

INTRODUCCIÓN A LA INFERENCIA ESTADÍSTICA DISTRIBUCIÓN EN EL MUESTREO

Objetivos generales del tema

En este tema se introducirá el concepto de estadístico como medio para extraer información acerca de la ley de probabilidad del fenómeno en estudio. Se verá que, al ser función de v.a. también es una v.a. y se verá cómo obtener su distribución. Por último, se estudiarán algunos casos concretos de estadístico.

Principales contenidos

- Introducción
- Estimación de la media de una población
- Estimación de la varianza de una población
- Estimación de una proporción
- Estadísticos ordenados

1 Introducción

El muestreo estadístico es la herramienta que la Matemática utiliza para el estudio de las características de una población a través de una determinada parte de la misma.

La muestra de estudio debe ser lo más pequeña posible ya que del hecho de que una muestra sea más grande, no se desprende necesariamente que la información sea más fiable.

Además, la muestra elegida debe serlo por un proceso aleatorio para que sea lo más representativa posible.

Términos usuales en un estudio estadístico

- Población: conjunto de todos los individuos que son objeto del estudio.
- Muestra: parte de la población en la que miden las características estudiadas.
- Muestreo: proceso seguido para la extracción de una muestra.
- Encuesta: proceso de obtener información de la muestra.

Los resultados obtenidos del estudio de **una muestra** no son del todo fiables, pero sí en buena medida. Los parámetros que obtienen de una muestra (estimadores estadísticos) nos permitirán arriesgarnos a predecir una serie de resultados para toda la población. De estas predicciones y del riesgo que conllevan se ocupa la Inferencia Estadística.

Las observaciones de una muestra se denotan por x_1, \dots, x_n .

Sin embargo, antes de hacer un muestreo o de experimentar, cualquier observación en particular estará sujeta a incertidumbre (por ejemplo, antes de saber cuál es el gasto medio de una familia de la muestra en alimentación, ésta podría ser $x_1 = 125.000$ ó $x_1 = 87.000$ ó muchos otros valores posibles).

Debido a esta incertidumbre, antes de que se disponga de datos numéricos concretos, consideramos las observaciones como variables aleatorias y las denotamos por letras mayúsculas X_1, \dots, X_n .

Esto a su vez implica que hasta que se hayan obtenido los datos, cualquier función de las observaciones (media muestral, varianza de la muestra, etc.) son funciones de variables aleatorias, y por tanto variables aleatorias con distribución de probabilidad propia, llamada *distribución en el muestreo*.

Definition 1 Se llama **espacio muestral** al conjunto de muestras posibles que pueden obtenerse al seleccionar una muestra aleatoria, de un determinado tamaño, de una cierta población.

Definition 2 Se llama **estadístico** a cualquier función $T(X_1, \dots, X_n)$ de la muestra (X_1, \dots, X_n) . El estadístico $T(X_1, \dots, X_n)$, como función de variables aleatorias (X_1, \dots, X_n) , es también una variable aleatoria, y tendrá por tanto una distribución de probabilidad, llamada **distribución en el muestreo**.

Algunos estadísticos importantes:

$$T(X_1, \dots, X_n) = \bar{X} = \frac{X_1 + \dots + X_n}{n} \quad \text{Media Muestral}$$

$$T(X_1, \dots, X_n) = \text{Min}(X_1, \dots, X_n)$$

$$T(X_1, \dots, X_n) = \text{Max}(X_1, \dots, X_n)$$

$$T(X_1, \dots, X_n) = S_X^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n} \quad \text{Varianza Muestral}$$

$$T(X_1, \dots, X_n) = \hat{S}_X^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1} \quad \text{Cuasivarianza Muestral}$$

Example 3 *Se está interesado en conocer la probabilidad θ de obtener cara con una moneda, es decir, se trata de estudiar la variable aleatoria*

$$X = \begin{cases} 1 & \text{si se obtiene cara} \\ 0 & \text{si se obtiene cruz} \end{cases}$$

cuya distribución está caracterizada por

| X_i | Probabilidad |
|-------|--------------|
| 1 | θ |
| 0 | $1 - \theta$ |

que depende del parámetro θ que varía en el espacio paramétrico $\Theta = [0, 1]$. Se realizan tres lanzamientos, con lo que se dispone de una m.a.s X_1, X_2, X_3 . Puesto que la muestra es aleatoria simple, se verifica que

$$\mathbb{P}(X_1 = x_1, X_2 = x_2, X_3 = x_3) = \mathbb{P}(X_1 = x_1) \cdot \mathbb{P}(X_2 = x_2) \cdot \mathbb{P}(X_3 = x_3)$$

La probabilidad de todas las muestras posibles es, por tanto, la siguiente:

| X_1 | X_2 | X_3 | Probabilidad |
|-------|-------|-------|------------------------|
| 1 | 1 | 1 | θ^3 |
| 1 | 1 | 0 | $\theta^2(1 - \theta)$ |
| 1 | 0 | 1 | $\theta^2(1 - \theta)$ |
| 0 | 1 | 1 | $\theta^2(1 - \theta)$ |
| 1 | 0 | 0 | $\theta(1 - \theta)^2$ |
| 0 | 1 | 0 | $\theta(1 - \theta)^2$ |
| 0 | 0 | 1 | $\theta(1 - \theta)^2$ |
| 0 | 0 | 0 | $(1 - \theta)^3$ |

Es decir, tenemos la distribución de la muestra a través de su función de

probabilidad. Si nos interesa el estadístico “media muestral”

$$\bar{X} = \frac{X_1 + X_2 + X_3}{3}$$

a partir de la distribución en el muestreo de la muestra, se tiene la siguiente distribución en el muestreo para el estadístico:

| \bar{X} | Probabilidad |
|-----------|-----------------------|
| 1 | θ^3 |
| 2/3 | $3\theta^2(1-\theta)$ |
| 1/3 | $3\theta(1-\theta)^2$ |
| 0 | $(1-\theta)^3$ |

Example 4 De una población X con distribución de Poisson de parámetro λ , $Poisson(\lambda)$, se obtiene una m.a.s. de tamaño n . Determinar la distribución en el muestreo de la media muestral \bar{X} .

Si $X \sim Poisson(\lambda)$, entonces $\mathbb{P}(X = k) = e^{-\lambda} \frac{\lambda^k}{k!}$, $k = 0, 1, 2, \dots$

Dada una m.a.s. (X_1, \dots, X_n) , $T(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}$. Su distribución es:

$$\mathbb{P}(\bar{X} = k) = \mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n X_i = k\right) = \mathbb{P}\left(\sum_{i=1}^n X_i = nk\right) = e^{-n\lambda} \frac{(n\lambda)^{nk}}{(nk)!}$$

\uparrow
 $\sum_{i=1}^n X_i \sim Poisson(n\lambda)$

Definition 5 Estimator Dada una muestra (X_1, \dots, X_n) , un estimador del parámetro θ es una función de la muestra $T(X_1, \dots, X_n)$ que aproxima el valor de θ . Se suele denotar por $T(X_1, \dots, X_n) = \hat{\theta}$

Definition 6 Estimator Insesgado: Sea $T(X_1, \dots, X_n)$ un estadístico que estima un parámetro θ , denotémoslo por $T(X_1, \dots, X_n) = \hat{\theta}$. El **sesgo** del estimador de θ es

$$Sesgo[\hat{\theta}] = E[\hat{\theta}] - \theta = E[T(X_1, \dots, X_n)] - \theta$$

El estimador es insesgado si $\text{Sesgo}[\hat{\theta}] = 0$, es decir $E[\hat{\theta}] = \theta$.

Definition 7 Es un estimador **asintóticamente insesgado** si

$$\text{Sesgo}[\hat{\theta}] \rightarrow 0 \text{ cuando } n \rightarrow \infty$$

Definition 8 Estimador Consistente: Sea $T(X_1, \dots, X_n)$ un estadístico que estima un parámetro θ , denotémoslo por $T(X_1, \dots, X_n) = \hat{\theta}$. El estimador es consistente si

$$\text{Var}[\hat{\theta}] = \text{Var}[T(X_1, \dots, X_n)] \rightarrow 0 \text{ cuando } n \rightarrow \infty$$

2 Estimación de la media de una población

Sea X_1, \dots, X_n una m.a.s. de una variable aleatoria X con media $\mathbb{E}(X) = \mu$ y varianza $\text{Var}(X) = \sigma^2$. El estimador más razonable de la media poblacional μ es la **media muestral**

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

que verifica las siguientes propiedades:

1. Es un **estimador insesgado** de μ .

$$\mathbb{E}(\bar{X}) = \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i) = \frac{1}{n} \sum_{i=1}^n \mu = \mu$$

2. Es un estimador **consistente en media cuadrática** de μ , puesto que es insesgado y su varianza tiende a 0 cuando $n \rightarrow \infty$.

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{n}{n^2} \sigma^2 = \frac{\sigma^2}{n}$$

3. La distribución de exacta de \bar{X} depende de la distribución de la población X . Por ejemplo, si X es normal la distribución de \bar{X} también lo será. Para muestras grandes, por el Teorema Central del Límite, la distribución de \bar{X} puede aproximarse por una normal ($n \geq 30$).

Theorem 9 Sea X_1, \dots, X_n una m.a.s. de una variable aleatoria $X \sim N(\mu, \sigma)$, entonces

$$\boxed{\bar{X} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)}$$

Example 10 Sea X una población con distribución $N(90, \sigma = 20)$.

a) Si se obtiene una m.a.s. de tamaño 16, ¿cuál es la probabilidad de que la media muestral \bar{X} sea mayor o igual que 92?

b) Determinar el tamaño muestral para que la probabilidad de que la media muestral sea menor o igual que 98 sea $\mathbb{P}(\bar{X} \leq 98) = 0,99$.

$$X \sim N(\mu = 90, \sigma = 20)$$

$$\left. \begin{aligned} \mathbb{E}(\bar{X}) &= \mu = 90 \\ \text{Var}(\bar{X}) &= \frac{\sigma^2}{n} = \frac{20^2}{16} \Rightarrow \sigma_{\bar{X}} = \sqrt{\frac{400}{16}} = \frac{20}{4} = 5 \end{aligned} \right\} \Rightarrow \bar{X} \sim N(90, 5)$$

$$a) \mathbb{P}(\bar{X} \geq 92) = \mathbb{P}(N(90, 5) \geq 92) = \mathbb{P}\left(N(0, 1) \geq \frac{92 - 90}{5}\right) =$$

$$= \mathbb{P}\left(Z \geq \frac{2}{5}\right) = \mathbb{P}(Z \geq 0.4) = \boxed{0,3446}$$

$$b) 0.99 = \mathbb{P}(\bar{X} \leq 98) = \mathbb{P}\left(N\left(90, \frac{20}{\sqrt{n}}\right) \leq 98\right) = \mathbb{P}\left(N(0, 1) \leq \frac{98 - 90}{\frac{20}{\sqrt{n}}}\right) =$$

$$= \mathbb{P}\left(Z \leq \sqrt{n} \frac{2}{5}\right) \Rightarrow \sqrt{n} \frac{2}{5} = 2,33 \Rightarrow n = \left(\frac{2,33 \cdot 5}{2}\right)^2 = 33,9 \Rightarrow \boxed{n \geq 34}$$

3 Estimación de la varianza de una población

Sea X_1, \dots, X_n una m.a.s. de una variable aleatoria X con media $\mathbb{E}(X) = \mu$ y varianza $Var(X) = \sigma^2$. El estimador más razonable de la varianza poblacional σ^2 es la **varianza muestral**

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

que verifica las siguientes propiedades:

1. Es un **estimador consistente** de σ^2 .
2. Es un estimador **asintóticamente insesgado** de σ^2 puesto que

$$\mathbb{E}(S_X^2) = \frac{\sigma^2}{n} + \sigma^2$$

Puesto que la varianza muestral S^2 es un estimador sesgado (aunque asintóticamente insesgado), para estimar la varianza σ^2 se usa la **cuasivarianza muestral**:

$$\widehat{S}_x^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{n}{n-1} S^2$$

que es un estimador **insesgado** y **consistente** de σ^2 .

$$\mathbb{E}(\widehat{S}^2) = \mathbb{E}\left(\frac{n}{n-1} S^2\right) = \frac{n}{n-1} \mathbb{E}(S^2) = \frac{n}{n-1} \left(\sigma^2 - \frac{\sigma^2}{n}\right) = \sigma^2$$

3. (**Teorema de Fisher**) Bajo hipótesis de normalidad ($X \sim N(\mu, \sigma)$) se

verifica que

$$\frac{nS^2}{\sigma^2} = \frac{(n-1)\widehat{S}^2}{\sigma^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} \sim \chi_{n-1}^2$$

4. Si conocemos la varianza poblacional μ , se verifica

$$\frac{\sum_{i=1}^n (X_i - \mu)^2}{\sigma^2} \sim \chi_n^2$$

5. Bajo hipótesis de normalidad ($X \sim N(\mu, \sigma)$), las variables aleatorias \bar{X} y \widehat{S}^2 son independientes.

Example 11 *Dada una población $X \sim N(6, \sigma = 2.5)$, y tomando una muestra aleatoria simple X_1, \dots, X_n de tamaño $n = 12$, calcular la probabilidad de que la varianza muestral sea mayor que 4.9.*

$$\begin{aligned} \mathbb{P}(S^2 > 4.9) &= 1 - \mathbb{P}(S^2 \leq 4.9) = 1 - \mathbb{P}\left(\frac{nS^2}{\sigma^2} \leq \frac{n \cdot 4.9}{\sigma^2}\right) = \\ &= 1 - \mathbb{P}\left(\chi_{n-1}^2 \leq \frac{12 \cdot 4.9}{2.5^2}\right) = 1 - \mathbb{P}(\chi_{11}^2 \leq 9.41) \end{aligned}$$

En la tabla de la distribución de la χ_n^2 aparecen los siguientes resultados

$$\mathbb{P}(\chi_{11}^2 \leq 7.58) = 0.25$$

$$\mathbb{P}(\chi_{11}^2 \leq 10.03) = 0.5$$

Mediante una interpolación numérica se calcula la probabilidad que in-

teresa

| x | y |
|-------|------|
| 7.58 | 0.25 |
| 9.41 | p |
| 10.03 | 0.5 |

$$y = y_1 + \frac{y_2 - y_1}{x_2 - x_1} (x - x_1) = 0.25 + \frac{0.5 - 0.25}{10.03 - 7.58} (x - 7.58)$$

$$y = 0.102x - 0.52347 \implies p = 0.102 \cdot 9.41 - 0.52347 = 0.43635$$

$$\mathbb{P}(S^2 > 4.9) = 1 - \mathbb{P}(\chi_{11}^2 \leq 9.41) = 1 - 0.43635 = \boxed{0.56365}$$

3.1 Estadístico t de Student. Estimación de la media poblacional cuando σ^2 es desconocida.

Bajo hipótesis de normalidad ($X \sim N(\mu, \sigma)$) se verifica que $\bar{X} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$, o equi-valentemente

$$\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$$

Este resultado puede ser de poca utilidad si la varianza poblacional σ^2 es desconocida, ya que entonces no se podrá usar esta conclusión para hacer previsiones acerca de \bar{X} . Cabe pensar que el resultado no será muy distinto si se sustituye σ por la cuasidesviación típica muestral \hat{S} , puesto que, al menos para muestras grandes, σ^2 y \hat{S}^2 tendrán valores semejantes. Tal idea llevó a Student (pseudónimo de W. Gosset) a considerar el estadístico

$$\frac{\bar{X} - \mu}{\frac{\hat{S}}{\sqrt{n}}} = \frac{\bar{X} - \mu}{\frac{S}{\sqrt{n-1}}} \sim t_{n-1}$$

Example 12 *Dada una población $X \sim N(-1, \sigma)$, se extrae una m.a.s. de tamaño $n = 10$ con los siguientes resultados:*

1.08, -1.79, -2.54, 0.37, -0.6, 0.53, 0.28, -2.21, -2.66, 1.45

Calcular la probabilidad de que la media muestral \bar{X} sea mayor que -1.2.

$$\mathbb{P}(\bar{X} > -1.2) = 1 - \mathbb{P}(\bar{X} \leq -1.2) = 1 - \mathbb{P}\left(\frac{\bar{X} - \mu}{\frac{\hat{S}}{\sqrt{n}}} \leq \frac{-1.2 - \mu}{\frac{\hat{S}}{\sqrt{n}}}\right)$$

$$S^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2 = \frac{1}{10}(1.08^2 + 1.79^2 + 2.54^2 + 0.37^2 + 0.6^2 + 0.53^2 + 0.28^2$$

$$+ 2.21^2 + 2.66^2 + 1.45^2) - \left[\frac{1}{10}(1.08 + 1.79 + 2.54 + 0.37 + 0.6 + 0.53 +$$

$$+ 0.28 + 2.21 + 2.66 + 1.45)\right]^2 = \frac{1}{10} \cdot 25.7405 - 0.609^2 = 2.20317$$

$$S = \sqrt{2.20317} = 1.4843 \implies \hat{S} = \sqrt{\frac{n}{n-1} S^2} = \sqrt{\frac{10}{9} 2.20317} = 1.5646$$

$$\mathbb{P}(\bar{X} > -1.2) = 1 - \mathbb{P}\left(\frac{\bar{X} - \mu}{\frac{\hat{S}}{\sqrt{n}}} \leq \frac{-1.2 - \mu}{\frac{\hat{S}}{\sqrt{n}}}\right) = 1 - \mathbb{P}\left(t_{n-1} \leq \frac{-1.2 - (-1)}{\frac{1.5646}{\sqrt{10}}}\right) =$$

$$= 1 - \mathbb{P}\left(t_9 \leq \frac{-1.2 - (-1)}{\frac{1.5646}{\sqrt{10}}}\right) = 1 - \mathbb{P}(t_9 \leq -0.4042) = \mathbb{P}(t_9 \leq 0.4042)$$

Se calcula la probabilidad anterior por interpolación numérica

| x | y |
|--------|-----|
| 0.261 | 0.6 |
| 0.4042 | p |
| 0.543 | 0.7 |

$$y = y_1 + \frac{y_2 - y_1}{x_2 - x_1} (x - x_1) = 0.6 + \frac{0.7 - 0.6}{0.543 - 0.261} (x - 0.4042)$$

$$y = 0.35462x + 0.456 \implies p = 0.35462 \cdot 0.4042 + 0.456 = 0.6$$

$$\mathbb{P}(\bar{X} > -1.2) = 0.6$$

4 Estimación de una proporción

Se desea estimar la proporción p de individuos de una población que tiene una determinada característica. Para ello se toma una m.a.s. de elementos de la población, anotando un 1 si dicho elemento tiene la característica, y 0 en otro caso, es decir, se tiene una m.a.s. X_1, \dots, X_n de una $B(1, p)$

$$X = \begin{cases} 1 & \text{(tiene la caract.) con probab. } p \\ 0 & \text{(no la tiene) con probab. } 1 - p \end{cases}, \quad \mathbb{E}(X) = p, \quad \text{Var}(X) = p(1 - p)$$

Un estimador razonable de p es la proporción de elementos de la muestra que tiene dicha característica, es decir

$$\hat{p} = \frac{X_1 + \dots + X_n}{n}$$

Propiedades

1. $\mathbb{E}(\hat{p}) = \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i) = \frac{1}{n} np = p$
2. $Var(\hat{p}) = Var\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n p(1-p) = \frac{p(1-p)}{n}$
3. La distribución de \hat{p} depende de la distribución de la población X , pero cuando n es grande

$$\hat{p} \sim N\left(p, \sqrt{\frac{p(1-p)}{n}}\right)$$

Example 13 *En el proceso de producción de una empresa, el 1% de los productos sale defectuoso. Para corroborarlo se obtiene una m.a.s. de tamaño $n = 25$ y se estima la proporción de productos defectuosos. Estimar la probabilidad de que la proporción estimada sea mayor que el 2%.*

$$\begin{aligned} \mathbb{P}(\hat{p} > 0.02) &= \mathbb{P}\left(\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} > \frac{0.02 - p}{\sqrt{\frac{p(1-p)}{n}}}\right) \simeq \mathbb{P}\left(Z > \frac{0.02 - \hat{p}}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}}\right) = \\ &= \mathbb{P}\left(Z > \frac{0.02 - 0.01}{\sqrt{\frac{0.01(1-0.01)}{25}}}\right) = \mathbb{P}(Z > 0.50) = 0.3085 \end{aligned}$$

$$\mathbb{P}(\hat{p} > 0.02) = 0.3085$$

5 Estadísticos ordenados

Sea X_1, \dots, X_n una m.a.s de una población X con función de distribución $F(x)$ y densidad $f(x)$. Es importante estudiar entre qué valores podrían estar los valores muestrales; se consideran entonces $X_{(1)}, \dots, X_{(n)}$ los estadísticos de orden (los valores muestrales ordenados de menor a mayor ($X_{(1)} \leq \dots \leq X_{(n)}$)). Aunque X_1, \dots, X_n son independientes idénticamente distribuidos (i.i.d.) por tratarse de una m.a.s., $X_{(1)}, \dots, X_{(n)}$ no lo son.

Ejemplos

$$X_{(1)} = \min(X_1, \dots, X_n)$$

$$X_{(n)} = \max(X_1, \dots, X_n)$$