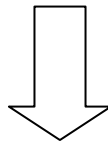


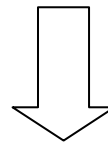
Tema \*: **Introducción a la Teoría de la Estimación**

## Introducción

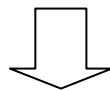
Sea  $X$  la variable aleatoria poblacional con distribución de probabilidad  $f_{\theta}(x)$ ,  
donde  $\theta \in \Theta$  es el parámetro poblacional desconocido



Objetivo: **Obtener valores para  $\theta$  a partir de una muestra aleatoria simple seleccionada de la v.a.  $X$**

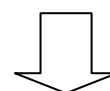


### Procedimientos de Estimación



#### Estimación puntual

Se obtienen valores numéricos para el valor del parámetro



#### Estimación por intervalos

Se construyen intervalos que contengan al parámetro con cierta seguridad

# 1 **Estimación puntual**

Sea la v.a. poblacional  $X$  con distribución de probabilidad  $f_{\theta}(x)$ , donde  $\theta$  es el parámetro poblacional desconocido.

Para hacer inferencia sobre el posible valor de  $\theta \Rightarrow$

{ Seleccionamos una m.a.s. de  $X$ ,  $X_1, X_2, \dots, X_n$ ,  
y definimos un **estimador para  $\theta$** ,  $T(X_1, X_2, \dots, X_n)$ , utilizando la información muestral.  
El valor concreto que toma el estimador para una realización muestral concreta,  
 $\hat{\theta} = T(x_1, x_2, \dots, x_n)$ , se le llama **estimación puntual de  $\theta$**  y  
es una solución particular a nuestro problema.

Notar que un **Estimador Puntual del parámetro  $\theta$**  es un estadístico que sirva para hacer inferencia sobre el parámetro  $\theta$ .

Dado que es posible definir muchos estadísticos que sirvan para hacer inferencias sobre el valor del parámetro  $\theta$ , el objetivo es seleccionar aquellas expresiones que proporcionen garantías en el proceso de estimación paramétrica. Para ello vamos a formular propiedades deseables para la determinación de un buen estimador.

## 1.1 **Características deseables de los estimadores**

Sea  $T(X_1, X_2, \dots, X_n)$  un estimador del parámetro  $\theta$ .

- Diremos que el estimador es **insesgado** si

$$E [T(X_1, X_2, \dots, X_n)] = \theta \Leftrightarrow \text{El estimador es } \mathbf{exacto}$$

Mide la *exactitud* del estimador ya que representa la concentración de las estimaciones entorno al verdadero valor del parámetro.

Cuando  $E [T(X_1, X_2, \dots, X_n)] \neq \theta$ , se dice que el estimador es **sesgado**.

- Diremos que un estimador insesgado es **consistente** si

$$Var [T(X_1, X_2, \dots, X_n)] \xrightarrow{n \rightarrow \infty} 0 \Leftrightarrow \text{El estimador es } \mathbf{preciso}$$

Mide la *precisión* del estimador ya que representa la concentración de las estimaciones entorno al valor medio del estimador.

NOTA: Veamos si los estadísticos media muestral, varianza muestral y proporción muestral verifican las anteriores propiedades deseables:

- Para una variable aleatoria poblacional  $X$  con  $E(X) = \mu$  y  $Var(X) = \sigma^2$ , definimos el estadístico **MEDIA MUESTRAL**:

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

que se emplea para hacer inferencias sobre la media poblacional  $\mu$ .

Se tiene que:

$$\begin{aligned} E(\bar{X}) &= \mu \\ Var(\bar{X}) &= \frac{\sigma^2}{n} \\ \Rightarrow \bar{X} &\text{ es un } \mathbf{\textit{estimador insesgado y consistente}} \text{ para } \mu \end{aligned}$$

- Cuando tenemos una realización muestral, la varianza muestral es una medida de la variabilidad de los datos entorno a la media muestral.

Para una variable aleatoria poblacional  $X$  con  $E(X) = \mu$  y  $Var(X) = \sigma^2$ , definimos el estadístico **VARIANZA MUESTRAL**

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

que se emplea para hacer inferencias sobre la varianza poblacional  $\sigma^2$ .

Además, hemos visto que para obtener el comportamiento aleatorio de  $\bar{X}$  (ya sea aplicando la aditividad del modelo normal o vía el Teorema Central del Límite) es necesario que  $\sigma^2$  sea conocida. Sin embargo, en la mayoría de los casos prácticos, la varianza  $\sigma^2$  es desconocida, por lo que es imprescindible estimarla de los datos incluso si únicamente queremos hacer inferencias sobre  $\mu$ .

Se puede comprobar que:

$$\begin{aligned} E(S^2) &= \sigma^2 \\ Var(S^2) &= \frac{2\sigma^4}{n-1} \\ \Rightarrow S^2 &\text{ es un } \mathbf{\textit{estimador insesgado y consistente}} \text{ para } \sigma^2 \end{aligned}$$

- Para una variable aleatoria poblacional  $X$  dicotómica,  $X \sim b(p)$ , definimos el estadístico **PROPORCIÓN MUESTRAL**:

$$\hat{p} = \frac{\sum_{i=1}^n X_i}{n}$$

que se emplea para hacer inferencias sobre la proporción de unidades que en la población verifican la propiedad de interés,  $p$ .

Se tiene que:

$$E(\hat{p}) = p$$

$$Var(\hat{p}) = \frac{p(1-p)}{n}$$

$\Rightarrow \hat{p}$  es un **estimador insesgado y consistente** para  $p$

## 2 **Estimación por intervalos**

En la estimación puntual, calculamos el valor del estimador del parámetro para una muestra concreta e inferimos ese valor a la población. El problema fundamental radica en que, aunque el estimador cumpla las propiedades deseables, no sabemos si la estimación puntual obtenida con él para una muestra concreta estará o no próxima al valor del parámetro ya que éste es desconocido antes y después del muestreo. Una desventaja de la estimación puntual es que no podemos establecer el error cometido al estimar el parámetro, ni la fiabilidad de nuestra estimación. En otras palabras, no queremos limitarnos a dar un valor concreto para aproximar un parámetro sino proponer una medida del error que pensamos cometer. Para ello, vamos a proporcionar un intervalo que contenga el valor del parámetro.

**Objetivo:** Basándose en la información contenida en las observaciones muestrales se busca un intervalo de valores dentro del cual se tiene cierta “seguridad” de que se encuentre el valor del parámetro. A este intervalo se le denomina **intervalo de confianza para el parámetro  $\theta$** .

### Definición:

Un intervalo aleatorio (intervalo en el que sus dos extremos son variables aleatorias)

$$I(X_1, X_2, \dots, X_n) = [I_1(X_1, X_2, \dots, X_n), I_2(X_1, X_2, \dots, X_n)]$$

se dice que es un **intervalo de confianza para  $\theta$  al nivel de confianza  $1 - \alpha$** , con  $\alpha \in (0, 1)$ , si se cumple que

$$P(\theta \in [I_1(X_1, X_2, \dots, X_n), I_2(X_1, X_2, \dots, X_n)]) = 1 - \alpha$$

## 2.1 Método para construir intervalos de confianza

1. Encontrar una variable aleatoria  $U = h(T, \theta)$ , función del parámetro y del estadístico  $T$ , cuya distribución de probabilidad sea independiente de  $\theta$  y de cualquier otro parámetro desconocido.

2. Seleccionar valores  $U_1(\alpha)$  y  $U_2(\alpha)$  en la distribución de probabilidad de v.a.  $U$  tal que

$$P(U_1(\alpha) \leq U \leq U_2(\alpha)) = 1 - \alpha$$

siendo  $\alpha$  un número  $\in (0, 1)$  y  $1 - \alpha$  el nivel de confianza pedido.

3. Despejar adecuadamente el parámetro  $\theta$  en el intervalo anterior de manera que obtengamos:

$$P(g_1(T, \alpha) \leq \theta \leq g_2(T, \alpha)) = 1 - \alpha$$

### 2.1.1 Intervalo de confianza nivel $1 - \alpha$ para la media $\mu$ de una población normal, $X \sim N(\mu, \sigma)$ , con $\sigma^2$ conocida

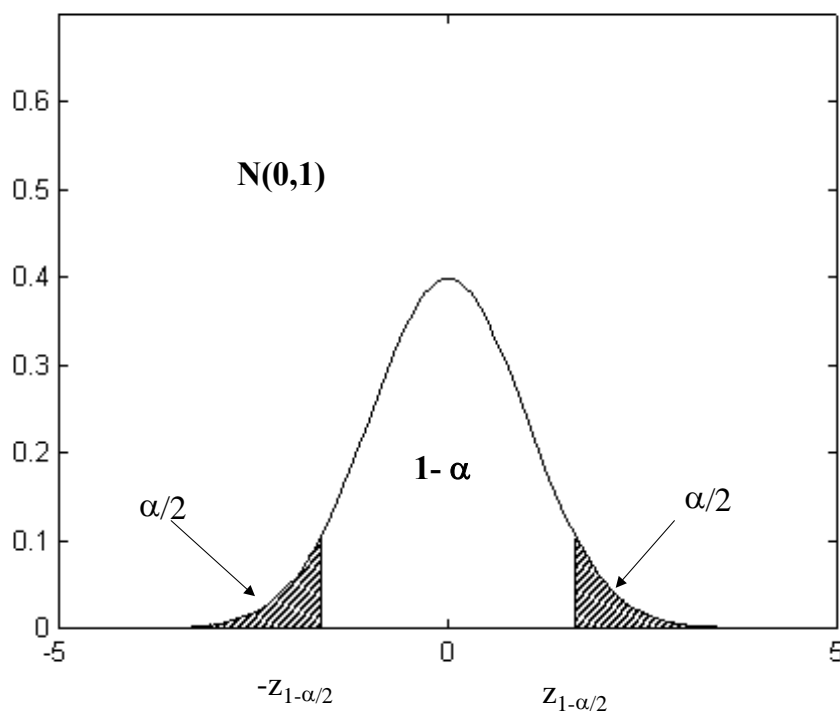
Sea  $(X_1, X_2, \dots, X_n)$  una m.a.s. de una v.a.  $X \sim N(\mu, \sigma)$

- Para hacer inferencias sobre  $\mu$  consideramos la media muestral  $\bar{X} \sim N(\mu, \frac{\sigma}{\sqrt{n}})$ , que al tipificarla tenemos que la variable aleatoria  $U \equiv Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$

- El intervalo será tal que

$$P\left[U_1(\alpha) \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq U_2(\alpha)\right] = 1 - \alpha$$

- Consideramos los valores  $U_1(\alpha) = -z_{1-\frac{\alpha}{2}}$  y  $U_2(\alpha) = z_{1-\frac{\alpha}{2}}$ , ya que en la distribución normal se verifica que:



- Entonces,

$$P(-z_{1-\frac{\alpha}{2}} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{1-\frac{\alpha}{2}}) = 1 - \alpha$$

- Operando adecuadamente obtenemos

$$P(\bar{X} - z_{(1-\frac{\alpha}{2})} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{(1-\frac{\alpha}{2})} \cdot \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$$

- Deducimos que un intervalo de confianza al nivel  $(1 - \alpha)$  para la media poblacional  $\mu$  es

$$\left[ \bar{X} - z_{(1-\frac{\alpha}{2})} \cdot \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{(1-\frac{\alpha}{2})} \cdot \frac{\sigma}{\sqrt{n}} \right] = \bar{X} \pm \underbrace{z_{(1-\frac{\alpha}{2})} \cdot \frac{\sigma}{\sqrt{n}}}_{\text{MARGEN de ERROR}}$$

Para cada realización muestral  $(x_1, x_2, \dots, x_n)$  se obtiene un intervalo de confianza distinto para  $\mu$  al nivel de confianza  $1 - \alpha$ :

$$\left[ \bar{x} - z_{(1-\frac{\alpha}{2})} \cdot \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{(1-\frac{\alpha}{2})} \cdot \frac{\sigma}{\sqrt{n}} \right]$$

**Nota:**

Cuando la v.a. poblacional  $X$  no sea necesariamente normal con  $E(X) = \mu$  y  $Var(X) = \sigma^2$ , el intervalo de confianza obtenido sigue siendo válido para muestras grandes vía el Teorema Central del Límite, ya que de esta manera podemos asegurar que  $\bar{X}$  se distribuye aproximadamente según una  $N(\mu, \frac{\sigma}{\sqrt{n}})$  cuando  $n$  es grande.

### 2.1.2 Determinación del tamaño muestral cuando $\sigma^2$ es conocida

En algún momento en el trabajo estadístico tendremos que decidir qué **tamaño muestral será seleccionado** de la población. Esto es, para estimar la media poblacional  $\mu$  me planteo el número de observaciones de la muestra son necesarias para garantizar, con una confianza dada que el margen de error sea menor que una cantidad prefijada.

El problema planteado es:

- Dado el nivel de confianza  $100(1 - \alpha)\%$  y fijado *el margen de error máximo que permito cometer*  $\equiv err$ .
- Me planteo el valor de  $n$  para que:

$$z_{(1-\frac{\alpha}{2})} \cdot \frac{\sigma}{\sqrt{n}} \leq err$$

- Despejando obtendremos el valor de  $n$ .

**2.1.3 Intervalo de confianza nivel  $1 - \alpha$  para la media  $\mu$  de una población normal,  $X \sim N(\mu, \sigma)$ , con  $\sigma^2$  desconocida**

Sea  $(X_1, X_2, \dots, X_n)$  una m.a.s. de una v.a.  $X \sim N(\mu, \sigma)$

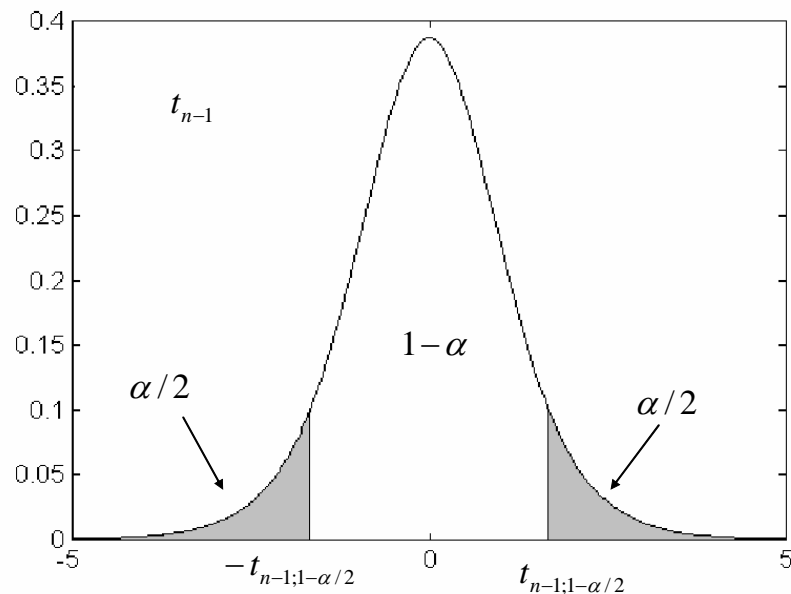
- Consideremos la variable aleatoria  $U \equiv T = \frac{\bar{X} - \mu}{S/\sqrt{n}}$

- Su distribución de probabilidad es  $U \equiv T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$

- El intervalo será tal que

$$P \left[ U_1(\alpha) \leq \frac{\bar{X} - \mu}{S/\sqrt{n}} \leq U_2(\alpha) \right] = 1 - \alpha$$

- Consideramos los valores  $U_1(\alpha) = -t_{n-1;1-\frac{\alpha}{2}}$  y  $U_2(\alpha) = t_{n-1;1-\frac{\alpha}{2}}$ , ya que en la distribución t de Student se verifica que:



- Entonces,

$$P(-t_{n-1;1-\frac{\alpha}{2}} \leq \frac{\bar{X} - \mu}{S/\sqrt{n}} \leq t_{n-1;1-\frac{\alpha}{2}}) = 1 - \alpha$$

- Operando adecuadamente obtenemos

$$P\left(\bar{X} - t_{n-1;1-\frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{n-1;1-\frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}}\right) = 1 - \alpha$$

- Deducimos que un intervalo de confianza al nivel  $(1 - \alpha)$  para la media poblacional  $\mu$  es:

$$\left[ \bar{X} - t_{n-1;1-\frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}}, \bar{X} + t_{n-1;1-\frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}} \right] = \bar{X} \mp \underbrace{t_{n-1;1-\frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}}}_{\text{MARGEN de ERROR}}$$

Para cada realización muestral  $(x_1, x_2, \dots, x_n)$  se obtiene un intervalo de confianza distinto para  $\mu$  al nivel de confianza  $1 - \alpha$ :

$$\left[ \bar{x} - t_{n-1;1-\frac{\alpha}{2}} \cdot \frac{s}{\sqrt{n}}, \bar{x} + t_{n-1;1-\frac{\alpha}{2}} \cdot \frac{s}{\sqrt{n}} \right]$$

Interpretación frecuentista de un intervalo de confianza a nivel  $1 - \alpha$ : Cada vez que tomemos una muestra de tamaño  $n$  tendremos que la media poblacional  $\mu$  se encontrará dentro de dicho intervalo con una confianza del  $100(1 - \alpha)\%$  y sólo un  $100\alpha\%$  de las veces, el intervalo propuesto por este método no recogerá al verdadero valor de  $\mu$ .

**Nota: Factores que influyen en la amplitud de un intervalo de confianza:**

- El nivel de confianza  $1 - \alpha$ : A mayor nivel de confianza mayor amplitud del intervalo
- El tamaño muestral  $n$ : A mayor tamaño muestral menor amplitud del intervalo
- Desviación típica de la población  $\sigma$  o estimación de la desviación típica poblacional,  $S$ : A mayor  $\sigma$  o  $S$  mayor amplitud del intervalo.



## 2.2 Caso de dos poblaciones independientes

Generalicemos lo visto anteriormente al caso de dos poblaciones.

Supongamos que tenemos dos poblaciones **independientes**

$$\begin{cases} X_1 \sim N(\mu_1, \sigma_1) \\ \text{y} \\ X_2 \sim N(\mu_2, \sigma_2) \end{cases}$$

y queremos **construir un estadístico para hacer inferencias sobre la diferencia de medias poblacionales**,  $\mu_1 - \mu_2$ :

- Sea  $(X_{11}, X_{12}, \dots, X_{1n_1})$  una m.a.s. de tamaño  $n_1$  de la primera población  $\Rightarrow$

$$\bar{X}_1 \sim N\left(\mu_1, \frac{\sigma_1}{\sqrt{n_1}}\right)$$

- Sea  $(X_{21}, X_{22}, \dots, X_{2n_2})$  una m.a.s. de tamaño  $n_2$  de la segunda población  $\Rightarrow$

$$\bar{X}_2 \sim N\left(\mu_2, \frac{\sigma_2}{\sqrt{n_2}}\right)$$

- Definimos el estadístico **DIFERENCIA DE MEDIAS MUESTRALES**:

$$\boxed{\bar{X}_1 - \bar{X}_2}$$

Para conocer la distribución muestral de este estadístico vamos a distinguir tres casos:

### 2.2.1 Caso de que las varianzas $\sigma_1^2$ y $\sigma_2^2$ sean conocidas

De manera análoga al caso de una población se tiene que:

$$\bar{X}_1 - \bar{X}_2 \sim N\left(\mu_1 - \mu_2, \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}\right) \Rightarrow Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1)$$

Por tanto, un intervalo de confianza al nivel  $(1 - \alpha)$  para la diferencia de medias poblacionales  $\mu_1 - \mu_2$  es:

$$\begin{aligned} & \left[ \bar{X}_1 - \bar{X}_2 - z_{(1-\frac{\alpha}{2})} \cdot \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, \bar{X}_1 - \bar{X}_2 + z_{(1-\frac{\alpha}{2})} \cdot \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right] = \\ & = \bar{X}_1 - \bar{X}_2 \mp \underbrace{z_{(1-\frac{\alpha}{2})} \cdot \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}_{\text{MARGEN de ERROR}} \end{aligned}$$

### 2.2.2 Caso de que las varianzas $\sigma_1^2$ y $\sigma_2^2$ sean desconocidas no iguales

Para estimar  $\mu_1 - \mu_2$  utilizamos el estadístico:

$$T = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}$$

cuya distribución de probabilidad exacta no es conocida y se utilizan aproximaciones asintóticas. Nosotros vamos a utilizar una aproximación sencilla y conservadora dada por

$$T = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}} \approx t_k$$

donde  $k = \min\{n_1 - 1, n_2 - 1\}$ , aunque la mayoría de los programas de ordenador implementan una expresión para  $k$  más sofisticada.

Notar que  $S_1^2$  es la varianza muestral obtenida de la muestra seleccionada de la primera población y  $S_2^2$  es la varianza muestral obtenida de la muestra extraída de la segunda población.

Por tanto, un intervalo de confianza al nivel  $(1 - \alpha)$  para la diferencia de medias poblacionales  $\mu_1 - \mu_2$  es:

$$\begin{aligned} & \left[ \bar{X}_1 - \bar{X}_2 - t_{k-1; 1-\frac{\alpha}{2}} \cdot \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}, \bar{X}_1 - \bar{X}_2 + t_{k-1; 1-\frac{\alpha}{2}} \cdot \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}} \right] = \\ & = \bar{X}_1 - \bar{X}_2 \mp \underbrace{t_{k-1; 1-\frac{\alpha}{2}} \cdot \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}_{\text{MARGEN de ERROR}} \end{aligned}$$

### 2.2.3 Caso de que las varianzas $\sigma_1^2$ y $\sigma_2^2$ sean desconocidas pero iguales, $\sigma_1^2 = \sigma_2^2 = \sigma^2$

La varianza común  $\sigma^2$  la estimamos mediante

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$

El estadístico utilizado para estimar  $\mu_1 - \mu_2$  es:

$$T = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1+n_2-2}$$

Por tanto, un intervalo de confianza al nivel  $(1 - \alpha)$  para la diferencia de medias poblacionales  $\mu_1 - \mu_2$  es:

$$\begin{aligned} & \left[ \bar{X}_1 - \bar{X}_2 - t_{n_1+n_2-2; 1-\frac{\alpha}{2}} \cdot \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}, \bar{X}_1 - \bar{X}_2 + t_{n_1+n_2-2; 1-\frac{\alpha}{2}} \cdot \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}} \right] = \\ & = \bar{X}_1 - \bar{X}_2 \mp \underbrace{t_{n_1+n_2-2; 1-\frac{\alpha}{2}} \cdot \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}_{\text{MARGEN de ERROR}} \end{aligned}$$