

Lección 9: Contrastes de hipótesis

1. Introducción

El contraste de hipótesis consiste en formular una hipótesis acerca de una población, y aceptar o rechazar esa hipótesis basándonos en la información proporcionada por una muestra.

En el contraste de hipótesis, se denomina **hipótesis nula** (H_0) a la hipótesis planteada y que vamos a comprobar usando una muestra. También se plantea la denominada **hipótesis alternativa** (H_1).

El método de contraste de hipótesis trata de reunir suficientes evidencias como para comprobar que la hipótesis nula es muy poco probable. Si conseguimos reunir estas evidencias, se dice que *rechazamos la hipótesis nula*, y por tanto aceptamos la hipótesis alternativa. Por el contrario, si no logramos reunir suficientes evidencias y la hipótesis nula es plausible, *no se rechaza la hipótesis nula*. Es importante entender que cuando no se rechaza la hipótesis nula, no implica que H_0 sea cierta; simplemente no hemos logrado demostrar que es falsa.

1.1. Tipos de contrastes

Vamos a distinguir dos tipos de contrastes:

- Contrastes paramétricos.
- Contrastes no paramétricos.

Realizamos un **contraste paramétrico** cuando la población sobre la que se establece la hipótesis pertenece a una familia de distribuciones conocida (como por ejemplo, la distribución Normal).

En este tipo de contrastes, normalmente las hipótesis H_0 y H_1 se refieren a algún parámetro de la población (por ejemplo, μ en una distribución Normal). Un caso típico consiste en comprobar si la media es igual a un valor supuesto, o es menor que ese valor. Es decir:

$$H_0 : \mu = 100$$

$$H_1 : \mu < 100$$

La hipótesis se realiza sobre los parámetros de la población, puesto que son valores desconocidos. Solo a través de las muestras podemos estimar o formular hipótesis acerca de esos parámetros poblacionales.

Por el contrario, realizamos un **contraste no paramétrico** cuando no podemos decir cuál es la distribución de probabilidad de la población de la que extraemos la muestra. En estos casos, solo podemos formular hipótesis acerca del tipo de distribución de la variable aleatoria que estamos estudiando. Por ejemplo, podemos formular la hipótesis de que el número de goles por jornada en las ligas europeas es una distribución normal, y comprobar mediante un contraste si es una hipótesis plausible o no.

En este tema solo trataremos los **contrastos paramétricos**. En el siguiente tema veremos algunos casos de **contrastos no paramétricos**.

1.2. Tipos de hipótesis

En un contraste siempre tendremos que decidir una hipótesis nula H_0 y una hipótesis alternativa H_1 . Estas hipótesis pueden ser de dos tipos:

- **Hipótesis simple**

Una hipótesis es simple si plantea que el parámetro de la población tiene un valor concreto. Por ejemplo,

$$H_0 : \mu = 100$$

es una hipótesis simple.

- **Hipótesis compuesta**

Una hipótesis es compuesta cuando plantea que el parámetro de la población puede tomar más de un valor. Por ejemplo,

$$H_1 : \mu < 100$$

es una hipótesis compuesta

En un contraste de hipótesis se puede dar cualquier combinación entre H_0 y H_1 . Es decir, podemos tener las siguientes combinaciones:

H_0	H_1
Simple	Simple
Simple	Compuesta
Compuesta	Simple
Compuesta	Compuesta

De todas formas, los casos más habituales son:

- H_0 simple y H_1 compuesta.
- H_0 compuesta y H_1 compuesta.

El primer caso (H_0 simple, H_1 compuesta) es además mucho más habitual, y los ejemplos posteriores se referirán casi siempre a este caso.

2. Errores en un contraste de hipótesis

Vamos a ilustrar los posibles errores que se pueden cometer en un contraste usando un ejemplo.

Tenemos dos urnas con la siguiente composición:

- Urna A: 4 bolas blancas y 2 bolas negras.
- Urna B: 2 bolas blancas y 4 bolas negras.

Elegimos al azar una urna y extraemos tres bolas, y vamos a realizar el siguiente contraste de hipótesis:

$$H_0 : \text{La urna elegida ha sido la A}$$
$$H_1 : \text{La urna elegida ha sido la B}$$

El resultado del contraste puede ser que rechazamos H_0 , y por tanto asumimos que las bolas salieron de la urna B. Por el contrario, si no logramos reunir evidencias suficientes, no podremos rechazar H_0 y tendremos que asumir que las bolas salieron de la urna A.

Imaginemos que sacamos las tres bolas de la urna A, pero nuestra muestra contiene 2 bolas negras y una bola blanca. Si enseñamos la muestra a otra persona que no sepa de dónde han salido las bolas, esa persona nos dirá que es más probable que hayan salido de B, y por tanto rechazará H_0 . Cuando la hipótesis nula es cierta, pero la rechazamos en el contraste, estamos cometiendo un **error de tipo I**.

Repitamos ahora el muestreo. Elegimos la urna B, sacamos tres bolas, y resultan ser 2 bolas blancas y 1 bola negra. Ahora enseñamos esa muestra a otra persona que no sabe de dónde han salido las bolas. Esa persona nos dirá que lo más probable es que vengan de A. Es decir, aceptará la hipótesis nula H_0 , por ser más probable que H_1 . Cuando aceptamos la hipótesis nula, pero la hipótesis nula es falsa, estamos cometiendo un **error de tipo II**.

Este tipo de errores son inherentes al contraste de hipótesis. En ocasiones, un tipo de error puede ser mucho más costoso que otro. Veamos algunos ejemplos:

- Afirmar que un paciente tiene cáncer cuando no lo tiene.
En este caso, si la hipótesis nula es que no tiene cáncer, y nuestro test rechaza H_0 , estamos cometiendo un error de tipo I (rechazamos H_0 , le decimos que tiene cáncer, pero resulta que H_0 es cierta).
- Afirmar que un paciente no tiene cáncer, pero sí lo tiene.
En este caso, hemos aceptado H_0 y hemos rechazado H_1 , pero H_0 es falsa; es decir, hemos cometido un error de tipo II.

¿Cuál es más costoso? En este caso, parece claro que el error de tipo II puede tener consecuencias más graves que el error de tipo I.

Veamos otro ejemplo:

- Suponemos que un puente viejo aguantará, pero acaba derrumbándose.
En este caso, nuestra H_0 era que el puente es seguro, la hemos aceptado, pero ha resultado ser falsa. Hemos cometido un error de tipo II.
- Suponemos que un puente viejo no aguantará, pero no se derrumba.
En este caso, hemos establecido H_0 como que el puente aguantará, y la rechazamos. Pero H_0 ha resultado ser cierta. Hemos cometido un error de tipo I.

De nuevo, en este caso el error de tipo II puede ser mucho más grave que el error de tipo I.

Otro ejemplo más:

- En un juicio, declaramos inocente al acusado, pero era culpable.
En este caso, H_0 es que el acusado es inocente. Aceptamos H_0 pero resulta ser falsa. Es decir, es un error de tipo II.
- En un juicio, declaramos culpable al acusado, pero era inocente.
En este caso, de nuevo H_0 es que el acusado es inocente, pero se rechaza, y se condena al acusado. Sin embargo, H_0 ha resultado ser cierta. Si rechazamos H_0 pero es cierta, cometemos un error de tipo I.

En este caso, el error de tipo I es más costoso que el error de tipo II. De hecho, el sistema legal funciona de este modo, intenta evitar condenar a inocentes. Esto es, el objetivo es evitar cometer error de tipo I (es decir, condenar a inocentes), a expensas de cometer errores de tipo II (dejar escapar a criminales cuya culpabilidad no se puede demostrar).

Por tanto, existen elementos externos al contraste de hipótesis que pueden hacer que tanto los errores de tipo I como de tipo II sean muy graves.

La probabilidad de cometer un error de tipo I se conoce como **nivel de significación** (α). La probabilidad de cometer un error de tipo II se representa como β , siendo $1 - \beta$ la **potencia del contraste**.

En resumen, existen cuatro posibles opciones en un contraste, resumidas en la siguiente tabla:

		Decisión	
		No rechazar H_0	Rechazar H_0
H_0	Cierta	Acierto	Error Tipo I
	Falsa	Error Tipo II	Acierto

3. Contraste de hipótesis

Como regla general, para realizar el contraste de hipótesis elegiremos un nivel de significación α . Es decir, elegiremos cuál es la probabilidad de cometer un error de tipo I. Cuanto más pequeño sea α , más improbable será cometer un error de tipo I.

Una vez fijado el nivel de significación, el método del contraste de hipótesis intenta maximizar la potencia del contraste $1 - \beta$, por lo que minimiza la probabilidad del error de tipo II β .

El método del contraste de hipótesis consta de los siguientes pasos:

1. Elegimos un estadístico $T(X_1, X_2, \dots, X_n, \theta)$ que dependa de la muestra y del parámetro de la población sobre el que vayamos a realizar el contraste
2. Dividimos el espacio de posibles valores de T en dos regiones R_0 (región de aceptación) y R_1 (región crítica), usando la distribución de probabilidad de la población, y el nivel de significación α .
3. Calculamos el valor de T para la muestra obtenida. Lo denominaremos T_0 .
4. Si $T_0 \in R_0$, entonces no rechazamos la hipótesis nula. Por el contrario, si $T_0 \in R_1$, entonces rechazamos la hipótesis nula.

En estos pasos hemos obviado dos cuestiones que son imprescindibles para realizar el contraste de hipótesis:

- ¿Cómo elegimos el estadístico $T(X_1, \dots, X_n, \theta)$?
- ¿Cómo calculamos cuál es la región de aceptación R_0 y la región crítica R_1 ?

Para elegir el estadístico adecuado, no existe una respuesta general. En este tema vamos a estudiar algunos tipos de contrastes paramétricos, para los que ya existen estadísticos “prefijados” que podemos emplear.

Respecto a la segunda pregunta, la región crítica es la que hace que el valor T_0 sea improbable. Por este motivo, la región crítica estará siempre en las colas de la distribución, ocupando un área de α dentro de la función de densidad. En algunos casos estará solo en una de las colas, y en otros estará en las dos colas (ocupando $\frac{\alpha}{2}$ en cada cola).

En la figura 1 se observa la función de densidad de una distribución t de Student con 25 grados de libertad. Se han coloreado las dos zonas que forman la región crítica. Se ha calculado el valor t_0 que hace que

$$P(t_n > t_0) = \frac{0,05}{2} = 0,025$$

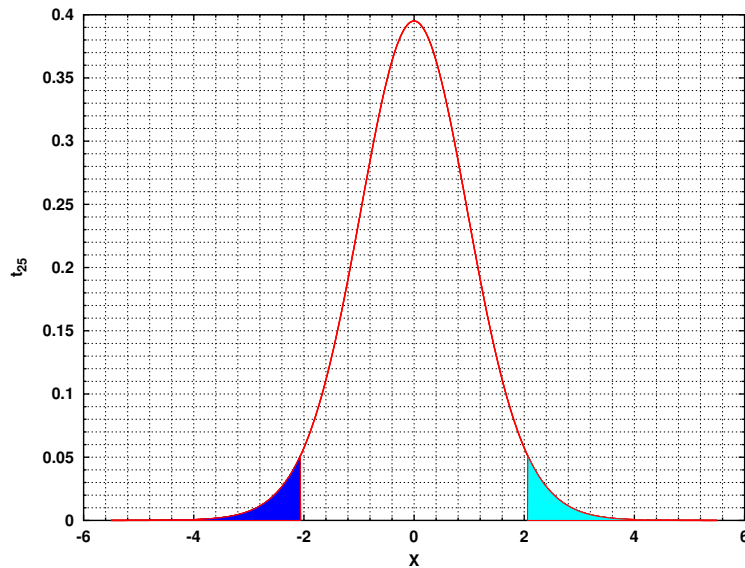


Figura 1: Región crítica formada por dos colas, de $\frac{\alpha}{2}$ cada una.

En este caso, como la función de densidad es simétrica, el valor para la cola de los negativos será $-t_0$. Es decir, el valor t_0 divide el espacio de posibles valores del estadístico T en dos regiones:

- La región crítica
En este caso $(-\infty, -t_0] \cup [t_0, +\infty)$. Si el valor del estadístico cae dentro de esta zona, rechazaríamos la hipótesis nula H_0 (y por tanto aceptaríamos la hipótesis alternativa H_1).
- La región de aceptación
En este caso $(-t_0, t_0)$. Si el valor del estadístico cae dentro de esta zona, aceptaríamos la hipótesis nula H_0 .

Ejemplo 1 *Un agricultor quiere vender su cosecha de frutas. Para calcular su valor aproximado, consulta a un experto que tasa la producción en 30 kg. de fruta por árbol.*

Para mayor seguridad decide tomar una muestra de 26 árboles al azar, obteniendo un rendimiento medio de 34 kg. de fruta por árbol, con una desviación típica de 4 kg. ¿Puede el agricultor aceptar la hipótesis del experto? (use un nivel de significación del 5%).

Solución

En este caso, tomaremos como H_0 que la media de producción por árbol es 30 kg. ¿Cuál será la hipótesis alternativa? En algunos casos, plantearemos la posibilidad de que la media sea superior a 30, o inferior. Pero en este caso, simplemente diremos que la media no

es 30, sin especificar si es mayor o menor. Por tanto, el contraste quedaría definido con las siguientes hipótesis nula y alternativa:

$$H_0 : \mu = 30$$

$$H_1 : \mu \neq 30$$

Para el contraste hemos obtenido una muestra de tamaño $n = 26$, en la que hemos obtenido $\bar{x} = 34$, y $s = 4$.

Ya tenemos todo lo necesario para empezar el contraste. Sigamos los pasos definidos en los párrafos anteriores:

1. Elegir el estadístico

En este caso, estamos realizando un muestreo para tomar decisiones sobre la media de la población, a partir de la media muestral. Como ya sabemos

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

Sin embargo, desconocemos el valor de σ . Podríamos intentar aproximar σ por s . Pero si realizamos esa aproximación, la variable aleatoria que tendríamos que usar (lo veremos en las próximas secciones) sería

$$\frac{\bar{x} - \mu}{s/\sqrt{n-1}} \in t_{n-1}$$

Es decir, tenemos que usar la distribución t de Student, en vez de la Normal.

2. Calcular las regiones de aceptación R_0 y crítica R_1

Para este contraste, aceptamos H_0 cuando el valor de \bar{x} esté próximo μ . El valor de la media muestral se podría alejar de μ bien con valores muy grandes por encima de μ , bien con valores muy pequeños por debajo de μ .

Por tanto, la región crítica estará presente en las dos colas de la distribución t_{n-1} . En cada cola, ocuparemos un área de $\frac{\alpha}{2}$ para calcular las regiones críticas.

Como $n = 26$ y $\alpha = 0,05$, tenemos que calcular t_0 tal que

$$P(t_{25} > t_0) = \frac{0,05}{2} = 0,025$$

La región crítica estará formada por $(-\infty, t_0] \cup [t_0, +\infty)$

En las tablas de la distribución t de Student, podemos comprobar que $t_0 = 2,060$. La región crítica está formada por las dos áreas coloreadas que se muestran en la figura 1.

3. Calcular el valor de T_0

Como el estadístico era

$$\frac{\bar{x} - \mu}{s/\sqrt{n-1}}$$

tenemos que

$$T_0 = \frac{\bar{x} - 30}{s/\sqrt{n-1}} = \frac{34 - 30}{4/\sqrt{25}} = 5$$

4. Aceptar o rechazar H_0

Como hemos obtenido que $T_0 > t_0$, la muestra se encuentra en la región crítica, por lo que rechazamos H_0 , y concluimos que $\mu \neq 30$. La hipótesis del experto es falsa.

□

4. Contrastes en poblaciones normales

4.1. Contraste sobre μ con σ conocida

En este caso, el contraste es del tipo

$$H_0 : \mu = \mu_0$$

$$H_1 : \mu \neq \mu_0$$

El estadístico que usaremos para este contraste es

$$\frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \in N(0, 1)$$

Por el tipo de hipótesis H_1 , es un contraste a dos colas, donde cada cola ocupa $\frac{\alpha}{2}$. Tendremos que calcular el límite de la región crítica $\lambda_{\alpha/2}$ usando la tabla de la distribución Normal

$$P(N(0, 1) > \lambda_{\alpha/2}) = \frac{\alpha}{2}$$

Por simetría, la región crítica será $(-\infty, -\lambda_{\alpha/2}] \cup [+ \lambda_{\alpha/2}, +\infty)$

En algunas ocasiones, incluso si σ es desconocida pero la muestra es muy grande, podemos aproximar σ por s y usar este mismo estadístico.

Ejemplo 2 De una población Normal con varianza 100, se desea contrastar la hipótesis de que la media poblacional es 40. Para ello, extraemos una muestra aleatoria de tamaño 25, resultando $\bar{x} = 41$. Si el nivel de significación es 10%, ¿es aceptable la hipótesis?

Solución

El contraste de hipótesis es:

$$H_0 : \mu = 40$$

$$H_1 : \mu \neq 40$$

con $\alpha = 0,1$

El estadístico para el contraste es

$$T(X_1, \dots, X_n, \mu) = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \in N(0, 1)$$

Para calcular la región crítica, buscamos en una tabla de la distribución Normal

$$P(Z > \lambda_{\alpha/2}) = 1 - \frac{0,1}{2} = 0,95$$

Resultando $\lambda_{\alpha/2} = 1,645$

Calculamos

$$T_0 = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{41 - 40}{10/\sqrt{25}} = 0,5$$

Como $T_0 \in (-\lambda_{\alpha/2}, +\lambda_{\alpha/2})$ no podemos rechazar la hipótesis nula.

□

4.2. Contraste sobre μ con σ desconocida

El contraste a realizar es:

$$H_0 : \mu = \mu_0$$

$$H_1 : \mu \neq \mu_0$$

El estadístico a usar es

$$\frac{\bar{x} - \mu}{s/\sqrt{n-1}} \in t_{n-1}$$

Para calcular la región crítica, buscamos t_0 en la tabla de la distribución t de Student, tal que

$$P(t_{n-1} > t_0) = \frac{\alpha}{2}$$

Por simetría, la región crítica será $(-\infty, -t_0] \cup [t_0, +\infty)$

Ejemplo 3 *El cava Puerta del Sol tiene un contenido medio de azúcar de 24 gr/l. Para contrastar dicha hipótesis al 1% se toma una muestra de 21 botellas, resultando $\bar{x} = 25$ gr/l y $s^2 = 4$. Suponiendo que la población es normal, ¿es aceptable la hipótesis del fabricante?*

Solución

El contraste a realizar es el siguiente:

$$H_0 : \mu = 24$$

$$H_1 : \mu \neq 24$$

con $\alpha = 0,01$

Como σ es desconocida, tendremos que usar es

$$\frac{\bar{x} - \mu}{s/\sqrt{n-1}} \in t_{n-1}$$

Buscamos t_0 en la tala de t_{21-1}

$$P(t_{20} > t_0) = \frac{0,01}{2} = 0,005$$

Resultando $t_0 = 2,845$

Calculamos

$$T_0 = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{25 - 24}{2/\sqrt{20}} = 2,234$$

Como $T_0 \in (-t_0, +t_0)$ no podemos rechazar la hipótesis nula al 1% de significación.

□

4.3. Contraste sobre σ^2 con μ desconocida

El contraste a realizar es

$$H_0 : \sigma^2 = \sigma_0^2$$

$$H_1 : \sigma^2 \neq \sigma_0^2$$

El estadístico para este contraste es

$$T = \frac{ns^2}{\sigma^2} \in \chi_{n-1}^2$$

Para calcular la región crítica a dos colas debemos tener en cuenta que la distribución χ_{n-1}^2 no es simétrica. Tendremos que calcular los dos extremos de la región crítica como sigue:

$$P(\chi_{n-1}^2 > \lambda_{\alpha/2}) = \frac{\alpha}{2}$$

y

$$P(\chi_{n-1}^2 < \lambda_{1-\alpha/2}) = \frac{\alpha}{2}$$

La región crítica estará formada por $(-\infty, \lambda_{1-\alpha/2}] \cup [\lambda_{\alpha/2}, +\infty)$

Ejemplo 4 *La varianza de una muestra de 17 valores de una población normal es $s^2 = 9$. ¿Se puede aceptar al 5% de significación que la varianza de la población es $\sigma^2 = 16$? ¿Y si la muestra tuviera 314 valores?*

Solución

El contraste a realizar es

$$H_0 : \sigma^2 = 16$$

$$H_1 : \sigma^2 \neq 16$$

El estadístico es

$$T = \frac{ns^2}{\sigma^2} \in \chi_{n-1}^2$$

Calculamos la región crítica para $\alpha = 0,05$. Por un lado tenemos

$$P(\chi_{17-1}^2 > \lambda_{\alpha/2}) = \frac{0,05}{2} = 0,025$$

Resultando $\lambda_{\alpha/2} = 28,85$

Para la otra cola calculamos

$$P(\chi_{17-1}^2 < \lambda_{1-\alpha/2}) = \frac{0,05}{2} = 0,025$$

Resultando $\lambda_{1-\alpha/2} = 6,91$

Por tanto la región crítica es $(-\infty, 6,91] \cup [28,85, +\infty)$

El valor del estadístico es

$$T_0 = \frac{ns^2}{\sigma^2} = \frac{17 \cdot 9}{16} = 9,57$$

Por tanto, dado que $9,57 \in (6,91, 28,85)$ no podemos rechazar la hipótesis nula, no hemos reunido evidencias suficientes para rechazar H_0 .

Ahora repetimos el contraste de hipótesis usando una muestra de tamaño $n = 314$, siendo el resto de valores iguales.

En este caso, tenemos que recalcular la región crítica. Por un lado tenemos

$$P(\chi_{314-1}^2 > \lambda_{\alpha/2}) = \frac{0,05}{2} = 0,025$$

Resultando $\lambda_{\alpha/2} = 364,98$

Para la otra cola calculamos

$$P(\chi_{314-1}^2 < \lambda_{1-\alpha/2}) = \frac{0,05}{2} = 0,025$$

Resultando $\lambda_{1-\alpha/2} = 266,80$

Por tanto la región crítica es $(-\infty, 266,80] \cup [364,98, +\infty)$

Recalculamos también

$$T_0 = \frac{ns^2}{\sigma^2} = \frac{314 \cdot 9}{16} = 176,63$$

Por tanto, en este caso T_0 no está en el intervalo $(266,80, 364,98)$, por lo que rechazamos la hipótesis nula.

¿Por qué en el primer caso no rechazamos H_0 y en el segundo caso sí? En el primer caso no hemos logrado reunir evidencias suficientes para rechazar H_0 . Pero eso no significa que tengamos que aceptar H_0 , simplemente que no tenemos evidencias para rechazar H_0 .

En el segundo caso, al incrementar el tamaño de la muestra, ya tenemos evidencias suficientes como para rechazar H_0 y aceptar H_1 .

Este ejemplo ilustra una propiedad muy importante de los contrastes de hipótesis, que a menudo se malinterpreta. Que no logremos rechazar H_0 no implica que H_0 sea verdadera. Simplemente indica que no tenemos evidencias para rechazar H_0 . Si estuviéramos realizando un estudio de Ingeniería, sería arriesgado tomar H_0 como cierta en un contraste que no logra rechazar H_0 . Si queremos comprobar si una hipótesis es cierta, tenemos que plantear un contraste de manera que esa hipótesis sea la hipótesis alternativa H_1 , e intentar falsificar H_0 para mostrar que H_1 es verdadera.

□

4.4. Contraste sobre σ^2 con μ conocida

De nuevo, el contraste a realizar es

$$\begin{aligned} H_0 : \sigma^2 &= \sigma_0^2 \\ H_1 : \sigma^2 &\neq \sigma_0^2 \end{aligned}$$

Para este contraste usamos el estadístico

$$\frac{n\hat{\sigma}^2}{\sigma^2} \in \chi_{n-1}^2$$

donde

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

La región crítica se calcula como en el caso con μ desconocida. La región crítica estará formada por dos colas, teniendo en cuenta que la distribución χ_{n-1}^2 no es simétrica. Tendremos que calcular los dos extremos de la región crítica como sigue:

$$P(\chi_{n-1}^2 > \lambda_{\alpha/2}) = \frac{\alpha}{2}$$

y

$$P(\chi_{n-1}^2 < \lambda_{1-\alpha/2}) = \frac{\alpha}{2}$$

La región crítica estará formada por $(-\infty, \lambda_{1-\alpha/2}] \cup [\lambda_{\alpha/2}, +\infty)$

Ejemplo 5 De una población Normal con $\mu = 16$ se desea contrastar la hipótesis de que la varianza poblacional vale 3 con un nivel de significación $\alpha = 0,1$. Para ello se obtiene una muestra aleatoria de tamaño 10. La muestra obtenida es $\{10, 11, 13, 14, 15, 15, 16, 17, 29, 20\}$

Solución

El contraste a realizar es

$$H_0 : \sigma^2 = 3$$

$$H_1 : \sigma^2 \neq 3$$

Para este caso podemos calcular que

$$n\hat{\sigma}^2 = \sum_{i=1}^{10} (x_i - \mu)^2 = 102$$

Por tanto

$$T_0 = \frac{n\hat{\sigma}^2}{\sigma^2} = \frac{102}{3} = 34$$

Calculamos ahora la región crítica. Por un lado tenemos

$$P(\chi_{10-1}^2 > \lambda_{\alpha/2}) = \frac{0,05}{2} = 0,025$$

Resultando $\lambda_{\alpha/2} = 18,31$

Para la otra cola calculamos

$$P(\chi_{10-1}^2 < \lambda_{1-\alpha/2}) = \frac{0,05}{2} = 0,025$$

Resultando $\lambda_{1-\alpha/2} = 3,94$

Por tanto la región crítica es $(-\infty, 3,94] \cup [18,31, +\infty)$

Como $T_0 = 34$ no está en el intervalo $(3,94, 18,31)$, tenemos que rechazar H_0 y concluir que $\sigma^2 \neq 3$.

□

5. Contrastes entre dos poblaciones normales

5.1. Contraste para $\mu_x - \mu_y$ con σ_x y σ_y conocidas

En este caso, el contraste consiste en

$$H_0 : \mu_x - \mu_y = \mu_0$$

$$H_1 : \mu_x - \mu_y \neq \mu_0$$

El estadístico en este caso es

$$\frac{\bar{x} - \bar{y} - (\mu_x - \mu_y)}{\sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}} \in N(0, 1)$$

Para la región crítica, que es a dos colas, tenemos que elegir un valor $\lambda_{\alpha/2}$ tal que

$$P(Z > \lambda_{\alpha/2}) = \frac{\alpha}{2}$$

Por simetría, la región crítica será $(-\infty, -\lambda_{\alpha/2}] \cup [+\lambda_{\alpha/2}, +\infty)$

Se puede emplear el mismo estadístico incluso si las varianzas poblacionales son desconocidas, pero la muestra es grande. En este caso, se pueden usar las varianzas muestrales como aproximación de las varianzas poblacionales sin cambiar el estadístico.

Ejemplo 6 *El consumo de cigarrillos entre los estudiantes de Filosofía e Ingeniería se ha extendido en los últimos años. Se pretende contrastar al 4% si los estudiantes de Filosofía fuman la misma cantidad diaria de cigarrillos que los de Ingeniería. Se ha tomado una muestra de 400 estudiantes de Filosofía, y resulta que suman de media 16 cigarrillos diarios. Se tomó también una muestra de 200 estudiantes de Ingeniería, resultando una media de*

15 cigarrillos diarios. Si el consumo diario de cigarrillos es una variable aleatoria Normal, y las desviaciones típicas son 40 para los estudiantes de Filosofía, y 20 para los estudiantes de Ingeniería. ¿Es el consumo de cigarrillos entre los estudiantes de Filosofía e Ingeniería igual?

Solución

Tomaremos las siguientes variables aleatorias

$X =$ núm. de cigarrillos diarios que fuma un estudiante de Filosofía

$Y =$ núm. de cigarrillos diarios que fuma un estudiante de Ingeniería

El contraste a realizar es

$$H_0 : \mu_x - \mu_y = 0$$

$$H_1 : \mu_x - \mu_y \neq 0$$

con $\alpha = 0,04$

Los datos de la población que conocemos son $\sigma_x = 40$ y $\sigma_y = 20$

Por su lado, los datos de las muestras que conocemos son $n_x = 400$, $n_y = 200$, $\bar{x} = 16$ e $\bar{y} = 15$.

El valor del estadístico es

$$T_0 = \frac{\bar{x} - \bar{y} - (\mu_x - \mu_y)}{\sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}} = \frac{16 - 15 - 0}{\sqrt{\frac{40^2}{400} + \frac{20^2}{200}}} = 0,408$$

Para calcular la región crítica, obtenemos el valor $\lambda_{\alpha/2}$ tal que

$$P(Z > \lambda_{\alpha/2}) = \frac{\alpha}{2} = \frac{0,04}{2} = 0,02$$

Resultando $\lambda_{\alpha/2} = 2,054$

La región crítica es por tanto $(-\infty, -2,054] \cup [2,054, +\infty)$

Como $T_0 = 0,408 \in (-2,054, 2,054)$, no podemos rechazar H_0 , y por tanto no hemos logrado encontrar evidencias de que haya una diferencia en el número medio de cigarrillos que fuman los estudiantes de Filosofía y los de Ingeniería.

□

5.2. Contraste sobre $\mu_x - \mu_y$ con varianzas desconocidas pero $\sigma_x = \sigma_y$

En este caso, el contraste consiste en

$$\begin{aligned}H_0 &: \mu_x - \mu_y = \mu_0 \\H_1 &: \mu_x - \mu_y \neq \mu_0\end{aligned}$$

El estadístico en este caso es

$$\frac{\bar{x} - \bar{y} - (\mu_x - \mu_y)}{\sqrt{n_x + n_y} \sqrt{n_x s_x^2 + n_y s_y^2}} \sqrt{n_x n_y} \sqrt{n_x + n_y - 2} \in t_{n_x + n_y - 2}$$

Para calcular la región crítica, buscamos t_0 en la tabla de la distribución t de Student, tal que

$$P(t_{n_x + n_y - 2} > t_0) = \frac{\alpha}{2}$$

Por simetría, la región crítica será $(-\infty, -t_0] \cup [+t_0, +\infty)$

Ejemplo 7 Queremos contrastar si el número medio de estudiantes por grupo en las Escuelas de Ingeniería de Caminos y Aeronáutica son iguales. Para ello hemos tomado una muestra de 14 grupos de Caminos y 18 de Aeronáutica. Para los grupos de Caminos hemos obtenido que

$$\sum_{i=1}^{14} x_i = 140$$

y

$$\sum_{i=1}^{14} (x_i - \bar{x})^2 = 490$$

Para los de Aeronáutica hemos obtenido

$$\sum_{i=1}^{18} y_i = 198$$

y

$$\sum_{i=1}^{18} (y_i - \bar{y})^2 = 486$$

Si suponemos que ambas poblaciones tienen la misma varianza, contraste la hipótesis de que el número medio de estudiantes por grupo es similar en ambas escuelas, con un nivel de significación del 5%.

Solución

El contraste a realizar es

$$H_0 : \mu_x - \mu_y = 0$$

$$H_1 : \mu_x - \mu_y \neq 0$$

Los datos muestrales son $n_x = 14$, $n_y = 18$, $\bar{x} = 140/14 = 10$, $\bar{y} = 198/18 = 11$, $s_x^2 = 490/14 = 35$ y $s_y^2 = 486/18 = 27$.

El estadístico resulta ser

$$\frac{\bar{x} - \bar{y} - (\mu_x - \mu_y)}{\sqrt{n_x + n_y} \sqrt{n_x s_x^2 + n_y s_y^2}} \sqrt{n_x n_y} \sqrt{n_x + n_y - 2} = \frac{10 - 11 - 0}{\sqrt{14 + 18} \sqrt{490 + 486}} \sqrt{14 + 18} \sqrt{14 + 18 - 2}$$

Por tanto, $T_0 = -0,492$

Para calcular la región crítica, buscamos t_0 en la tabla de la distribución t de Student, tal que

$$P(t_{14+18-2} > t_0) = \frac{\alpha}{2} = \frac{0,05}{2} = 0,025$$

Resultando $t_0 = 2,042$.

Por simetría, la región crítica será $(-\infty, -2,042] \cup [2,042, +\infty)$

Por tanto $T_0 = -0,492 \in (-2,042, 2,042)$, por lo que no podemos rechazar la hipótesis nula, es decir, no tenemos evidencia para afirmar que el número medio de alumnos por grupo en ambas escuelas sea diferente.

□

5.3. Contraste sobre $\frac{\sigma_x^2}{\sigma_y^2}$

Si queremos contrastar una hipótesis de la forma $\sigma_x^2 = K\sigma_y^2$, podemos transformar la hipótesis para plantear el contrastes de la siguiente forma:

$$H_0 : \frac{\sigma_x^2}{\sigma_y^2} = K$$

$$H_1 : \frac{\sigma_x^2}{\sigma_y^2} \neq K$$

El estadístico para este contraste es

$$\frac{n_x}{n_x - 1} \frac{n_y - 1}{n_y} \frac{\sigma_y^2 s_x^2}{\sigma_x^2 s_y^2} \in F_{n_x - 1, n_y - 1}$$

Para calcular la región crítica debemos tener en cuenta que la distribución F de Fisher no es simétrica, pero tiene la propiedad

$$F_{\alpha/2}^{a,b} = \frac{1}{F_{1-\alpha/2}^{b,a}} \quad (1)$$

Por tanto, calculamos $F_{\alpha/2}$ tal que

$$P(F_{n_x-1, n_y-1} > F_{\alpha/2}) = \frac{\alpha}{2}$$

Para calcular $F_{1-\alpha/2}$ obtenemos el valor que verifica

$$P(F_{n_y-1, n_x-1} > \frac{1}{F_{1-\alpha/2}}) = \frac{\alpha}{2}$$

La región crítica estará formada entonces por

$$(-\infty, F_{1-\alpha/2}] \cup [F_{\alpha/2}, +\infty)$$

Ejemplo 8 Para contrastar la hipótesis de que las varianzas de dos poblaciones normales son iguales, se toman dos muestras independientes de tamaños 31 y 21. Las varianzas muestrales resultan ser 400 y 300, respectivamente. Contraste la hipótesis al 2% de significación.

Solución

El contraste que nos pide el enunciado del problema es

$$\begin{aligned} H_0 : \sigma_x^2 &= \sigma_y^2 \\ H_1 : \sigma_x^2 &\neq \sigma_y^2 \end{aligned}$$

que podemos transformar en el siguiente contraste

$$\begin{aligned} H_0 : \frac{\sigma_x^2}{\sigma_y^2} &= 1 \\ H_1 : \frac{\sigma_x^2}{\sigma_y^2} &\neq 1 \end{aligned}$$

Los datos muestrales son $n_x = 31$, $n_y = 21$, $s_x^2 = 400$ y $s_y^2 = 300$.

El valor del estadístico es

$$T_0 = \frac{n_x}{n_x - 1} \frac{n_y - 1}{n_y} \frac{\sigma_y^2 s_x^2}{\sigma_x^2 s_y^2} = \frac{31}{31 - 1} \frac{21 - 1}{21} \cdot 1 \cdot \frac{400}{300} = 1,312$$

Para calcular la región crítica, obtenemos el valor $F_{\alpha/2}$ tal que

$$P(F_{30,20} > F_{\alpha/2}) = \frac{\alpha}{2} = \frac{0,02}{2}$$

El valor resulta ser $F_{\alpha/2} = 2,779$.

Por otro lado, calculamos

$$P(F_{20,30} > \frac{1}{F_{1-\alpha/2}}) = \frac{\alpha}{2} = \frac{0,02}{2}$$

Obteniendo

$$\frac{1}{F_{1-\alpha/2}} = 2,549$$

Por tanto $F_{1-\alpha/2} = \frac{1}{2,549} = 0,392$.

Nótese que para calcular $F_{1-\alpha/2}$ hemos usado la distribución F de parámetros $(20, 30)$, en vez de $(30, 20)$. Hemos aprovechado la propiedad de la distribución F que dice que $F_{\alpha/2}^{a,b} F_{1-\alpha/2}^{b,a} = 1$, tal y como se muestra en la ecuación (1).

La región crítica es pues $(-\infty, 0,392] \cup [2,779, +\infty)$

Como el valor $T_0 = 1,312 \in (0,392, 2,779)$, no podemos rechazar H_0 , y no podemos decir que las varianzas de ambas poblaciones sean diferentes.

□

6. Contrastes sobre el coeficiente de correlación de Pearson

Cuando ajustamos un modelo por mínimos cuadrados, podemos contrastar hipótesis acerca del coeficiente de correlación. El contraste más habitual trata de averiguar si podemos decir que el coeficiente de correlación es significativamente diferente de cero.

El contraste consiste en

$$H_0 : \rho = 0$$

$$H_1 : \rho \neq 0$$

donde ρ es el coeficiente de correlación de Pearson poblacional.

Para este contraste, y solo en el caso en el que las muestras x e y hayan sido extraídas de poblaciones normales, podemos usar el siguiente estadístico

$$\sqrt{\frac{r^2}{1-r^2}}(n-2) \in t_{n-2}$$

siendo r el coeficiente de correlación muestral.

Como la distribución t de Student es simétrica, para calcular la región crítica basta con obtener el valor t_0 que hace que

$$P(t_{n-2} > t_0) = \frac{\alpha}{2}$$

La región crítica será $(-\infty, -t_0] \cup [t_0, +\infty)$

Ejemplo 9 *En una granja están probando un nuevo tipo de maíz, que supuestamente hace que las gallinas pongan más huevos. Para contrastar si existe correlación entre la cantidad de maíz y el número de huevos, se toma una muestra formada por 42 gallinas. El coeficiente de correlación muestral ha resultado ser $r = 0,2$. Suponga que las poblaciones son normales, y contraste la hipótesis de que hay correlación al 1%.*

Solución

El contraste es

$$H_0 : \rho = 0$$

$$H_1 : \rho \neq 0$$

El valor del estadístico es

$$T_0 = \sqrt{\frac{r^2}{1-r^2}}(n-2) = \sqrt{\frac{0,2^2}{1-0,2^2}}(42-2) = 1,291$$

Para la región crítica, calculamos

$$P(t_{42-2} > t_0) = \frac{\alpha}{2} = \frac{0,01}{2}$$

Resultando $t_0 = 2,705$

Por tanto, la región crítica es $(-\infty, -2,705] \cup [2,705, +\infty)$

Como $T_0 = 1,291 \in (-2,705, 2,705)$ no podemos rechazar H_0 , y por tanto no podemos decir que exista correlación entre ambas variables. No podemos afirmar que exista una relación entre la cantidad de maíz y el número de huevos que ponen las gallinas.

□

7. Contrastes sobre proporciones

En este tipo de contrastes, se plantean hipótesis acerca del parámetro p de una variable de Bernoulli o de una variable binomial. Veremos dos tipos de contrastes: p toma un valor concreto, y acerca de la diferencia entre proporciones.

7.1. Contraste sobre la proporción de una característica

En este contraste planteamos las hipótesis:

$$H_0 : p = p_0$$

$$H_1 : p \neq p_0$$

Para este contraste podemos usar el estadístico

$$\frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \in N(0, 1)$$

donde p es el valor poblacional y \hat{p} es el valor de la proporción muestral.

Para calcular la región crítica, tendremos que encontrar un valor $\lambda_{\alpha/2}$ tal que

$$P(Z > \lambda_{\alpha/2}) = \frac{\alpha}{2}$$

Por simetría, la región crítica estará formada por $(-\infty, -\lambda_{\alpha/2}] \cup [+\lambda_{\alpha/2}, +\infty)$

Ejemplo 10 *En la Escuela de Caminos se hace un sondeo a 1600 alumnos y profesores. 960 de ellos se confiesan colchoneros, mientras que el resto dice ser madridista. A un nivel de significación del 5%, ¿podemos decir que el 65% de los profesores y alumnos son colchoneros?*

Solución

Tal y como está planteado el enunciado, podemos calcular que la proporción de colchoneros en la muestra es

$$\hat{p} = \frac{960}{1600} = 0,6$$

El contraste consistirá de

$$H_0 : p = 0,65$$

$$H_1 : p \neq 0,65$$

El valor del estadístico es

$$T_0 = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{0,6 - 0,65}{\sqrt{\frac{0,65(1-0,65)}{1600}}} = -4,19$$

Para calcular la región crítica, hallamos $\lambda_{\alpha/2}$ tal que

$$P(Z > \lambda_{\alpha/2}) = \frac{\alpha}{2} = \frac{0,05}{2} = 1,96$$

Por lo que la región crítica resulta ser $(-\infty, -1,96] \cup [1,96, +\infty)$

Por tanto, como T_0 no está en $(-1,96, 1,96)$, tenemos que rechazar la hipótesis nula, y podemos decir que la proporción de colchoneros en la Escuela no es 0,65.

□

7.2. Contraste sobre la diferencia de proporciones

En este caso queremos realizar un contraste como el siguiente:

$$H_0 : p_x - p_y = p_0$$

$$H_1 : p_x - p_y \neq p_0$$

El estadístico para este tipo de contrastes es

$$\frac{\hat{p}_x - \hat{p}_y - p_0}{\sqrt{\frac{\hat{p}_x(1-\hat{p}_x)}{n_x} + \frac{\hat{p}_y(1-\hat{p}_y)}{n_y}}} \in N(0, 1)$$

siempre y cuando n_x y n_y sean valores grandes.

Los parámetros \hat{p}_x y \hat{p}_y se refieren a las proporciones muestrales, mientras que p_0 es el valor de la diferencia que se va a contrastar.

De nuevo, para calcular la región crítica, tendremos que encontrar un valor $\lambda_{\alpha/2}$ tal que

$$P(Z > \lambda_{\alpha/2}) = \frac{\alpha}{2}$$

Por simetría, la región crítica estará formada por $(-\infty, -\lambda_{\alpha/2}] \cup [+\lambda_{\alpha/2}, +\infty)$

Ejemplo 11 *Para cuantificar la influencia de una campaña publicitaria se toman dos muestras independientes antes y después de la campaña, resultando que:*

	<i>Antes</i>	<i>Después</i>
<i>Conocedores del producto</i>	145	282
<i>No conocedores</i>	355	318

Al 5% de significación, ¿se puede aceptar que el incremento de conocedores del producto ha sido del 25%?

Solución

El contraste a realizar es:

$$H_0 : p_x - p_y = 0,25$$

$$H_1 : p_x - p_y \neq 0,25$$

donde X e Y se refieren a los encuestados después y antes de la campaña.

Antes de la campaña, se encuestó a $n_y = 145 + 355 = 500$ personas. Después de la campaña, el número de encuestados fue $n_x = 282 + 318 = 600$.

Por tanto, las proporciones muestrales son $\hat{p}_x = \frac{282}{600} = 0,47$, después de la campaña. Por el contrario, antes de la campaña la proporción de conocedores era $\hat{p}_y = \frac{145}{500} = 0,29$.

El valor del estadístico es por tanto

$$T_0 = \frac{\hat{p}_x - \hat{p}_y - p_0}{\sqrt{\frac{\hat{p}_x(1-\hat{p}_x)}{n_x} + \frac{\hat{p}_y(1-\hat{p}_y)}{n_y}}} = \frac{0,47 - 0,29 - 0,25}{\sqrt{\frac{0,47(1-0,47)}{600} + \frac{0,29(1-0,29)}{500}}} = -2,43$$

Por su lado, la región crítica está definida por $\lambda_{\alpha/2}$, que se obtiene al calcular

$$P(Z > \lambda_{\alpha/2}) = \frac{\alpha}{2} = \frac{0,05}{2}$$

Resultando $\lambda_{\alpha/2} = 1,96$.

Por tanto la región crítica es $(-\infty, -1,96] \cup [1,96, +\infty)$

Como $T_0 = -2,43$ no está en $(-1,96, 1,96)$ rechazamos la hipótesis nula, y por tanto el incremento en la proporción de conocedores no ha sido del 25%.

□

8. Contrastes con una sola cola

Los contrastes a una sola cola pueden tener alguna de las formas siguientes

$$H_0 : \theta \geq \theta_0$$

$$H_1 : \theta < \theta_0$$

o bien

$$\begin{aligned}H_0 &: \theta \leq \theta_0 \\H_1 &: \theta > \theta_0\end{aligned}$$

En el caso general, no es posible obtener estadísticos para este tipo de contrastes, ya que no existen estadísticos para hipótesis compuestas. Por tanto, es necesario transformar la hipótesis nula a una hipótesis simple, mientras que se conserva el sentido del contraste.

Por ejemplo, si estamos contrastando la media de una población normal con las siguientes hipótesis:

$$\begin{aligned}H_0 &: \mu \leq 5 \\H_1 &: \mu > 5\end{aligned}$$

Podemos plantear el siguiente contraste, que nos proporcionará los mismos resultados

$$\begin{aligned}H_0 &: \mu = 5 \\H_1 &: \mu > 5\end{aligned}$$

ya que si la muestra aconseja rechazar la hipótesis nula con $\mu = 5$, también aconsejará rechazar H_0 con $\mu \leq 5$.

En el caso general, usando una hipótesis simple en la hipótesis nula, podemos dejar la hipótesis alternativa como una hipótesis compuesta. La precaución que habrá que tener es que la región crítica tendrá una única cola. En el ejemplo anterior, la región crítica estará en la cola para valores más altos. Por el contrario, para una hipótesis alternativa del tipo $\mu < 5$, la región crítica estaría en la zona de los valores más bajos de la distribución.

Normalmente el sentido de la hipótesis alternativa se plantea según el problema que tengamos que resolver. Por ejemplo, tenemos un producto cuyo fabricante afirma que solo contiene defectos en el 0,5% de los casos. En este contraste, la afirmación del fabricante sería la hipótesis nula, y la alternativa es que la proporción de defectos sea mayor. No tiene sentido contrastar que la proporción de defectos sea menor.

Por tanto, el contraste sería

$$\begin{aligned}H_0 &: p = 0,005 \\H_1 &: p > 0,005\end{aligned}$$

En general, usaremos los mismo estadísticos que hemos visto hasta ahora, con las mismas distribuciones de probabilidad para cada estadístico, pero tomando una única cola en la que repartir el nivel de significación α . La cola estará en la zona de los valores más bajos o más altos, dependiendo de cómo planteemos la hipótesis alternativa H_1 .

Ejemplo 12 *Un fabricante comercializa tablets, y afirma que como mucho un 5% de los tablets contiene grietas en la pantalla. Un comerciante al que el fabricante provee de dispositivos electrónicos ha acumulado 500 tablets y ha comprobado que 28 de ellos tenían grietas en la pantalla. A un 1% de nivel de significación, ¿es correcta la especificación de defectuosos del fabricante?*

Solución

En este caso, contrastaremos si la proporción de defectuosos es mayor a lo que el fabricante reporta. Será un contraste a una única cola, usando la especificación del fabricante como hipótesis nula. Por tanto, el contraste será:

$$H_0 : p = 0,05$$

$$H_1 : p > 0,05$$

Según la muestra obtenida por el comerciante, sabemos que

$$\hat{p} = \frac{28}{500} = 0,056$$

Recordamos que el estadístico para contrastes con proporciones es

$$T_0 = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{0,056 - 0,05}{\sqrt{\frac{0,05(1-0,05)}{500}}} = 0,616$$

En cuanto a la región crítica, como el estadístico sigue una distribución $N(0, 1)$, tendremos que seleccionar solo los valores altos de la distribución, usando una única cola para la región crítica. La frontera entre la región de aceptación y la crítica, λ_α , se calculará como

$$P(Z > \lambda_\alpha) = \alpha = 0,01$$

En este caso usamos α en vez de $\frac{\alpha}{2}$, ya que todo el nivel de significación se distribuye en una única cola.

El valor que resulta es $\lambda_\alpha = 2,326$

Por tanto, la región crítica para el contraste a una única cola es $[2,326, +\infty)$

En este caso, como $T_0 = 0,616 < 2,326$, no podemos rechazar H_0 , y a pesar de las evidencias recogidas por el comerciante, no podemos afirmar que la especificación del fabricante no sea correcta.

□