



Muestras aleatorias y distribuciones de muestreo.

Definición 1.

Si las variables aleatorias X_1, X_2, \dots, X_n tienen la misma función de probabilidad que la distribución de la población y su distribución de probabilidad conjunta es igual al producto de las marginales, entonces X_1, X_2, \dots, X_n forma un conjunto de n variables aleatorias independientes e idénticamente distribuidas (IID) que constituye una muestra aleatoria de la población.

La función (densidad) conjunta de probabilidad de X_1, X_2, \dots, X_n es la función de verosimilitud de la muestra dada por:

$$L(\underline{x}, \theta) = \prod_{i=1}^n f(x_i; \theta)$$

en donde

$$\underline{x} = \{x_1, x_2, \dots, x_n\}$$

denota los datos muestreados.

Cuando las realizaciones de

$$\underline{x} = \{x_1, x_2, \dots, x_n\}$$

se conocen, la función

$$L(\underline{x}, \theta) = \prod_{i=1}^n f(x_i; \theta)$$

depende sólo del parámetro desconocido θ .

Definición 2

Un parámetro es una característica numérica de la distribución de la población de manera que cuando se conoce esta, la distribución, queda descrita, sino total, al menos parcialmente.

Los parámetros o funciones de los parámetros se estiman a partir de la información contenida en una muestra.

Definición 3.

Un estadístico es cualquier función de las variables aleatorias que se observaron en la muestra de manera que esta función no contiene cantidades desconocidas.

Si se utiliza un estadístico T para estimar un parámetro desconocido θ , entonces T recibe el nombre de estimador de θ y el valor específico que tome T, por ejemplo t, se denomina estimación de θ .

Definición 4

La distribución de muestreo de un estadístico T es la distribución de probabilidad de T que puede obtenerse como resultado de un número infinito de muestras aleatorias independientes, cada una de tamaño n, provenientes de la población de interés.

Teorema 1

Sea X_1, X_2, \dots, X_n un conjunto de n variables aleatorias independientes cada una con función generadora de momentos

$$m_{X_1}(t), m_{X_2}(t), \dots, m_{X_n}(t)$$

Si $Y = a_1X_1 + a_2X_2 + \dots + a_nX_n$ en donde a_1, a_2, \dots, a_n son constantes, entonces:

$$m_Y(t) = m_{X_1}(a_1t) \cdot m_{X_2}(a_2t) \cdot \dots \cdot m_{X_n}(a_nt)$$

Teorema 2

Sea X_1, X_2, \dots, X_n un conjunto de variables aleatorias independientes y normalmente distribuidas con medias $E(X_i)$ y varianzas $\text{Var}(X_i) = \sigma_i^2$ para $i = 1, 2, \dots, n$. Si $Y = a_1X_1 + a_2X_2 + \dots + a_nX_n$ en donde a_1, a_2, \dots, a_n son constantes, entonces Y es una variable aleatoria normal con media

$$E(Y) = a_1\mu_1 + a_2\mu_2 + \dots + a_n\mu_n$$

y con varianza

$$\text{var}(Y) = a_1^2\sigma_1^2 + a_2^2\sigma_2^2 + \dots + a_n^2\sigma_n^2$$

(La hipótesis de normalidad puede quitarse)

La distribución de muestreo de la media muestral \bar{X} .

Sea X_1, X_2, \dots, X_n una muestra aleatoria que consiste en un conjunto de variables aleatorias independientes e idénticamente distribuidas (v.a IID) tales que $E(X_i)=\mu$ y $\text{Var}(X_i)=\sigma^2$ para todo i .

Entonces el estadístico

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

se define como la media de las n v.a IID.

Aplicando el teorema 2 se tiene :

$$(a_i = \frac{1}{n}, \forall i)$$

$$E(\bar{X}) = \frac{1}{n} \sum_i E(X_i) = \frac{1}{n} \sum_i \mu = \frac{1}{n} (n\mu) = \mu$$

$$\text{Var}(\bar{X}) = \sum_i \frac{1}{n^2} \sigma^2 = \frac{\sigma^2}{n}$$

de donde

$$d.e.(\bar{X}) = \frac{\sigma}{\sqrt{n}}$$

Teorema 3

Sea X_1, X_2, \dots, X_n una muestra aleatoria que consiste en un conjunto de variables aleatorias independientes y normalmente distribuidas tales que $E(X_i)=\mu$ y $\text{Var}(X_i)=\sigma^2$ para $i = 1, 2, \dots, n$.

Entonces la distribución de la media muestral \bar{X} es normal con media μ y varianza σ^2/n .

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

Un problema de máximo interés consiste en saber lo que ocurre si no se especifica la distribución de probabilidad de la población a partir de la cual se extrae la muestra.

Teorema central del límite.

Sea X_1, X_2, \dots, X_n una muestra aleatoria que consiste en un conjunto de variables aleatorias independientes e idénticamente distribuidas (v.a IID) en una distribución de probabilidad no especificada y tales que $E(X_i)=\mu$ y $\text{Var}(X_i)=\sigma^2$ para todo i . Entonces el promedio muestral

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

tiene una distribución de media μ y varianza σ^2/n que tiende a una normal conforme n tiende a ∞

Debe de notarse el hecho de que si el modelo de probabilidad de la población es semejante a una distribución normal, la aproximación normal será buena aun para muestras pequeñas.

En general, para $n > 30$, la aproximación normal será relativamente buena y puede emplearse

$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$$

para hacer inferencias sobre μ cuando se conoce el valor de la varianza poblacional σ^2 .

Distribución en el muestreo de S^2

Supongamos que la población se encuentra normalmente distribuida con μ conocida y σ^2 desconocida. Se define S^2 como

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

donde para cada i

$$X_i \rightarrow N(\mu, \sigma)$$

Teorema 5

Sea X_1, X_2, \dots, X_n una muestra aleatoria de una distribución normal de media μ y varianza σ^2 . La distribución de la variable aleatoria

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

es del tipo de la Chi-cuadrado con n grados de libertad.

Desde un punto de vista práctico, la varianza muestral tal y como se encuentra definida anteriormente tiene poco uso, pues rara vez se conoce la media poblacional μ . En su lugar se emplea la varianza muestral, definida por:

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

Más adelante se verá porqué se emplea el divisor $n-1$ en lugar de n .

El reemplazo de la media desconocida μ por la media muestral \bar{x} da origen a la presencia de otro estadístico en la definición de S^2 . Como consecuencia se tiene que la distribución de muestreo

$$\frac{(n-1)S^2}{\sigma^2} \rightarrow \chi_{n-1}^2$$

Teorema 6

Sea X_1 y X_2 son variables aleatorias independientes y cada una tiene una distribución Chi-cuadrado con ν_1 y ν_2 grados de libertad, entonces $Y = X_1 + X_2$ tiene también una

Chi-cuadrado con $\nu_1 + \nu_2$ grados de libertad.

Teorema 7

$$\frac{(n-1)S^2}{\sigma^2} \rightarrow \chi_{n-1}^2$$

donde

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

Demostración

$$(n-1)S^2 = \sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n [(X_i - \mu) + (\mu - \bar{X})]^2 = \sum_{i=1}^n [(X_i - \mu) - (\bar{X} - \mu)]^2 = \dots = \sum_{i=1}^n [(X_i - \mu)]^2 - n(\bar{X} - \mu)^2$$

de donde

$$(n-1)S^2 + n(\bar{X} - \mu)^2 = \sum_{i=1}^n (X_i - \mu)^2$$

Dividiendo los dos miembros por σ^2

$$\frac{(n-1)S^2}{\sigma^2} + \left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right)^2 = \frac{\sum_{i=1}^n (X_i - \mu)^2}{\sigma^2}$$

Por el teorema 5

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

sigue una Chi-cuadrado con n grados de libertad y

$$\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right)^2$$

de manera similar sigue una chi-cuadrado con 1 grado de libertad, dado que

$$\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right) \rightarrow N(0, 1)$$

En virtud del teorema 6 se sigue que

$$\frac{(n-1)S^2}{\sigma^2} \rightarrow \chi_{n-1}^2$$

La distribución t de Student

Se sabe, cuando una muestra proviene de una distribución normal con desviación estándar conocida σ , que la distribución de

$$Z = \left(\frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \right) \rightarrow N(0, 1)$$

Generalmente el valor de σ no se conoce y lo que se hace es reemplazar σ por un estimador s , que es el valor de la desviación estándar muestral S .

Desafortunadamente

$$\left(\frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \right)$$

ya no es $N(0, 1)$ aún cuando la muestra provenga de una distribución normal.

Sin embargo, es posible determinar la distribución muestral exacta de

$$\left(\frac{\bar{X} - \mu}{S / \sqrt{n}} \right)$$

cuando la población es normal $N(\mu, \sigma)$ con μ y σ desconocidos.

Teorema 8

Sea Z normal $N(0, 1)$ y X una Chi-cuadrado con ν grados de libertad. Si Z y X son independientes, entonces la variable aleatoria

$$T = \frac{Z}{\sqrt{X/\nu}}$$

sigue una t-Student con ν grados de libertad.

La similitud de la t-Student y la $N(0,1)$ es alta para valores grandes de ν , sobre todo para $\nu \geq 30$.

Teorema 9

Cuando se muestrea una población normal $N(\mu, \sigma)$, el estadístico

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

sigue una t-Student con n-1 grados de libertad.

Distribución de la diferencia de dos medias muestrales

Supongamos que X sigue una $N(\mu_x, \sigma)$ e Y una $N(\mu_y, \sigma)$ donde X e Y son variables aleatorias independientes con **varianzas iguales conocidas**.

Sabemos que

$$\bar{X} \rightarrow N(\mu_x, \sigma/\sqrt{n_x}) \quad e \quad \bar{Y} \rightarrow N(\mu_y, \sigma/\sqrt{n_y})$$

Entonces

$$\bar{X} - \bar{Y} = \bar{X} + (-1)\bar{Y}$$

sigue una distribución normal de media $\mu_x - \mu_y$ y varianza $1^2 \cdot \sigma^2/n_x + (-1)^2 \cdot \sigma^2/n_y$

Por tanto si se conoce el valor de σ^2 , el estadístico

$$Z = \frac{\bar{X} - \bar{Y} - (\mu_x - \mu_y)}{\sigma \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}}$$

se distribuye según una normal de $N(0; 1)$

Se ha supuesto que σ es conocido. Sin embargo, es poco probable que esto suceda. Por tanto para el caso en el que el muestreo se lleve a cabo sobre dos poblaciones normales independientes con **varianzas iguales pero desconocidas**, para cada una de las muestras obtenidas pueden definirse las varianzas muestrales definidas S_x^2 y S_y^2 y dado que

$$\frac{(n_x - 1)S_x^2}{\sigma^2} \rightarrow \chi^2_{n_x - 1}$$

$$\frac{(n_y - 1)S_y^2}{\sigma^2} \rightarrow \chi^2_{n_y - 1}$$

y teniendo en cuenta el teorema 6

$$W = \frac{(n_x - 1)S_x^2}{\sigma^2} + \frac{(n_y - 1)S_y^2}{\sigma^2}$$

sigue también una chi-cuadrado con $n_x + n_y - 2$ grados de libertad

Luego

$$T = \frac{Z}{\sqrt{\frac{W}{(n_x + n_y - 2)}}}$$

sigue una t-Student con $n_x + n_y - 2$ grados de libertad.

Por tanto

$$T = \frac{\frac{\bar{X} - \bar{Y} - (\mu_x - \mu_y)}{\sigma \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}}}{\frac{1}{\sigma} \sqrt{\frac{(n_x - 1)S_x^2 + (n_y - 1)S_y^2}{n_x + n_y - 2}}}$$

simplificando

$$T = \frac{\bar{X} - \bar{Y} - (\mu_x - \mu_y)}{\sqrt{\frac{(n_x - 1)S_x^2 + (n_y - 1)S_y^2}{n_x + n_y - 2} \left(\frac{1}{n_x} + \frac{1}{n_y} \right)}} = \frac{\bar{X} - \bar{Y} - (\mu_x - \mu_y)}{S_p \sqrt{\left(\frac{1}{n_x} + \frac{1}{n_y} \right)}}$$

donde

$$S_p = \sqrt{\frac{(n_x - 1)S_x^2 + (n_y - 1)S_y^2}{n_x + n_y - 2}}$$

es un estimador combinado de la varianza.

Si suponemos que las **varianzas poblacionales son distintas pero conocidas**, entonces se tiene:

$$\bar{X} \rightarrow N\left(\mu_x, \frac{\sigma_x^2}{n_x}\right)$$

$$\bar{Y} \rightarrow N\left(\mu_y, \frac{\sigma_y^2}{n_y}\right)$$

De donde

$$\bar{X} - \bar{Y} \rightarrow N\left(\mu_x - \mu_y, \sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}\right)$$

y por tanto

$$Z = \frac{\bar{X} - \bar{Y} - \mu_x - \mu_y}{\sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}} \rightarrow N(0,1)$$

En el **supuesto de que σ_x^2 y σ_y^2 sean desconocidas** y haya que estimarlas a partir de S_x^2 y S_y^2 el problema se complica, en tal caso se tiene:

$$t = \frac{(\bar{X} - \bar{Y}) - (\mu_x - \mu_y)}{\sqrt{\frac{S_x^2}{n_x} + \frac{S_y^2}{n_y}}}$$

sigue una distribución t-Student con f grados de libertad, en donde f es la aproximación de Welch:

$$f = \frac{\left(\frac{S_x^2}{n_x} + \frac{S_y^2}{n_y}\right)^2}{\frac{\left(\frac{S_x^2}{n_x}\right)^2}{n_x + 1} + \frac{\left(\frac{S_y^2}{n_y}\right)^2}{n_y + 1}}$$

expresado en número entero.

La distribución F

La inferencias con respecto a la varianza σ^2 cuando se muestrea una población normal se formula con base a

$$\frac{(n-1)S^2}{\sigma^2} \rightarrow \chi^2_{n-1}$$

En esta sección se formularán inferencias con respecto a la varianza de dos distribuciones normales independientes con base en las muestras aleatorias de cada una.

Teorema 10

Sea X una variable aleatoria que se distribuye según una Chi-cuadrado con v_1 grados de libertad e Y otra variable aleatoria independiente de X, que se distribuye según una Chi-cuadrado con v_2 grados de libertad. Entonces la variable aleatoria

$$F = \frac{X/v_1}{Y/v_2}$$

tiene una función de densidad de probabilidad dada por

$$g(f; v_1, v_2) = \begin{cases} \frac{\Gamma[(v_1 + v_2)/2] v_1^{v_1/2} v_2^{v_2/2}}{\Gamma(v_1/2)\Gamma(v_2/2)} f^{(v_1-2)/2} (v_2 + v_1 f)^{-(v_1+v_2)/2} & \text{Si } f > 0 \\ 0 & f \leq 0 \end{cases}$$

Es fácil ver que si

$$F = \frac{X/v_1}{Y/v_2}$$

sigue una distribución F con v_1, v_2 grados de libertad, entonces $F' = 1/F$ sigue una F con v_2, v_1 grados de libertad

$$F' = \frac{1}{F} = \frac{Y/v_2}{X/v_1}$$

Es por esto que en las tablas sólo aparecen los valores cuantiles $f_{1-\alpha;v_1;v_2}$ para $\alpha < 0.5$. Si se desean los valores cuantiles para $\alpha > 0.5$

$$P(F \leq f_{1-\alpha;v_1;v_2}) = P\left(\frac{1}{F} > \frac{1}{f_{1-\alpha;v_1;v_2}}\right) = 1 - \alpha$$

o bien

$$P(F' \leq f'_{\alpha;v_2;v_1}) = \alpha$$

siendo

$$f'_{\alpha;v_2;v_1} = \frac{1}{f_{1-\alpha;v_1;v_2}}$$

Volviendo al problema de desarrollar estadísticos para formular inferencias con respecto a las varianzas de dos distribuciones normales independientes.

Sean X_1, X_2, \dots, X_{n_x} variables aleatorias independientes $N(\mu_x, \sigma_x)$

Sean Y_1, Y_2, \dots, Y_{n_y} variables aleatorias independientes $N(\mu_y, \sigma_y)$

Si X e Y son independientes

$$\frac{(n_x - 1)S_x^2}{\sigma_x^2} \rightarrow \chi^2_{n_x - 1}$$

$$\frac{(n_y - 1)S_y^2}{\sigma_y^2} \rightarrow \chi^2_{n_y - 1}$$

entonces por el teorema 10

$$\frac{\frac{(n_x - 1)S_x^2}{\sigma_x^2}}{\frac{(n_y - 1)S_y^2}{\sigma_y^2}} = \frac{S_x^2 / \sigma_x^2}{S_y^2 / \sigma_y^2} \rightarrow F_{n_x - 1, n_y - 1}$$