

- A set of navigation icons typically found in Beamer presentations, including symbols for back, forward, search, and other slide controls.

Comparando vector de medias de dos poblaciones

La estadística T^2 es apropiada para comparar las respuestas de un conjunto de experimentos con ciertas especificaciones (población 1) con las respuestas de otro conjunto de experimentos con distintas especificaciones (población 2).

considere una muestra aleatoria de dimensión n_1 proveniente de la población 1, y una muestra de dimensión n_2 de la población 2. Las observaciones sobre p variables se pueden arreglar de la siguiente manera:

| población 1 | estadística suficiente | |
|---|---|---|
| $\mathbf{x}_{11}, \dots, \mathbf{x}_{1n_1}$ | $\bar{\mathbf{x}}_1 = \frac{1}{n_1} \sum_{j=1}^{n_1} \mathbf{x}_{1j}$ | $\mathbf{S}_1 = \frac{1}{n_1-1} \sum_{j=1}^{n_1} (\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)(\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)'$ |
| $\mathbf{x}_{21}, \dots, \mathbf{x}_{2n_2}$ | $\bar{\mathbf{x}}_2 = \frac{1}{n_2} \sum_{j=1}^{n_2} \mathbf{x}_{2j}$ | $\mathbf{S}_2 = \frac{1}{n_2-1} \sum_{j=1}^{n_2} (\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)(\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)'$ |

Nos interesa realizar inferencia sobre:

(vector promedi de población 1) - (vector promedio de población 2) = $\mu_1 - \mu_2$

Por ejemplo, nos gustaría responder las preguntas:

¿Es $\mu_1 = \mu_2$?

Si $\mu_1 - \mu_2 \neq 0$, ¿cuáles componentes son diferentes?

- Cuando $\mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \mathbf{\Sigma}$, $\sum_{j=1}^{n_1} (\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)(\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)'$ es una estimación de $(n_1 - 1)\mathbf{\Sigma}$, y $\sum_{j=1}^{n_2} (\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)(\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)'$ es una estimación de $(n_2 - 1)\mathbf{\Sigma}$.
- Por lo que podemos agrupar la información de ambas muestras con el proposito de estimar la matriz de covarianza común $\mathbf{\Sigma}$.
- $$\mathbf{S}_{pooled} = \frac{\sum_{j=1}^{n_1} (\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)(\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)' + \sum_{j=1}^{n_2} (\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)(\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)'}{n_1 + n_2 - 2}$$
$$= \frac{n_1 - 1}{n_1 + n_2 - 2} \mathbf{S}_1 + \frac{n_2 - 1}{n_1 + n_2 - 2} \mathbf{S}_2$$
- Puesto que $\sum_{j=1}^{n_1} (\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)(\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)'$ tiene $n_1 - 1$ d.f., y $\sum_{j=1}^{n_2} (\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)(\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)'$ tiene $n_2 - 1$ d.f., el divisor de la ec. anterior se obtiene al combinar los grados de libertad de cada componente.

Para probar la hipótesis $\mu_1 - \mu_2 = \delta_0$ (arbitrario), consideramos la distancia de $\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2$ a δ_0 .

Sabemos que $E(\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2) = \mu_1 - \mu_2$, además por la suposición de independencia se tiene que

$$\text{Cov}(\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2) = \left(\frac{1}{n_1} + \frac{1}{n_2} \right) \mathbf{\Sigma} \quad (1)$$

Ahora, dado que \mathbf{S}_{pooled} estima $\mathbf{\Sigma}$, vemos que $\left(\frac{1}{n_1} + \frac{1}{n_2} \right) \mathbf{S}_{pooled}$

es un estimador de $\text{Cov}(\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2)$

Por último, prueba de la razón de verosimilitud de $H_0 : \mu_1 - \mu_2 = \delta_0$, se basa en el cuadrado de la distancia estadística T^2 , y esta dada por

$$T^2 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2 - \delta_0)' \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) \mathbf{S}_{pooled} \right]^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2 - \delta_0) > c^2, \quad (2)$$

donde el valor crítico de c^2 se determina mediante la distribución de dos muestras del estadístico T^2

Distribución de dos muestras del estadístico T^2

Resultado

Si $\mathbf{X}_{11}, \mathbf{X}_{12}, \dots, \mathbf{X}_{1n_1}$ es una muestra aleatoria de dimensión n_1 proveniente de $N_p(\mu_1, \Sigma)$, y $\mathbf{X}_{21}, \mathbf{X}_{22}, \dots, \mathbf{X}_{2n_2}$ es una muestra aleatoria de dimensión n_2 proveniente de $N_p(\mu_2, \Sigma)$, entonces

$$T^2 = [\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2 - (\mu_1 - \mu_2)]' \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) \mathbf{S}_{pooled} \right]^{-1} [\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2 - (\mu_1 - \mu_2)] \quad (3)$$

se distribuye como $\frac{(n_1+n_2-2)p}{(n_1+n_2-p-1)} F_{p, n_1+n_2-p-1}$

Por lo que

$$P[\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2 - (\mu_1 - \mu_2)]' \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) \mathbf{S}_{pooled} \right]^{-1} [\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2 - (\mu_1 - \mu_2) \leq c^2] = 1 - \alpha, \text{ donde}$$

$$c^2 = \frac{(n_1 + n_2 - 2)p}{(n_1 + n_2 - p - 1)} F_{p, n_1 + n_2 - p - 1}(\alpha) \quad (4)$$

Distribución de dos muestras del estadístico T^2

- Nos interesa principalmente la región de confianza para $\mu_1 - \mu_2$.
- Del resultado anterior podemos concluir que todas las $\mu_1 - \mu_2$ dentro de una distancia estadística c^2 de $\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2$ constituye una región de confianza.
- Esta región es un elipsoide centrado en la diferencia muestral $\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2$, y cuyos ejes son determinados por los eigenvalores y eigenvectores de \mathbf{S}_{pooled}

Comparación de dos muestras cuando $\Sigma_1 \neq \Sigma_2$

- Cuando $\Sigma_1 \neq \Sigma_2$ no es posible encontrar una medida de distancia del tipo T^2 , cuya distribución no dependa de Σ_1 y Σ_2
- La magnitud de las discrepancias (respecto a la suposición de normalidad) que son críticas en la situación multivariada dependen probablemente, en gran medida, sobre el número de variables p .
- No obstante, para n_1 y n_2 grandes, se evita la complejidad debido a matrices de covarianza distintas.

Observación

Si, $n_1 = n_2 = n$, entonces $(n - 1)/(n + n - 2) = 1/2$, así

$$\frac{1}{n_1} \mathbf{S}_1 + \frac{1}{n_2} \mathbf{S}_2 = \frac{1}{n} (\mathbf{S}_1 + \mathbf{S}_2) = \frac{(n-1)\mathbf{S}_1 + (n-1)\mathbf{S}_2}{n+n-2} \left(\frac{1}{n} + \frac{1}{n} \right) =$$

$$\mathbf{S}_{pooled} \left(\frac{1}{n} + \frac{1}{n} \right)$$

con muestras iguales, el procedimiento para muestras grandes es esencialmente el mismo al procedimiento basado en la matriz de varianza agregada (*pooled*).

Ejemplo 3: Procedimiento para muestras grandes para inferencias acerca de diferencia en medias

Nos interesa analizar los datos discutidos en el ejemplo 2 utilizando la aproximación para muestras grandes (sobreponga los resultados).

- Encontrar la combinación lineal más crítica que nos lleva al rechazo de $H_0 : \mu_1 - \mu_2 = 0$. Interprete y grafique sus resultados.

Residuos

Con un poco de álgebra se obtiene de los residuos la expresión:

$$\sum_{l=1}^g \sum_{j=1}^{n_l} (\mathbf{x}_{lj} - \bar{\mathbf{x}})(\mathbf{x}_{lj} - \bar{\mathbf{x}})' = \sum_{l=1}^g (\bar{\mathbf{x}}_l - \bar{\mathbf{x}})(\bar{\mathbf{x}}_l - \bar{\mathbf{x}})' + \sum_{l=1}^g \sum_{j=1}^{n_l} (\mathbf{x}_{lj} - \bar{\mathbf{x}}_l)(\mathbf{x}_{lj} - \bar{\mathbf{x}}_l)', \quad (12)$$

donde, el primer termino se refiere a la *suma total corregida de los cuadrados y productos cruzados*, el segundo término a la *suma de los cuadrados y productos cruzados entre tratamientos*, y el tercer término a la *suma de cuadrados y productos cruzados entre residuos*.

La *suma de cuadrados y productos cruzados entre residuos* se puede expresar como

$$\mathbf{W} = \sum_{l=1}^g \sum_{j=1}^{n_l} (\mathbf{x}_{lj} - \bar{\mathbf{x}} - l)(\mathbf{x}_{lj} - \bar{\mathbf{x}} - l)' = (n_1 - 1)\mathbf{S}_1 + (n_2 - 1)\mathbf{S}_2 + \cdots + (n_g - 1)\mathbf{S}_g, \quad (13)$$

donde \mathbf{S}_l es la matriz de covarianza muestral para la muestra l .

Esta matriz es una generalización de $(n_1 + n_2 - 2)\mathbf{S}_{pooled}$

De manera analoga al caso univariado, la hipótesis de no efectos de tratamiento

$$H_0 : \tau_1 = \tau_2 = \cdots = \tau_g = \mathbf{0}. \quad (14)$$

se prueba al considerar la razones relativas entre la suma de tratamientas y la suma de residuos.

De manera equivalente se puede considerar la razón entre la suma de residuos y la suma total.

Lambda de Wilks

Un test para $H_0 : \tau_1 = \dots = \tau_g = \mathbf{0}$ toma en cuenta varianzas generalizadas.

Se rechaza H_0 si la razón de varianzas generalizadas es pequeña:

$$\lambda^* = \frac{|\mathbf{W}|}{|\mathbf{B} + \mathbf{W}|} = \frac{|\sum_{l=1}^g \sum_{j=1}^{n_l} (\bar{\mathbf{x}}_{lj} - \bar{\mathbf{x}}_l)(\bar{\mathbf{x}}_{lj} - \bar{\mathbf{x}}_l)'|}{|\sum_{l=1}^g \sum_{j=1}^{n_l} (\bar{\mathbf{x}}_{lj} - \bar{\mathbf{x}})(\bar{\mathbf{x}}_{lj} - \bar{\mathbf{x}})'|}, \quad (15)$$

Esta cantidad se le conoce como lambda de Wilks, la cual está relacionada con el criterio de la razón de verosimilitud. Además, se puede expresar en términos de los eigenvalores $\hat{\lambda}_1, \dots, \hat{\lambda}_s$ de $\mathbf{W}^{-1}\mathbf{B}$ como:

$$\lambda^* = \prod_{i=1}^s \left(\frac{1}{1 + \hat{\lambda}_i} \right), \quad (16)$$

, donde $s = \min(p, g - 1)$, el rango de \mathbf{B} .

Otros estadísticos para testear la igualdad de múltiples medias multivariadas, como son el *estadístico de Pillai*, el estadístico de *Lawley-Hotelling*, y el *estadístico de Roy*, también pueden ser escritos como funciones particulares de los eigenvalores de $\mathbf{W}^{-1}\mathbf{B}$.

Distribución de λ^*

La distribución exacta de λ^* se puede derivar para casos especiales:

| p | g | distribución muestral para datos normales multivariados |
|------------|------------|---|
| $p = 1$ | $g \geq 2$ | $\left(\frac{\sum n_l - g}{g-1} \right) \left(\frac{1-\lambda^*}{\lambda^*} \right) \sim F_{g-1, \sum n_l - g}$ |
| $p = 2$ | $g \geq 2$ | $\left(\frac{\sum n_l - g - 1}{g-1} \right) \left(\frac{1-\sqrt{\lambda^*}}{\sqrt{\lambda^*}} \right) \sim F_{2(g-1), 2(\sum n_l - g - 1)}$ |
| $p \geq 1$ | $g = 2$ | $\left(\frac{\sum n_l - p - 1}{p} \right) \left(\frac{1-\lambda^*}{\lambda^*} \right) \sim F_{p, \sum n_l - p - 1}$ |
| $p \geq 1$ | $g = 3$ | $\left(\frac{\sum n_l - p - 2}{p} \right) \left(\frac{1-\sqrt{\lambda^*}}{\sqrt{\lambda^*}} \right) \sim F_{2p, 2(\sum n_l - p - 2)}$ |

Distribución de λ^* para muestras grandes.

Para muestras grandes Bartlett (1938) propuso una modificación de λ^* para testear H_0 . Específicamente, ha demostrado que si H_0 es verdad y $\sum n_l = n$ es grande,

$$-\left(n-1-\frac{(p+g)}{2}\right) \ln \lambda^* = -\left(n-1-\frac{(p+g)}{2}\right) \ln \left(\frac{|\mathbf{W}|}{|\mathbf{B}+\mathbf{W}|}\right) \quad (17)$$

se aproxima a una distribución chi-cuadrada con $p(g - 1)$ grados de libertad.

De esta manera, para $\sum n_l = n$ grande, rechazamos H_0 a nivel de significancia α si

$$-\left(n-1-\frac{(p+g)}{2}\right) \ln \left(\frac{|\mathbf{W}|}{|\mathbf{W}+\mathbf{B}|}\right) > \chi_{p(g-1)}^2(\alpha), \quad (18)$$

donde $chi^2_{p(g-1)}(\alpha)$ es percentil superior (100α) de la distribución chi-cuadrada con $p(g-1)$ d.f.

