

## MANOVA y regresión lineal

Wednesday, March 9, 2022 7:23 AM

### ## MANOVA y regresión lineal ##

Voy a comparar un vec. de medias

$X_{ij} = \mu + \tau_i + \epsilon_{ij}$  (todos vecs.)  $\rightarrow$  obs  $j$  de la pobla<sup>o</sup>  $i$

$\tau_i \rightarrow$  efecto del tratamiento en  $i$ .

$\epsilon_{ij}$  indep  $N_p(0, \Sigma)$

Cada componente de  $X_{ij}$  satisface el modelo univariado.

$$X_{ij} = \underbrace{\bar{x}}_{\mu} + (\underbrace{\bar{x}_i - \bar{x}}_{\tau_i}) + (\underbrace{X_{ij} - \bar{x}_i}_{\epsilon_{ij}})$$

Hipótesis:

$H_0: \tau_1 = \dots = \tau_g = 0$  (no hay efecto en ninguna pobla<sup>o</sup>).

Fuente var<sup>o</sup>

SSP

Tratamiento  $B = \sum_{i=1}^g n_i (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})'$

$g-1$

Wilk's  $\Lambda$   
 $|W|/|T|$

Residual  $W = \sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)(x_{ij} - \bar{x}_i)'$

$n-g$

Total

$T = W + B$

$n-1$

$$\Lambda^* = \frac{|SSCP_e|}{|SSCP_e + SSCP_n|} \quad \Lambda^* = \frac{|W|}{|B+W|}$$

Si  $\Lambda^* > F$ , rechazo  $H_0$ .

Intervalos de C. Simultáneos.

Hay  $p$  vars

Hay  $\frac{p(p-1)}{2}$  diferencias

En Bonferroni:

$$t_{n-g} \left( \frac{\alpha}{2m} \right) \text{ con } m = \frac{p(p-1)}{2}$$

Sean  $n = \sum_{i=1}^g n_i$  con confianza al menos  $1-\alpha$

$\tau_{li} - \tau_{ki}$  pertenece a:

$$\bar{x}_{li} - \bar{x}_{ki} \pm t_{n-g} \left( \frac{\alpha}{p(p-1)} \right) \sqrt{\frac{W_{ii}}{n-g} \left( \frac{1}{n_i} + \frac{1}{n_k} \right)}$$

Con  $\tau_{li}$  el efecto del tratam<sup>l</sup> en la var  $i$

### Regresión lineal $\rightarrow$ métodos de predic<sup>o</sup>

$$Y = \beta_0 + \beta_1 x$$

$$Y = \beta_0 + \alpha_1 \beta_1 + \alpha_2 \beta_2 + \dots \text{ (varias vars.)}$$

Método estadístico para predecir vals de una o varias vars respuesta, dados los vals de una o varias vars predictoras.

Regr lin clásicas:

Sean  $z_1, \dots, z_r$ ,  $r$  vars predictoras.

Sea  $Y$  una var respuesta (presuntam. relacionada)

Ej:  $Y$  val en el mercado

$z_1 = \text{Zona}$ ,  $z_2 = \text{Área}$ ,  $z_3 = \text{Antigüedad} \dots$

Hay q' limpiar los datos

El modelo clásico, supone q'  $Y$  se compone de media q' depende de las  $z_i$  + un error aleatorio  $\epsilon$ .

Obs:

1. Los vals de las vars indep ( $z_i$ ) se toman fijos.

2. La v.a respuesta y el error son v.a

Tenemos:

$$Y = \beta_0 + \beta_1 z_1 + \dots + \beta_r z_r + \epsilon \text{ error.}$$

Si han  $n$  obs indep de  $Y$  asoc. a respectivos vals de  $z_i$

$z_i$ , entonces:

$$y_i = \beta_0 + \beta_1 z_{i1} + \dots + \beta_r z_{ir} + \epsilon_i$$

$$y_n = \beta_0 + \beta_1 z_{n1} + \dots + \beta_r z_{nr} + \epsilon_n$$

$z_{ab}$ : val que toma en la obs a la var  $b$ .

Se asume:

$$1. E[\epsilon_i] = 0$$

$$2. \text{var}(\epsilon_i) = \sigma^2 \rightarrow \text{cte}$$

$$-E[y_i] = \beta_0 + \beta_1 z_{i1} + \dots + \beta_r z_{ir} \quad 3. \text{Cov}(\epsilon_j, \epsilon_k) = 0 \quad \forall j \neq k \text{ (son indep)}$$

Matricial:

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & z_{11} & \dots & z_{1r} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & z_{n1} & \dots & z_{nr} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_r \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix} \rightarrow Y = Z \beta + \epsilon$$

$n \times 1 \quad n \times (r+1) \quad (r+1) \times 1 \quad n \times 1$

$$\bullet E[\epsilon] = 0 \quad \bullet \text{Cov}(\epsilon) = \sigma^2 I \rightarrow \text{las cov son 0.}$$

$\square$  C/ col de  $Z$  son las n obs de las vars predictoras  
 $\square$  C/ fila  $j$  de  $Z$  son los vals de las vars predictoras en la obs  $j$ .

Estimador de mínimos cuadrados

Construir estimaciones de  $\beta$  ( $\beta_0 + \dots + \beta_r$ ) usando  $Z$  y  $Y$  (datos) para hacer predicciones.

Sea  $lb$  vals de  $\beta$ , las difs:

$$y_i - \underbrace{\beta_0 - \beta_1 z_{i1} - \dots - \beta_r z_{ir}}_{\substack{\text{val observado} \\ \text{en } y}} \rightarrow \text{predic}^\circ, lb \leftarrow \text{correcto.}$$

Busco  $lb$  que minimice las sumas cuadráticas.

$$S(lb) = \sum_{i=1}^n (y_i - \beta_0 - \dots - \beta_r z_{ir})^2 = (Y - Z lb)' (Y - Z lb) = \text{escalar}$$

Las entradas de  $lb$  q minimizan  $S(lb)$  se obtienen con mínimos cuadrados (estima<sup>o</sup> de min cuadrados de  $\beta$ ) y se denotan  $\hat{\beta}$ .

obs:

1.  $\hat{\beta}$  son los vals. "más consistentes" con rela<sup>o</sup> lineal de  $Y$  y  $Z$ .

$$2. \hat{\epsilon}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 z_{i1} - \dots - \hat{\beta}_r z_{ir} \quad i = 1, \dots, n$$

$\rightarrow$  residuos.

Teo:

Sea  $Z$  con rango completo  
 El estimador de mínimos cuadrados de  $\beta$  está dado por:

$$\hat{\beta} = (Z'Z)^{-1} Z'Y$$

$$\hat{Y} = Z \hat{\beta} \rightarrow \text{con base en esto, luego estimo el error}$$

$$\hat{Y} = Z \hat{\beta} = H Y \quad \text{con } H = Z (Z'Z)^{-1} Z'$$

$$\hat{\epsilon} = Y - \hat{Y} = Y - Z \hat{\beta} = Y - H Y$$

$$= (I - Z (Z'Z)^{-1} Z') Y = (I - H) Y$$

