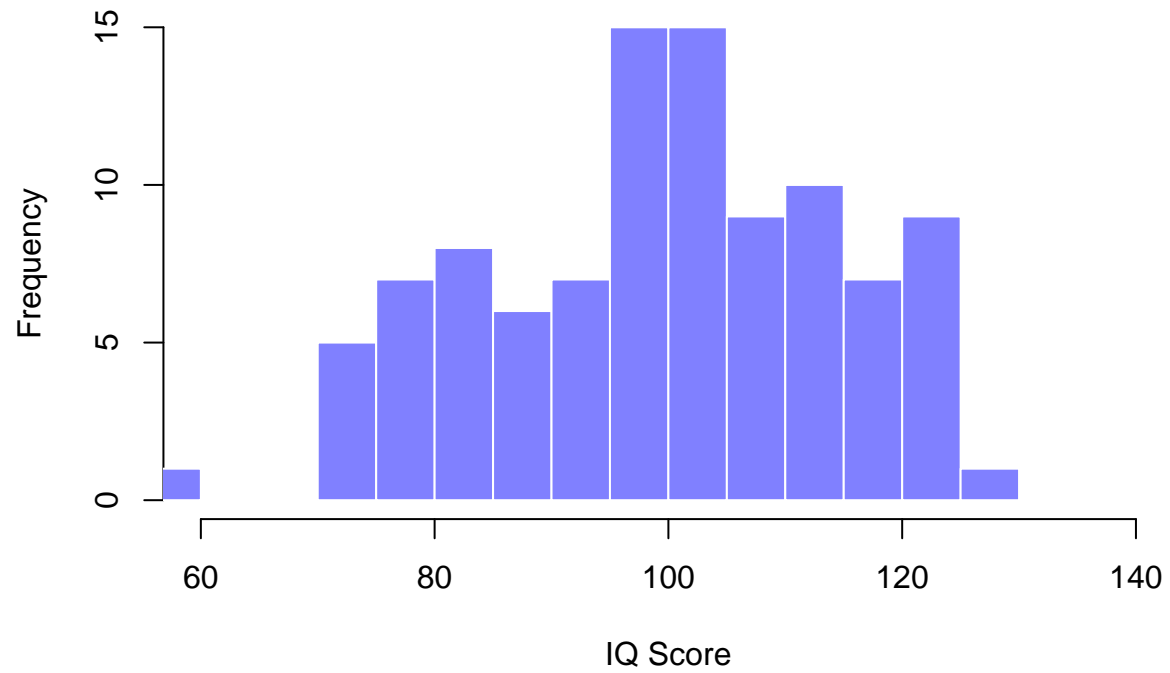
Alimi Garmendia

6/4/2021

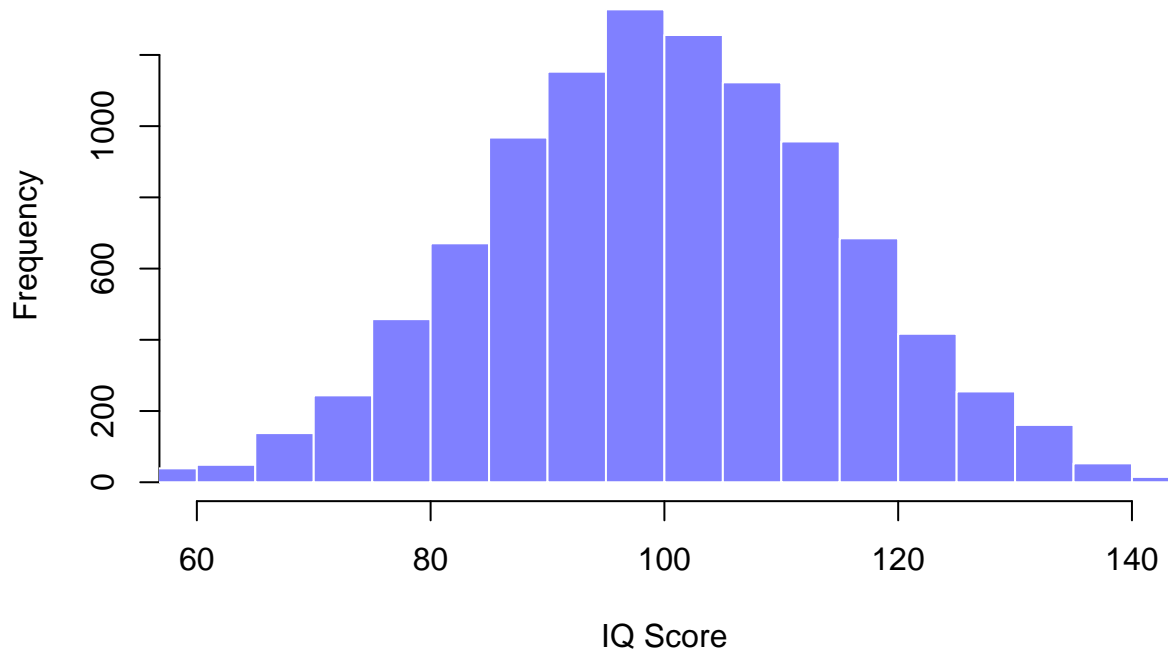
Muestreo, realmente importa?



```
## [1] "n= 100 mean= 99.2874785062804 sd= 14.6595766086536"
```



```
## [1] "n= 100 mean= 99.3400026658619 sd= 15.0275890703367"
```



```
## [1] "n= 10000 mean= 99.9208557174926 sd= 15.1281200571434"
```

No siempre vamos a tener acceso a muchas muestras de la poblacion, es por ello que debemos muestrear de la mejor manera posible, que las muestras, en conjunto, representen de mejor manera la poblacion.

Una vez tomada nuestra muestra debemos *estimar* las medidas estadisticas, una vez mas, escogiendo los aquellos estimadores que tengan menor sesgo o *bias*

Estimadores

Para estimar la media de la poblacion, podemos usar el mismo procedimiento que usamos para calcular la media de la muestra, es decir

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N x_i$$

Asi \bar{X} es un estimador insesgado de la media poblacional μ .

simbolo	que es?	Sabemos que es?
\bar{X}	Media de la muestra	Si, lo calculamos de los datos
μ	Media real de la poblacion	Casi nunca sabremos su valor con certeza
$\bar{\mu}$	Estimado de la media de la poblacion	Si, es identico a la media de la muestra

Para estimar la desviacion estandar de la muestra tenemos que realizar un pequeños ajuste de lo que conocemos

Habiamos definido la varianza de una muestra como

$$S^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{X})^2$$

Sin embargo, no podemos asumir este valor como estimador de la varianza de la poblacion pues este estimador es un estimador cegado. Se puede demostrar que, para obtener un estimador incesgado de la varianza poblacional podemos usar como estimador

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{X})^2$$

Luego podemos estimar la desviacion estandar usando

$$\hat{\sigma} = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{X})^2}$$

simbolo	que es?	Sabemos que es?
S^2	Varianza de la muestra	Si, la calculamos de los datos
σ^2	Varianza de la poblacion	Casi nunca sabremos su valor real con certeza
$\hat{\sigma}^2$	Estimador de la varianza de la Poblacion	Si, es casi igual que la varianza de la muestra pero con un ligero ajuste

Como hemos visto, no es difícil calcular los estimadores de la población, sin embargo, como podemos corroborar, estos estimadores dependen de la muestra que tome. Si tomo muchas muestras y calculo varias veces estos estimadores, veremos que varían de muestra en muestra. Es por eso que debemos buscar la manera de asegurar el rango en que nuestros estimadores se encontraran, con algún grado de certeza.

Intervalos de confianza

Suponiendo que nuestras muestras tienen una distribución aproximadamente normal. Que la media de la población es μ y la SD σ . Imaginemos que he hecho un estudio que incluye N y que la media del IQ de estos participantes es \bar{X} . Recordemos que, en una distribución normal, hay una prob de 95% que una cantidad con distribución normal se encuentre dentro de dos desviaciones* estándar de la media. Es decir:

```
qnorm(p = c(0.025,0.975))
```

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

Es decir, nuestra media muestral cumple que:

$$\mu - 1.96 \frac{\sigma}{\sqrt{n}} \leq \bar{X} \leq \mu + 1.96 \frac{\sigma}{\sqrt{n}}$$

Así, trabajando un poco el álgebra, podemos llegar a

$$\bar{X} - 1.96 \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + 1.96 \frac{\sigma}{\sqrt{n}}$$

Así, este rango de valores tiene 95% de probabilidad de contener la media de la población.

Veamos que este procedimiento depende de que conozcamos la varianza de la **población**, cosa que no siempre conoceremos. Por lo que debemos usar un estimador tomado de la **muestra**, $\hat{\sigma}$. Esto es directo de hacer sin embargo esto implica que debemos hacer una corrección. En lugar de usar la distribución normal, usaremos la distribución t, que depende del número de observaciones N . Cuando N es suficientemente grande, los valores obtenidos en la distribución t y en la normal serán iguales.

```
qt(p = .975, df = 10000)
```

```
## [1] 1.960201
```

Ejercicios

1. Imaginemos que quieres rentar un apartamento de una habitación en Caracas. La media de la renta mensual para una muestra aleatoria de 60 apartamentos que viste en el periódico es de \$1000. Asuma que la varianza de la población es de \$200. Construya un intervalo de confianza del 95%

```
cuantil = qnorm(c(0.025,0.975))
sqrt_n = sqrt(60)
mean_muestra = 1000
varianza = 200

IC = 1000+cuantil*varianza/sqrt_n
IC
```

```
## [1] 949.3939 1050.6061
```

2. Sobre cual poblacion de los apartamentos de Caracas podemos inferir, dados los resultados del apartado anterior?

Podemos usar el resultado anterior para estimar la media de los apartamentos de que vismo en el periodico. No podemos asegurar nada de todos los apartamentos de Caracas, pues no podemos asegurar que la muestra del periodico sea representativa

3. Que tan grande debe ser la muestra de los apartamentos de los apartados anteriores si queremos estiumar la media de la poblacion dentro de un margen de \$50 con una confianza del 90%?

```
alpha = 1-0.9
cuantil = qnorm(c(1-alpha/2))
cuantil
```

```
## [1] 1.644854
```

$$50 = Z * \frac{\sigma}{\sqrt{n}}$$

```
n = (cuantil*varianza/50)**2
n
```

```
## [1] 43.2887
```

Es decir, para tener un margen de error de +- \$50 necesitaríamos una muestra aleatoria de 44 apartamentos

4. En una amplia muestra de aceite de oliva italianos se observan, entre otras cosas, la variable “porcentaje de acido oleico en la fracciion lipidica”. Consideraremos dos de las areas de procedencia de los aceites: *Liguria Occidental* y *Liguria Oriental*



Se analizaron 50 muestras procedentes de Liguria Occidental y otras 50 de Liguria Oriental obteniendo los siguientes resultados:

Region	\bar{x}	s^2	s
Liguria Occidental	76.742	1.356	1.164
Liguria Oriental	77.46	2.47	1.572

Suponiendo normalidad y homocedasticidad, a un nivel de significación del 1 %, calcular el intervalo de confianza para la diferencia de la cantidad media de ácido oleico en los aceites de ambas regiones. Especificar las suposiciones previas para que el procedimiento empleado sea válido

Lo primero es elegir el intervalo adecuado para la **diferencia de medias**. Podemos suponer que las muestras en ambas regiones son independientes entre sí. En base a las hipótesis dadas y al hecho de que no conocemos las varianzas, el intervalo será el siguiente:

$$IC_{1-\alpha}(\mu_1 - \mu_2) = (\bar{x}_1 - \bar{x}_2 \pm t_{n_1+n_2-2; \frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}})$$

con

$$s_p = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}}$$

Asi

```
x_1_var = 77.46
x_2_var = 76.742
s_1 = 2.47
s_2 = 1.356
n = 50
sp = sqrt( ((n-1)*s_1 + (n-1)*s_2)/(2*n-2) )
alpha = 0.01
sumy = sqrt(2/n)

t = qt(c(alpha/2,1-alpha/2), 98)
```

Ahora podemos calcular el IC

```
IC = x_1_var - x_2_var + t*sp*sumy
IC
```

```
## [1] -0.008668211 1.444668211
```

```
#IC2 = x_1_var - x_2_var + t*sqrt(s_1/n + s_2/n)
```

5. A partir de los resultados anteriores, es razonable suponer que la cantidad mediada de ácido oleico es diferente en Liguria Oriental y Liguria Occidental? Y con una confianza del 90 %?

Como el cero se encuentra dentro del intervalo, con los datos que tenemos, y con dicho nivel de confianza. No podemos asegurar de que haya diferencias en la cantidad de ácido oleico entre ambas regiones.

Si cambiáramos el nivel de confianza a 90% lo único que cambia es el percentil que cambia a 1.66, lo que dejaría al cero fuera del intervalo. Por lo tanto, con una confianza del 90% si podríamos decir que los niveles del ácido son distintos entre las dos regiones.

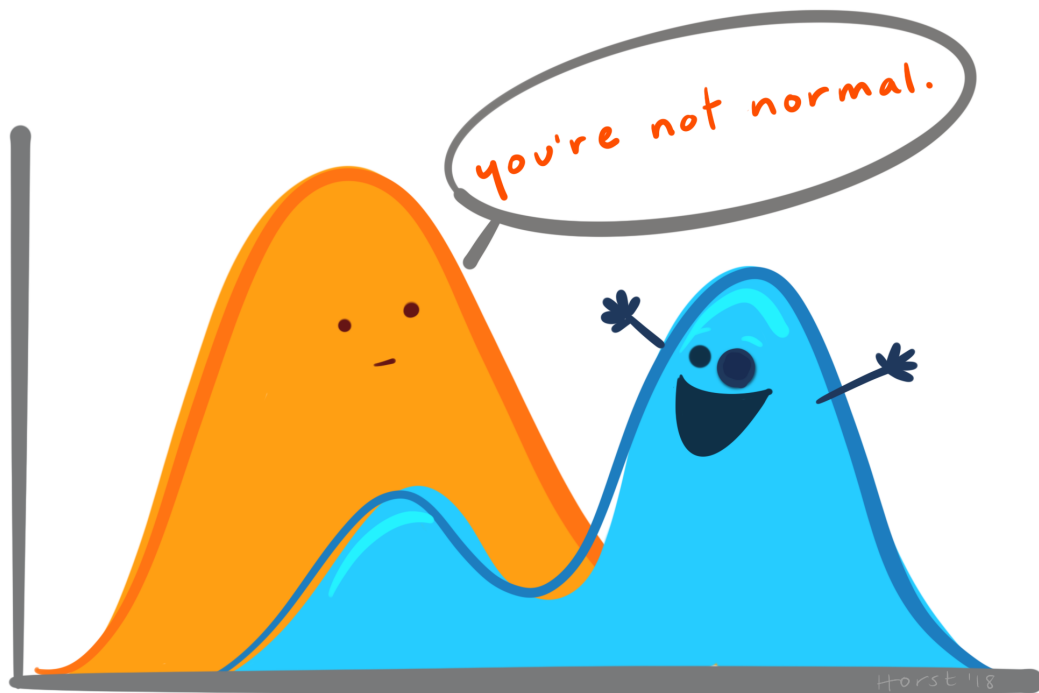
6. Calcular el intervalo de confianza al 95 % para la varianza de la cantidad de ácido oleico en Liguria Oriental. Especificar las suposiciones previas para que el procedimiento empleado sea válido.

Asumiendo normalidad, el intervalo de confianza a calcular es el siguiente

$$IC_{1-\alpha}(\sigma^2) = \left(\frac{(n-1)s^2}{\chi_{n-1;\alpha/2}^2}, \frac{(n-1)s^2}{\chi_{n-1;1-\alpha/2}^2} \right)$$

```
s = 2.47
n = 50
alpha = 0.05
sup = qchisq(alpha/2, n-1)
inf = qchisq(1-alpha/2, n-1)

numerador = (n-1)*s
IC = c(numerador/inf, numerador/sup)
```



Referencias

learn r from Danielle Navarro

Illustrations by Allison Horst