

“Lies, damned lies, and statistics”

“Hay tres tipos de mentiras:

- ▶ mentiras,
- ▶ grandes mentiras y
- ▶ estadísticas.”

Frase popularizada por Mark Twain (Samuel Langhorne Clemens, escritor norteamericano 1835-1919)

Probabilidad - Clase 11

Repaso de Estadística

Intervalos de confianza.

Ernesto Mordecki

Facultad de Ciencias, Universidad de la República. Montevideo, Uruguay

Curso de la Licenciatura en Matemática - 2020

Contenidos

Estadística (repaso): estimación de p

Intervalo de confianza

Ejemplo: un plebiscito

Ejemplo: el plebiscito (versión TCL)

Estadística (repaso): estimación de máxima verosimilitud

Problema: Observamos los resultados de un esquema de Bernoulli con n experimentos. Desconocemos $p \in (0, 1)$. Nos proponemos *estimar* el parámetro p . Mas precisamente:

- ▶ Sabemos que observamos un esquema de Bernoulli
- ▶ Observamos una realización ω . En particular sabemos $\mu(\omega)$
- ▶ Desconocemos $p \in (0, 1)$
- ▶ ¿Cómo podemos calcular (o al menos aproximar) p ?
- ▶ una propuesta intuitiva es tomar:

Estimación de máxima verosimilitud

$$\hat{p} \sim \frac{\mu(\omega)}{n}$$

¿Podemos dar un fundamento a esta intuición?

La propuesta es: Calculemos la probabilidad de la observación ω y elijamos el p que haga máxima esa probabilidad. Tenemos

- ▶ $\mathbf{P}(\omega, p) = p^k q^{n-k}$
- ▶ Ahora conocemos k pero desconocemos p (por eso escribimos la probabilidad como función de p).
- ▶ Resolvemos el problema

$$\max \mathbf{P}(\omega, p) \text{ para } p \in (0, 1)$$

Hagamos las cuentas

- ▶ El máximo de $\mathbf{P}(\omega, p)$ es el mismo que el de $\log \mathbf{P}(\omega, p)$.
Tenemos entonces

$$\ell(p) = \log \mathbf{P}(\omega, p) = \log p^k q^{n-k} = k \log p + (n - k) \log q$$

- ▶ Derivamos

$$\frac{\partial \ell(p)}{\partial p} = \frac{k}{p} - \frac{n - k}{1 - p} = 0.$$

- ▶ Operando: $k(1 - p) - p(n - k) = k - kp - pn + pk = 0$
- ▶ De aquí:

$$\hat{p} = \frac{k}{n} \text{ es el } \textit{estimador de máxima verosimilitud}.$$

Repasando

- ▶ Tenemos una estrategia de estimación, la *máxima verosimilitud* (máxima “creencia”)
- ▶ Calculamos la probabilidad de obtener la *muestra* que observamos, como función del parámetro (p en este caso) a estimar
- ▶ Calculamos el máximo de esa probabilidad como función de p
- ▶ Obtenemos el estimador de máxima verosimilitud \hat{p} .
- ▶ El sombrero $\hat{\theta}$ se usa para indicar que se está *estimando* un parámetro θ desconocido.

Conclusión

- ▶ En probabilidades tenemos un espacio

$$(\Omega, \mathcal{A}, \mathbf{P})$$

con \mathbf{P} conocida.

- ▶ A partir de esto calculamos probabilidades de sucesos que nos interesan. Después ocurre el azar (experimento)
- ▶ En estadística los experimentos se realizan primero, surgen *datos*, el ω , pero se desconoce \mathbf{P}
- ▶ Con esos datos se *estima* la probabilidad \mathbf{P} .

Intervalo de confianza

- ▶ Obtuvimos un valor aproximado de p , que es $\hat{p} = \mu(\omega)/n$, que se aproxima a p cuando n es grande.
- ▶ Queremos ahora dar un intervalo, de forma de tener una cierta certeza de no equivocarnos.
- ▶ Supongamos que esa certeza (o confianza) es de un 95%
- ▶ ¿Cómo hacemos?

- ▶ Los teoremas que vimos vienen en nuestra ayuda:
- ▶ Sabemos acotar y aproximar la probabilidad de un intervalo de la forma:

$$|\hat{p} - p| < \varepsilon$$

- ▶ Pero ahora lo que conocemos es \hat{p} , entonces el intervalo de confianza para p es de la forma

$$\hat{p} - \varepsilon < p < \hat{p} + \varepsilon.$$

- ▶ El problema es calcular ε para que el intervalo contenga 0.95 de probabilidad o más.

Intervalo de confianza: Teorema de Bernoulli

- ▶ El teorema de Bernoulli es una primera herramienta para construir intervalos de confianza.
- ▶ De su demostración, obtuvimos que

$$\mathbf{P}(|\hat{p} - p| \geq \varepsilon) \leq \frac{p(1-p)}{n\varepsilon^2} \leq \frac{1}{4n\varepsilon^2}.$$

- ▶ Entonces, calculamos ε para que

$$\frac{1}{4n\varepsilon^2} = 0.05.$$

- ▶ Es nos da

$$\varepsilon = \sqrt{\frac{1}{4n \times 0.05}} = \frac{2.24}{\sqrt{n}}.$$

- ▶ El intervalo de confianza para p es de la forma

$$\hat{p} - \frac{2.24}{\sqrt{n}} \leq p \leq \hat{p} + \frac{2.24}{\sqrt{n}}.$$

Ejemplo: un plebiscito

- ▶ Nos interesa conocer la opinión de una población respecto de un tema que se va a plebiscitar.
- ▶ El plebiscito se aprueba si recibe más de la mitad de los votos.
- ▶ Supongamos entonces que hacemos una encuesta a 1000 ciudadanos, y tenemos un porcentaje de afirmativos de 450 personas.
- ▶ Construir intervalos de confianza del 95 % y otro del 99% para la proporción verdadera en la población. Se usa decir que $\alpha = 0.05$ (el error admisible del intervalo).

Solución

- ▶ Tenemos $\hat{p} = \frac{450}{1000} = 0.45$.

- ▶ Por otra parte

$$\varepsilon = \frac{2.24}{\sqrt{1000}} = 0.071$$

- ▶ El intervalo para p es entonces de la forma

$$[0.45 - 0.071, 0.45 + 0.071] = [0.379, 0.521].$$

- ▶ Observemos que el intervalo contiene el valor 0.5.

Teorema integral de De-Moivre Laplace (Teorema Central del Límite: TCL)

Del teorema se obtiene, para $a > 0$, que

$$\mathbf{P} \left(-a < \frac{\mu - np}{\sqrt{npq}} \leq a \right) \sim \frac{1}{\sqrt{2\pi}} \int_{-a}^a e^{-x^2/2} dx.$$

Tenemos

$$\mathbf{P} \left(-a < \frac{\mu - np}{\sqrt{npq}} \leq a \right) = \mathbf{P} \left(-\frac{a}{\sqrt{n}} \sqrt{pq} < \frac{\mu}{n} - p \leq \frac{a}{\sqrt{n}} \sqrt{pq} \right)$$

Es decir

$$\varepsilon = \frac{a}{\sqrt{n}} \sqrt{pq} \sim \frac{a}{\sqrt{n}} \sqrt{\hat{p}(1 - \hat{p})}$$

Nos queda determinar a tal que

$$\frac{1}{\sqrt{2\pi}} \int_{-a}^a e^{-x^2/2} dx = 0.95.$$

Calculamos con la distribución normal. Queremos a tal que

$$\frac{1}{\sqrt{2\pi}} \int_{-a}^a e^{-x^2/2} dx = 0.95.$$

Es decir, por simetría, a verifica

$$\frac{1}{\sqrt{2\pi}} \int_0^a e^{-x^2/2} dx = \frac{0.95}{2} = 0.475$$

Entonces,

$$\Phi(a) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^a e^{-x^2/2} dx = 0.5 + 0.475 = 0.975.$$

Entonces, tenemos que hallar a tal que

$$\Phi(a) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^a e^{-x^2/2} dx = 0.975.$$

Utilizamos la función inversa de Φ en R :

```
> qnorm(0.975)  
[1] 1.959964
```

Obtuvimos que $a = 1.96$. El intervalo era

$$\varepsilon = \frac{a}{\sqrt{n}} \sqrt{pq} \sim \frac{a}{\sqrt{n}} \sqrt{\hat{p}(1 - \hat{p})} = \frac{1.96}{\sqrt{n}} \sqrt{\hat{p}(1 - \hat{p})}$$

Ejemplo: un plebiscito

- ▶ Nos interesa conocer la opinión de una población respecto de un tema que se va a plebiscitar.
- ▶ El plebiscito se aprueba si recibe más de la mitad de los votos.
- ▶ Supongamos entonces que hacemos una encuesta a 1000 ciudadanos, y tenemos un porcentaje de afirmativos de 450 personas.
- ▶ Construir intervalos de confianza del 95 % y otro del 99% para la proporción verdadera en la población. Se usa decir que $\alpha = 0.05$, el error admisible del intervalo.

Solución (II)

▶ Tenemos $\hat{p} = \frac{450}{1000} = 0.45$.

▶ Por otra parte

$$\varepsilon = \frac{1.96}{\sqrt{n}} \sqrt{\hat{p}(1 - \hat{p})} \leq \frac{1.96}{2\sqrt{n}} \sim \frac{1}{\sqrt{1000}} = 0.032.$$

▶ El intervalo para p es entonces de la forma

$$[0.45 - 0.032, 0.45 + 0.032] = [0.418, 0.482].$$

▶ Y no contiene a 0.5.

Observaciones

- ▶ Se pueden calcular intervalos de confianza con la acotación de grandes desvíos (pero el TCL es mejor)
- ▶ Se pueden calcular intervalos para otros niveles de confianza. Por ejemplo si $\alpha = 0.99$ tenemos que $a = 3$, y la amplitud del intervalo de confianza dado por el TCL es

$$\varepsilon = \frac{1.5}{\sqrt{n}}.$$

Recordemos que la amplitud para $\alpha = 0.05$ es $\varepsilon = \frac{1}{\sqrt{n}}$.

Ejercicio

Calcular el intervalo de confianza para $\alpha = 0.10$. Repasemos:

$$\mathbf{P} \left(-a < \frac{\mu - np}{\sqrt{npq}} \leq a \right) \sim \frac{1}{\sqrt{2\pi}} \int_{-a}^a e^{-x^2/2} dx.$$

equivalente a

$$\mathbf{P} \left(-\frac{a}{\sqrt{n}} \sqrt{pq} < \frac{\mu}{n} - p \leq \frac{a}{\sqrt{n}} \sqrt{pq} \right)$$

Lo que nos da una amplitud de intervalo

$$\varepsilon = \frac{a}{\sqrt{n}} \sqrt{pq} \leq \frac{a}{2\sqrt{n}}.$$

Calculamos con la distribución normal. Queremos a tal que

$$\frac{1}{\sqrt{2\pi}} \int_{-a}^a e^{-x^2/2} dx = 0.99.$$

Es decir, por simetría, a verifica

$$\frac{1}{\sqrt{2\pi}} \int_0^a e^{-x^2/2} dx = \frac{0.99}{2} = 0.495$$

Entonces,

$$\Phi(a) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^a e^{-x^2/2} dx = 0.5 + 0.495 = 0.995.$$

Entonces, tenemos que hallar a tal que

$$\Phi(a) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^a e^{-x^2/2} dx = 0.995.$$

Utilizamos la función inversa de Φ en R :

```
> qnorm(0.995)
[1] 2.575829
```

Obtenemos entonces que $a = 2.57$. El intervalo tiene amplitud

$$\varepsilon = \frac{a}{\sqrt{n}} \sqrt{pq} \leq \frac{2.57}{2\sqrt{n}}.$$

Buen finde!

