

Tema 4: INFERENCIA ESTADISTICA - V

El tamaño de la muestra: relevancia científica frente a significatividad estadística.
Contraste de hipótesis sobre la varianza y la proporción

Biología sanitaria 2017/18. Universidad de Alcalá

M. Marvá. Actualizado: 2017-11-16

El tamaño de la muestra: relevancia estadística vs relevancia científica

Un fabricante garantiza que produce comprimidos de diámetro medio de 13mm. En una muestra de 50 unidades tiene un diámetro medio $\bar{X} = 13.2$ mm, cuasidesviación típica $s = 0.6$ mm. ¿Es aceptable esta afirmación al nivel de significación $\alpha = 1\%$?

Se contrasta $H_0 : \mu = 13$, $H_1 : \mu \neq 13$

Región de no rechazo: $\left(13 - z_{0.995} \frac{0.6}{\sqrt{50}}, 13 + z_{0.995} \frac{0.6}{\sqrt{50}} \right)$

```
signif(13 + c(-1, 1) * qnorm(0.995) * 0.6/sqrt(50), digits = 4)
```

```
## [1] 12.78 13.22
```

El contraste NO es significativo.

¿Y si la muestra tiene tamaño 1000? Región de no rechazo:

```
signif(13 + c(-1, 1) * qnorm(0.995) * 0.6/sqrt(1000), digits = 4)
```

```
## [1] 12.95 13.05
```

y ahora el **contraste es significativo!!**

Relevancia estadística vs relevancia científica

- Todo contraste no significativo puede ser significativo si se aumenta el tamaño de la muestra lo suficientemente
- Cuando sólo se conocen los valores de los estadísticos (\bar{X}_* , s , n) la **d de Cohen** compara la diferencia entre la media muestral observada y H_0 con la cuasidesviación típica:

$$d = \frac{\bar{X}_* - \mu_0}{s}$$

sin considerar el tamaño de la muestra. Interpretación:

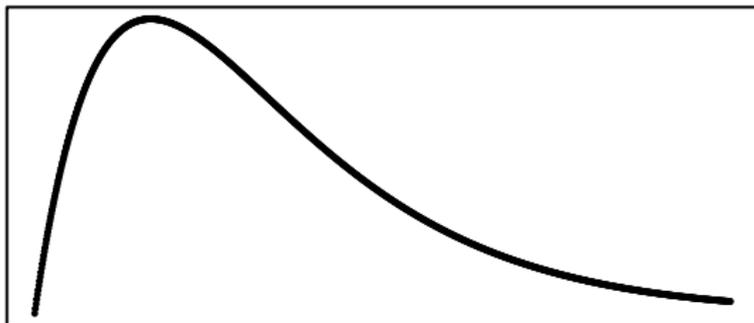
- ▶ $d < 0.2$, apunta a diferencia no relevante
 - ▶ $d > 0.8$, apunta a diferencia relevante
 - ▶ $0.2 < d < 0.8$ **conocimiento experto!!**
- Cuando se tiene acceso a los datos en bruto hay alternativas (que veremos en prácticas)

Sobre el ejemplo anterior $d = \frac{13.2 - 13}{0.6} = \frac{1}{3}$

IMPORTANTE: la d de Cohen **no sustituye** al contraste sobre la media. Es una herramienta *adicional* para el caso en el que se sospecha que el contraste es significativo debido a que el tamaño de la muestra es muy grande

Contraste sobre la varianza en poblaciones normales

El estadístico adecuado para entender la distribución muestral de σ^2 es $(n-1) \frac{s^2}{\sigma_0^2} \sim \chi_{n-1}^2$



El esquema del contraste: establecer H_0 ($\sigma = \sigma_0$) y H_1 ; obtener cuasidesviación típica s_* ,

- **Región de rechazo:** fijar ns α , usar localizar la región con valores del estadístico más contradictorios con H_0 , ubicar el valor del estadístico de contraste s_*^2
- **P-valor** Calcular la probabilidad de obtener valores s^2 del estadístico de contraste más contradictorios con H_0 que s_*^2

Ejemplo: Un laboratorio farmacéutico garantiza que produce comprimidos de diámetro uniforme, porque la desviación típica de su diámetro es 0.5mm. Una muestra de 15 unidades dio una cuasidesviación típica $s_* = 0.7\text{mm}$. ¿Es aceptable la afirmación del laboratorio al nivel de significación del 5%?

Para cada uno de los contrastes:

$$\begin{array}{ll} \textcircled{1} H_0 : \sigma_X^2 \leq \sigma_0^2 = 0.5 & H_1 : \sigma_X^2 = \sigma_0^2 > 0.5 \\ \textcircled{2} H_0 : \sigma_X = 0.5 = \sigma_0 & H_1 : \sigma_X \neq 0.5 = \sigma_0 \end{array}$$

se pide

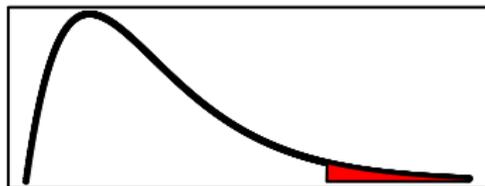
- 1 Dibujar, de forma aproximada, la región de rechazo
- 2 "localizar" el estadístico de contraste.

Contraste sobre la varianza con nivel de significación α .

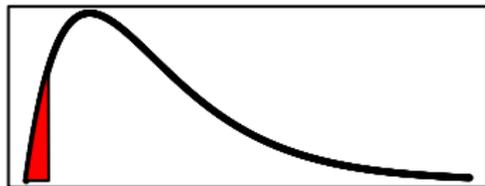
Región de rechazo: se usa el estadístico

$$Y = (n - 1) \frac{s^2}{\sigma_0^2}$$

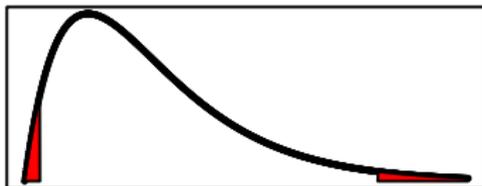
(a) $H_0 : \sigma^2 \leq \sigma_0^2$ $H_1 : \sigma^2 > \sigma_0^2$. Región de rechazo: $Y > \chi_{n-1, \alpha}^2$



(b) $H_0 : \sigma^2 \geq \sigma_0^2$ $H_1 : \sigma^2 < \sigma_0^2$. Región de rechazo: $Y < \chi_{n-1, 1-\alpha}^2$



(c) $H_0 : \sigma^2 = \sigma_0^2$ $H_1 : \sigma^2 \neq \sigma_0^2$. Región de rechazo: $Y \notin (\chi_{k, 1-\alpha/2}^2, \chi_{k, \alpha/2}^2)$



Ejemplo: Un laboratorio farmacéutico garantiza que produce comprimidos de diámetro uniforme, porque la desviación típica de su diámetro es 0.5mm. Una muestra de 15 unidades dio una cuasidesviación típica $s_* = 0.7$ mm. ¿Es aceptable la afirmación del laboratorio al nivel de significación del 5%?

Para cada uno de los contrastes:

$$\textcircled{1} H_0 : \sigma_X^2 \leq \sigma_0^2 = 0.5 \qquad H_1 : \sigma_X^2 = \sigma_0^2 > 0.5$$

$$\textcircled{2} H_0 : \sigma_X = 0.5 = \sigma_0 \qquad H_1 : \sigma_X \neq 0.5 = \sigma_0$$

se pide calcular el p-valor

Contraste sobre la varianza p-valor

Dados $H_0 (\sigma^2 = \sigma_0^2)$ y la varianza muestral s_*^2

(a) $H_0 : \sigma^2 \leq \sigma_0^2 \quad H_1 : \sigma^2 > \sigma_0^2.$

$$\text{P-valor} = P \left(\chi_{n-1}^2 > (n-1) \frac{s_*^2}{\sigma_0^2} \right),$$

cola derecha del estimador

(b) $H_0 : \sigma^2 \geq \sigma_0^2 \quad H_1 : \sigma^2 < \sigma_0^2.$

$$\text{P-valor} = P \left(\chi_{n-1}^2 < (n-1) \frac{s_*^2}{\sigma_0^2} \right),$$

cola izquierda del estimador.

(c) $H_0 : \sigma^2 = \sigma_0^2 \quad H_1 : \sigma^2 \neq \sigma_0^2.$

1 Si $s_*^2 > \sigma_0^2$, entonces

$$\text{P-valor} = 2P \left(\chi_{n-1}^2 > (n-1) \frac{s_*^2}{\sigma_0^2} \right)$$

2 Si $s_*^2 < \sigma_0^2$, entonces

$$\text{P-valor} = 2P \left(\chi_{n-1}^2 < (n-1) \frac{s_*^2}{\sigma_0^2} \right)$$

Inferencia sobre la proporción

Nos interesa la proporción p de individuos que tiene cierta característica.

Consideramos la **proporción muestral**

$$\hat{p} = \frac{X_1 + \dots + X_n}{n}, \quad X_i \sim \text{Bernoulli}(p)$$

Si se cumplen a la vez las condiciones:

$$n > 30, \quad n \cdot \hat{p} > 5, \quad n \cdot \hat{q} > 5.$$

Entonces, conforme crece n , la distribución de \hat{p} se aproxima a una normal

$$\hat{p} \sim N \left(p, \sqrt{\frac{\hat{p} \cdot \hat{q}}{n}} \right)$$

Ejemplo: Un estudio reciente sobre la mutación asociada con la brida en el arao común afirma que la proporción es $p = 0.35$. Un muestreo realizado en 2010 en la isla de Thresnish (Escocia) arrojó 139 individuos embridados y 317 no embridados. ¿Contradicen esta muestra la conclusión del estudio arriba citado?

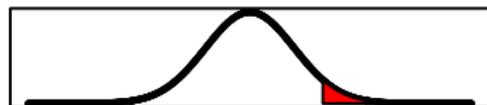
Realiza el contraste con al región de rechazo (con un nivel de significación del 1%) y calcula el p-valor del contraste.

Inferencia sobre la proporción Cálculo de la región de rechazo al ns α .

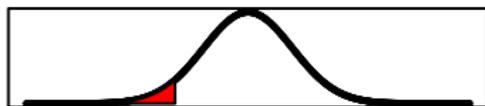
Fijada la hipótesis nula $H_0 (p = p_0)$ el TLC dice que

$$\hat{p} \sim N \left(p_0, \sqrt{\frac{\hat{p} \cdot \hat{q}}{n}} \right)$$

(a) $H_0 : p \leq p_0$ $H_a : p > p_0$. Región de rechazo: $\hat{p} > p_0 + z_\alpha \sqrt{\frac{p_0 \cdot q_0}{n}}$



(b) $H_0 : p \geq p_0$ $H_a : p < p_0$. Región de rechazo: $\hat{p} < p_0 - z_\alpha \sqrt{\frac{p_0 \cdot q_0}{n}}$



(c) $H_0 : p = p_0$ $H_a : p \neq p_0$. Región de rechazo: $|\hat{p} - p_0| > z_{\alpha/2} \sqrt{\frac{p_0 \cdot q_0}{n}}$



¡Ojo! se usa p_0 y q_0 en lugar de \hat{p} y \hat{q} para la desviación típica muestral.

Inferencia sobre la proporción Cálculo del p-valor. Fijada la hipótesis nula H_0 ($p = p_0$) y obtenida la proporción muestral \hat{p}_* el TLC dice así

$$\hat{p} \sim N\left(p_0, \sqrt{\frac{p_0 \cdot q_0}{n}}\right)$$

(a) $H_0 : p \leq p_0$ $H_a : p > p_0$. P-valor

$$P\left(Z > \frac{\hat{p} - p_0}{\sqrt{p_0 \cdot q_0/n}}\right) = P\left(\hat{p} > \hat{p}_* \mid \hat{p} \sim N\left(p_0, \sqrt{p_0 \cdot q_0/n}\right)\right)$$

cola a la derecha del estadístico. La igualdad es análoga para el resto de p-valores

(b) $H_0 : p \geq p_0$ $H_a : p < p_0$. P-valor

$$P\left(Z < \frac{\hat{p} - p_0}{\sqrt{p_0 \cdot q_0/n}}\right)$$

cola a la izquierda del estadístico.

(c) $H_0 : p = p_0$ $H_a : p \neq p_0$. P-valor

$$2 \cdot P\left(Z > \frac{|\hat{p} - p_0|}{\sqrt{p_0 \cdot q_0/n}}\right)$$

colas (bilateral) más alejadas de p_0 que el estadístico.