

Regresión lineal 2

Wednesday, March 23, 2022 2:11 PM

Regresión lineal 2

Descomposición de suma de cuadrados.

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum \hat{e}_i^2$$

Suma de cuadr. totales Suma de cuadr. de la regre.^o Suma de cuadr. residuales

$$R^2 = 1 - \frac{\sum \hat{e}_i^2}{\sum (y_i - \bar{y})^2} = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2}$$

R^2 me dice q' tanto la variabilidad de una var me explica a la otra.
Entre más significativa la var escogida, más va a mejorar mi R^2 .

$$R^2 = 1 \text{ si Error} = 0 = \hat{e}_i$$

$R^2 = 0$ cuando mi predic^o es exactam. \bar{y} . (no se ve nada)

$$\beta_i = 0$$

Inferencia en el Modelo de regresión

$$E[Y] = \beta_0 + \beta_1 z_1 + \dots + \beta_n z_n$$

Estudiamos la ~ muestral de $\hat{\beta}$ y de $\hat{e}'\hat{e}$ suponiendo E normal.

$$Y = Z\beta + E, \quad Z \text{ rango completo}$$

$$E \sim N(0, \sigma^2 I) : \text{el estimador es max}$$

homocedastidad: varianzas iguales.
errores normales, indep y con homocedastidad.

El estimador de máx verosimilitud de β es $\hat{\beta}$

$$\hat{\beta} = (Z'Z)^{-1} Z'Y \sim N(\beta, \sigma^2 (Z'Z)^{-1})$$

Región confianza $(1-\alpha)$ para β

n : # obs
 r : # var.

$$(\beta - \hat{\beta})' Z' Z (\beta - \hat{\beta}) \leq (r-1) s^2 F_{r-1, n-r-1}(\alpha)$$

$$s^2 = \frac{\sum \hat{e}_i^2}{(n-r-1)}$$

I.C. simultáneo para los β_i

$$\hat{\beta}_i \pm \sqrt{\text{var}(\hat{\beta}_i)} \sqrt{(r-1) F_{r-1, n-r-1}(\alpha)}$$

$\text{var}(\hat{\beta}_i)$ = i-ésimo elem de la diagonal de $s^2 (Z'Z)^{-1}$

I.C. 1 a 1 con confianza $(1-\alpha)$

$$\hat{\beta}_i \pm t_{n-r-1}(\frac{\alpha}{2}) \sqrt{\text{var}(\hat{\beta}_i)}$$

Si el 0 está incluido en el intervalo, β_i puede ser 0 y por tanto, no es tan significativa y podría sacarla de mi modelo.

Región de confianza

Está centrada en $\hat{\beta}$ y la orientación y la distancia dependen de los vcs. y vls. propios de $Z'Z$.

Ejemplo 2 vars.

$$\hat{\beta} = \begin{bmatrix} 30.76 \\ 6.63 \\ 0.045 \end{bmatrix}$$

$\hat{\beta}_1$ es la pendiente de la var 1, mientras la var 2 está constante.

Modelo ANOVA (lm)

$$H_0: \beta_0 = \beta_1 = \dots = \beta_r = 0$$

Hago pruebas de hip. para c/ $\beta_i = 0$

c/	T para H_0	Prob > T
int.	3.93	0.001
β_1	3.53	0.004
β_2	0.16	0.87

β_2 no me ayuda a explicar el comportam.
 β_2 no significativa. C#no! \hat{e}_i
Si hiciera el I.C. el 0 va a estar ahí.

Inferencia sobre parámetros del modelo.

Parte del análisis de reg. incluye inferencias sobre algunos paráms. del mod.

Ej. $\beta_i = 0$ (no influyen esa var.)

Denotamos esos predictores z_{11}, \dots, z_{1r} . $1/r, 2, \dots, (q+1) \times 1$

$$P.H: H_0: \beta_1 = \dots = \beta_r = 0 \Rightarrow \beta(2) = 0 \quad \beta = \begin{pmatrix} \beta(1) \\ \beta(2) \end{pmatrix} \rightarrow (r-q) \times 1$$

$$Z = \begin{pmatrix} Z(1) \\ Z(2) \end{pmatrix} \begin{matrix} n \times (q-1) \\ n \times (r-q) \end{matrix}$$

q+1 pues incluye a β_0

$$\text{El modelo lineal se vuelve: } Y = Z\beta + \varepsilon = (Z(1) \ Z(2)) \begin{pmatrix} \beta(1) \\ \beta(2) \end{pmatrix} + \varepsilon$$

Evaluar H_0 con extra sum of squares (ESS)

$$0 \leq ESS = SS_{res}(Z(1)) - SS_{res}(Z)$$

suma de cuad. residuales del mod solo con las vars de $Z(1)$

$$\hookrightarrow \sum \hat{\varepsilon}_i^2 Z(1) - \sum \hat{\varepsilon}_i^2$$

$$\Delta ESS = \text{grande}$$

este modelo restringido es peor, hay muchas cosas q' no se están explicando
Es mejor el mod ob.

$$Y = \beta_0 + \beta_1 Z_1 + \beta_2 Z_2 + \beta_3 Z_3 + \beta_4 Z_4 + \beta_r Z_r + \varepsilon$$

$$Y(1) = \beta_0 + \beta_1 Z_1 + \beta_2 Z_2 + \beta_r Z_r + \varepsilon$$

$$\hookrightarrow SS_{res}(Z(1))$$

$$ESS = SS_{res}(Z(1)) - SS_{res}(Z)$$

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

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

Sea Z con rango completo, $\varepsilon \sim N(0, \sigma^2 I)$

La prueba $H_0: \beta(2) = 0$ equivale a una prueba q' usa ESS

$$S^2 = \frac{(Y - Z\hat{\beta})'(Y - Z\hat{\beta})}{n-r-1} \rightarrow \text{grande cuando } n \text{ es pequeño}$$

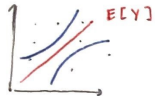
Concretam., se rechaza H_0 si

$$\frac{ESS}{\frac{r-q}{S^2}} > F_{r-q, n-r-1}(\alpha)$$

Para realizar la P.H. se construye el mod. con o sin estas vars (o terminos)

Las sumas de cuad. residuales se comparan en ambos casos.

Inferencia de las predicciones



El $E[Y]$ lo queremos predecir y este también tiene su I.C.

Una vez establecido el mod. se pueden hacer predicciones.

$$\text{Sea } z_0 = \begin{bmatrix} 1 \\ z_{01} \\ \vdots \\ z_{0r} \end{bmatrix} \text{ vals de vars. predictoras.}$$

Pocos esos vals z_0 y $\hat{\beta}$ se puede estimar el val de $E[Y_0]$ o de Y_0

Para $E[Y]$: Sea Y_0 la v.a. respuesta correspondiente a z_0

$$E[Y_0 | z_0] = \beta_0 + \beta_1 z_{01} + \dots + \beta_r z_{0r} = z_0' \beta$$

* El MLE (max likely. estimator) de $z_0' \beta$ es $z_0' \hat{\beta}$

* $z_0' \hat{\beta}$ es el MVUE (min varianza unbiased estimator) de $E[Y_0 | z_0]$

* Su varianza $\text{var}(z_0' \hat{\beta}) = z_0' (Z'Z)^{-1} z_0 \sigma^2$

Si ε normales, el IC $(1-\alpha)100\%$ para $E[Y_0 | z_0] = z_0' \beta$ es:

$$z_0' \hat{\beta} \pm t_{n-r-1}(\alpha/2) \sqrt{(z_0' (Z'Z)^{-1} z_0) S^2} \rightarrow \text{plano/recta de reg. entre q' y y}$$

Para Y_0 : Tenemos $Y_0 = z_0' \hat{\beta} + \varepsilon_0$ $\varepsilon_0 \sim N(0, \sigma^2)$

$$\text{var}(\varepsilon_0) = \sigma^2 (1 + z_0' (Z'Z)^{-1} z_0)$$

El IC con conf. $(1-\alpha)$

$$z_0' \hat{\beta} \pm t_{n-r-1}(\alpha/2) \sqrt{S^2 (1 + z_0' (Z'Z)^{-1} z_0)} \rightarrow \text{siempre es mayor q' el de } E[Y]$$

Donde está estival, varía más q' donde está más recta. \hookrightarrow Dado un val, entre q' y y' es pero q' está

- I Predic° para Y
- I.C para $E[Y]$



Verifica° de los mod.

Podemos verificar si el modelo es adecuado antes de usarlo para realizar predicciones, inferencias, etc.

Si el modelo es válido, $\hat{\varepsilon}$ es un estimador de

$$\varepsilon \sim N(0, \sigma^2 I)$$

en este caso,

$$\hat{\varepsilon} \sim N(0, \sigma^2 (I - Z(Z'Z)^{-1}Z'))$$

Pueden variar las varianzas de $\hat{\varepsilon}$ si los elems de la diag. de H , hii son muy difs.

Si... en algunos los E.

si al graficar los errores veo un patron, no son independientes.
 si hay zonas con mayores errores \rightarrow II.

• se puede usar:

1. $\widehat{\text{var}}(\hat{\epsilon}_i) = s^2(1 - h_{ii})$

2. $\hat{\epsilon}_i^* = \frac{\hat{\epsilon}_i}{\sqrt{s^2(1 - h_{ii})}}$

Para verificar la validez de la regresión.

- Métodos de verificación.

1. Graficar residuales $\hat{\epsilon}_i$ vs \hat{y}_i
(no debe haber dependencia)
2. Graficar $\hat{\epsilon}_i$ vs z_i
(si hay dep. funcional, indica q' hacen falta términos).
3. QQ plots e histogramas residuales de $\hat{\epsilon}_i$ (evaluar normalidad)