

Lezione 15 (a cura di Giovanni Mariani)

Regressione Multivariata

Consideriamo y_1, \dots, y_r , con r = numero variabili risposta
 x_1, \dots, x_k , con k = numero variabili esplicative

Nel modello di regressione Multivariata abbiamo più variabili risposta (tipicamente poche), in particolare avremo:

y_{ij} dove $i = 1, 2, \dots, n$ indica l' i -esimo soggetto
 $j = 1, 2, \dots, r$ indica la j -esima variabile risposta

Avremo quindi un'equazione lineare del tipo:

$$y_{ij} = \beta_{0j} + x_{i1}\beta_{1j} + x_{i2}\beta_{2j} + \dots + x_{ik}\beta_{kj} + \epsilon_{ij}$$

con: x_{ik} dove i indica l' i -esimo soggetto
e k indica la k -esima covariata

E' come se replicassi un modello di regressione multipla r volte, con la differenza, come vedremo, che si fanno assunzioni diverse in merito ai termini di errore.

ESEMPIO

y_1 = spesa alimentare

y_2 = spesa per vacanze

in questo caso allora $r=2$, voglio spiegare le y rispetto a:

x_1 = reddito famiglia

x_2 = numero di figli

x_3 = capofamiglia laureato

in questo caso $k=3$, avrò quindi $2 \cdot 4 = 8$ parametri, ovvero in generale si ha:

$$\# \text{ parametri (coefficienti)} = r \cdot (k+1)$$

nel nostro esempio specifico avrò: $\beta_{01}, \beta_{11}, \beta_{21}, \beta_{31}$ coefficienti relativi a y_{11}
e $\beta_{02}, \beta_{12}, \beta_{22}, \beta_{32}$ coefficienti relativi a y_{12} .

I coefficienti si interpretano nel modo usuale, tenendo a mente che il 2° indice ci dice a quale variabile risposta si riferiscono. Quindi β_{32} sarà l'effetto del capofamiglia laureato sulla spesa in vacanze della famiglia e β_{02} l'intercetta per y_2 , ovvero quando tutte le covariate sono uguali a 0.

Ora, l'**assunzione di base del modello** è che per ogni i , $E[\boldsymbol{\varepsilon}_i]=\mathbf{0}$.

dove $\boldsymbol{\varepsilon}_i$ è un vettore colonna di dimensioni $r \times 1$ e contiene tutti i termini di errore per l' i -esimo soggetto, ovvero:

$$\boldsymbol{\varepsilon}_i = \begin{pmatrix} \varepsilon_{i1} \\ \varepsilon_{i2} \\ \cdot \\ \cdot \\ \cdot \\ \varepsilon_{ir} \end{pmatrix}$$

con $i=1,2,\dots,n$

e la $\mathbf{Var}[\boldsymbol{\varepsilon}_i]=\boldsymbol{\Sigma}$ matrice varianza covarianza, ammetto così che ci può essere una correlazione tra gli errori.

Ritornando al nostro ESEMPIO con $r=2$ avremo che:

$$\mathbf{Var}[\boldsymbol{\varepsilon}_i]=\boldsymbol{\Sigma} = \begin{pmatrix} \mathbf{Var}[\varepsilon_{i1}] & \mathbf{Cov}[\varepsilon_{i1}, \varepsilon_{i2}] \\ \mathbf{Cov}[\varepsilon_{i1}, \varepsilon_{i2}] & \mathbf{Var}[\varepsilon_{i2}] \end{pmatrix}$$

In $\boldsymbol{\Sigma}$ le covarianze non devono necessariamente essere uguali a 0, ammettiamo infatti che ci può essere correlazione, e nella pratica infatti spesso si trova una correlazione maggiore di 0. Questa è una importante differenza rispetto a analizzare separatamente le variabili risposta.

Nel nostro ESEMPIO trovare una covarianza positiva significa che se la famiglia i -esima spende di più per alimenti allora ci aspetteremo che spenda di più anche per le vacanze.

Possiamo esprimere la nostra equazione lineare in una forma più sintetica, in **notazione matriciale**, ovvero nella forma:

$$Y=X*B + E$$

con Y,X,B,E matrici

in particolare avremo che:

Y è una matrice di dimensioni $n \times r$

$$Y = \begin{pmatrix} y_{11} & y_{12} & \dots & \dots & \dots & y_{1r} \\ y_{21} & y_{22} & \dots & \dots & \dots & y_{2r} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ y_{n1} & y_{n2} & \dots & \dots & \dots & y_{nr} \end{pmatrix}$$

contenente tutti i valori delle variabili risposta.

X è una matrice di dimensioni $n \times (k+1)$, è la matrice del disegno:

$$X = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & \dots & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & \dots & \dots & x_{2k} \\ 1 & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & x_{n1} & x_{n2} & \dots & \dots & \dots & x_{nk} \end{pmatrix}$$

B è una matrice di dimensione $(k+1) \times r$, in cui ad ogni colonna corrisponde una variabile risposta, matrice dei coefficienti:

$$B = \begin{pmatrix} \beta_{01} & \beta_{02} & \dots & \dots & \dots & \beta_{0r} \\ \beta_{11} & \beta_{12} & \dots & \dots & \dots & \beta_{1r} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \beta_{k1} & \beta_{k2} & \dots & \dots & \dots & \beta_{kr} \end{pmatrix}$$

E è una matrice di dimensione $n \times r$, per ogni riga a tutti gli errori dell i-esimo soggetto (deve essere coerente con la Y)

$$E = \begin{pmatrix} \epsilon_{11} & \epsilon_{12} & \dots & \dots & \dots & \epsilon_{1r} \\ \epsilon_{21} & \epsilon_{22} & \dots & \dots & \dots & \epsilon_{2r} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \epsilon_{n1} & \epsilon_{n2} & \dots & \dots & \dots & \epsilon_{nr} \end{pmatrix}$$

si può scrivere anche come vettore colonna dei vettori $\underline{\epsilon}_i$ trasposti.

Ora se vogliamo ricavarci l'equazione per il singolo soggetto i avremo:

$$y_{ij} = \underline{x}_i' \underline{\beta}_j + \epsilon_{ij}$$

con \underline{x}_i vettore riga trasposto, $\underline{\beta}_j$ vettore colonna e ϵ_{ij} vettore riga.

Per **stimare il modello** utilizziamo il metodo OLS (minimi quadrati) avremo:

$$Y_{\text{hat}} = X * B$$

dove B sono coefficienti che io ipotizzo e Y_{hat} matrice di dimensione $n \times r$.

Per ESEMPIO se prendiamo: $y_{\text{hat}52} = 3000$, indica la previsione per la quinta famiglia per le spese per le vacanze.

Andando a confrontare la mia previsione con il valore osservato posso vedere l'errore di previsione:

$$Y - Y_{\text{hat}} = Y - X * B$$

con $Y - X * B$ è la matrice degli errori di previsione, questa matrice nel nostro ESEMPIO mi dice per ogni famiglia e per ogni tipo di spesa l'errore di previsione.

Per avere una **misura complessiva degli errori** definisco S , come:

$$S = ||Y - X * B||^2$$

utilizzo l'operatore *norma* (somma degli elementi al quadrato sotto radice, se considero la matrice generica A allora $||A|| = (a_{11}^2 + a_{12}^2 + \dots)^{1/2}$).

Ora facendo la norma al quadrato elimino la radice per cui in pratica è come se facessi una sommatoria:

$$S = \sum_i \sum_j (y_{ij} - y_{\text{hat } ij})^2$$

avrò quindi la sommatoria rispetto alle famiglie (\sum_i) e rispetto alle variabili risposta (\sum_j).

Il passo successivo è minimizzare rispetto a B la quantità S che è funzione di B , quindi:

$$\min_B S(B) \text{ la cui soluzione è } B_{\text{hat}} = (X'X)^{-1}X'Y$$

con X matrice di dimensione $n * (k+1)$

X' matrice di dimensione $(k+1) * n$

Y matrice di dimensione $n * r$

e quindi B_{hat} matrice di dimensione $(k+1) * r$.

Ora per capire se c'è **correlazione tra gli errori** bisogna stimare la matrice di varianza e covarianza degli errori Σ . Andremo quindi a calcolare la matrice dei residui:

$$E_{\text{hat}} = Y - Y_{\text{hat}}$$

NB: ora la matrice Y_{hat} non è più quella con i coefficienti ipotizzati bensì quella con i coefficienti B_{hat} ottimali, ricavati col metodo dei minimi quadrati ($Y_{\text{hat}} = X * B_{\text{hat}}$).

Quindi avremo la Σ stimata di dimensione $r * r$:

$$\Sigma_{\text{hat}} = (E_{\text{hat}}' * E_{\text{hat}}) / (n - (k+1)) = \begin{pmatrix} \text{Var}_{\text{hat}}[\epsilon_{i1}] & \text{Cov}_{\text{hat}}[\epsilon_{i1}, \epsilon_{i2}] \\ \text{Cov}_{\text{hat}}[\epsilon_{i1}, \epsilon_{i2}] & \text{Var}_{\text{hat}}[\epsilon_{i2}] \end{pmatrix}$$

Per fare **inferenza** (verifica delle ipotesi e intervalli di confidenza) devo assumere una distribuzione per gli errori, tale distribuzione sarà un normale, in particolare avremo che:

il vettore $\underline{\epsilon}_i \sim N_r(0, \Sigma)$.